

Durham E-Theses

Continuous-Time Macro-Finance

MATTHEW SEKERKE

How to cite:

SEKERKE, MATTHEW (2021) Continuous-Time Macro-Finance. Doctoral thesis, Durham University.

Use policy

The full-text may be used and/or reproduced, and given to third parties in any format or medium, without prior permission or charge, for personal research or study, educational, or not-for-profit purposes provided that:

- a full bibliographic reference is made to the original source
- a <https://etheses.durham.ac.uk/id/eprint/14105/> is made to the metadata record in Durham E-Theses
- the full-text is not changed in any way

The full-text must not be sold in any format or medium without the formal permission of the copyright holders.

Please consult the [full Durham E-Theses policy](#) for further details.

Continuous-Time Macro-Finance



A thesis submitted in partial fulfilment
of the requirements for the degree of
DOCTOR OF PHILOSOPHY IN ECONOMICS

Matthew Sekerke

Supervisors:

Dr. Kevin Dowd

Dr. Abderrahim Taamouti

Department of Economics and Finance

Durham University

July 2021

Abstract

I characterize asset prices in general equilibrium with risky production. I develop a macroeconomic model that generalizes the framework of Cox, Ingersoll and Ross (1985a) to include elastic labor supply and production with multiple constant elasticity of substitution technologies. I solve for the time- and state-dependent equilibrium of the model using a novel deep neural network technique to approximate optimal policy rules. The calibrated model produces empirically-plausible risk-free rates, equity risk premia, and volatility surfaces from deep microfoundations and reasonable parameter values for uncertainty in production. The numerical solution procedure is far more flexible than standard methods and reveals previously-unknown features of equilibrium behavior and asset prices.

My model leads to a novel understanding of aggregate fluctuations. I find that technological shocks are not sufficient to generate aggregate fluctuations. Instead, I locate their origin in shocks to the *yield* of productive factors within *individual* production processes. Recovery from an adverse shock involves complex adjustments to the quantity and *allocation* of productive factors which bring the economy to a *different* equilibrium state than the one prevailing prior to the shock. We explain why such behaviors cannot be observed in standard macroeconomic models, and discuss their consequences for economic policy.

Declaration

The content of this doctoral dissertation is based on research work completed by the author at Durham University Business School, UK. No material contained in the thesis has previously been submitted for a degree in this or any other university.

Copyright ©2021 by Matthew Sekerke. All rights reserved.

The copyright of this dissertation rests with the author. No quotation from it should be published in any format without the author's prior written consent and information from it should be acknowledged appropriately.

Contents

Abstract	i
Declaration	ii
1 Introduction: The State of Macro-Finance	1
1.1 Consumption Risk and Asset Prices	2
1.2 Asset Prices and Production	8
1.2.1 Production-based asset pricing models	8
1.2.2 Asset pricing in general equilibrium	12
1.3 A Preliminary Appraisal	14
1.4 An Alternative Paradigm	19
1.4.1 General equilibrium with risky production	21
1.4.2 Numerical solutions by deep learning	23
1.4.3 Benefits	25
1.5 Outline of the Thesis	26
2 Production and Asset Prices in Continuous Time: A Portfolio-Theoretic Foundation for Macro-Finance	29
2.1 Intertemporal Equilibrium	31
2.1.1 Risky production from capital	33
2.1.2 Evolution of wealth and the control problem	38
2.1.3 First-order conditions and equilibrium	42
2.2 The Equilibrium Risk-Free Rate	45
2.2.1 Interpreting the risk-free rate	46
2.2.2 The market rate of return	49
2.3 Contingent Claims Valuation	51
2.3.1 Contingent claims risk premia	51
2.3.2 Master valuation equation	54
2.3.3 Specific contingent claims	55
2.4 Conclusions	56

3	A Deep-Learning Solution Procedure for Continuous-Time Macro-Finance Models	59
3.1	Computing Equilibrium in CIR Models	64
3.1.1	Overview of the computational strategy	66
3.1.2	Discretizing the state evolution	70
3.1.3	Learning optimal controls	71
3.1.4	Loss function and optimization	75
3.2	Specification and Calibration	77
3.2.1	Parameterization of the state process	79
3.2.2	Parameterization of the production processes	81
3.2.3	Initialization	83
3.3	Solution of the Model	84
3.3.1	Equilibrium and control profiles	86
3.3.2	Impulse responses	94
3.3.3	Comparison to DSGE methods	96
3.4	Conclusions	101
4	A Continuous-Time Macro-Finance Model with Labor and Capital	105
4.1	Introduction	105
4.1.1	Multifactor production with risk	106
4.1.2	Elastic labor supply	109
4.1.3	Challenges to the asset-pricing literature	111
4.1.4	Plan of the chapter	114
4.2	Production with Capital and Labor	115
4.2.1	Labor supply and factor input processes	116
4.2.2	Output processes	118
4.2.3	Technological change	120
4.2.4	The production function	122
4.2.5	Incomes, expected returns and risk premia	124
4.3	Parameterization and Calibration	127
4.3.1	Parameters determined by empirical studies	127
4.3.2	Calibrated parameters and targets	132
4.4	Solution of the Model	134
4.4.1	Dynamic equilibrium	136
4.4.2	Responses to factor yield shocks	145
4.4.3	Response to technology shocks	151
4.5	Conclusions	154

5	General Conclusions: A Paradigm for Macro-Finance	159
5.1	Contributions to the Literature	160
5.1.1	Dynamic numerical equilibrium	161
5.1.2	Equilibrium with risky production	163
5.1.3	Aggregate fluctuations	165
5.2	Roads Not Taken: Potential Extensions	166
5.2.1	Production risk	167
5.2.2	Heterogeneity	168
5.2.3	Structural estimation	170
5.2.4	Extensions to include money and banking	171
5.2.5	Jumps	173
5.3	Limitations and Open Questions	175
5.3.1	Methodological limitations	175
5.3.2	Theoretical limitations	177
5.3.3	Empirical limitations	178
5.4	Implications for Economic Policy	179
	References	183

Chapter 1

Introduction: The State of Macro-Finance

Few economic phenomena command as much practical and professional attention as the movements of financial asset prices. As an increasing share of the world's growing wealth comes to be held in the form of financial assets, a sound understanding of financial asset price movements becomes ever more important.¹

The capitalist order rests on a belief that the outcomes of competitive markets produce the greatest possible welfare for society, and most reward those individuals contributing the most to that welfare. With financial asset prices seemingly decoupling from everyday economic experience in these first decades of the twenty-first century, the time is ripe for a critical re-examination of how asset prices are determined in general equilibrium, and to ask whether large-scale developments in asset prices over the past decades are consistent with an Arrow-Debreu conception of the economy as an essentially competitive, stable, and allocatively-efficient process.

What counts as an explanation of asset price movements has undergone considerable refinement within the discipline of financial economics. The Copernican revolution ushered in by Markowitz's (1952) portfolio theory and the Capital Asset Pricing Model (CAPM) of Sharpe (1964), Lintner (1965), and Mossin (1966) framed financial market equilibrium as a set of re-

¹Piketty (2014) charts the 'metamorphoses' of capital from agricultural land and public debt to shares and real estate.

relationships that must prevail among asset prices to rule out the possibility of riskless gains, or arbitrage. To define equilibrium asset prices in this way was no small advance, and indeed, the majority of financial professionals can do their work competently using nothing more than the principle of no-arbitrage pricing.

But for modern economists, no-arbitrage conditions are not a complete explanation. Economists seek a causal explanation that connects movements in economic quantities to financial asset price movements. Though a long tradition of research traces movements in asset prices back to developments in the aggregate ‘real’ economy, such as the interwar studies of Keynes and Schumpeter and the seminal general equilibrium models of Arrow and Debreu, more recently such research has developed into a distinct sub-field of macroeconomics known as *macro-finance*, which “studies the relationship between asset prices and economic fluctuations.”² (Cochrane 2017: 945)

The purpose of this thesis is to advance and critically examine the discipline of macro-finance by (1) developing a paradigm for the joint equilibrium of the economy and financial markets that emphasizes risks intrinsic to production processes, and (2) introducing an innovative numerical solution method that permits such models to be analyzed in the absence of analytical solutions. The basis of the paradigm is provided by the continuous-time general equilibrium model of Cox, Ingersoll and Ross (1985a). Solutions are found by adapting a deep learning-based method introduced by Han and E (2016) and further developed by Han, Jentzen and E (2018) to find the optimal dynamic policy responses under uncertainty that characterize equilibrium in the model.

1.1 Consumption Risk and Asset Prices

Because the opportunity cost of any investment is measured in terms of foregone consumption, fluctuations in asset prices must be traceable in large

²I am indebted to Cochrane’s (2008, 2017) surveys of the literature in charting a narrative for the development of macro-finance. I use his work as a foil throughout this thesis owing to his leadership in the field.

part to variations in the utility of consumption. Cochrane (2005) shows how much of asset pricing theory can be unified by the concept of the stochastic discount factor, a ratio of marginal utilities of consumption that falls out of a consumer's intertemporal optimization problem in a natural way. A typical specification of the stochastic discount factor in consumption-based models is

$$m_{t+k} = \beta^k \frac{U'(C_{t+k})}{U'(C_t)} \quad (1.1)$$

where U is a utility function, β is a subjective discount factor, t is the time period at which valuations are obtained, and k indicates how far in the future a payout on investment is to be received. The existence of a risk-free rate, risk premia, the efficient frontier for portfolios, and many other central concepts in finance may all be deduced from the variability of consumption and the covariance of consumption with asset payoffs using the stochastic discount factor. Indeed, the consumption-based approach to asset pricing theory maintains that marginal utilities of consumption are the sole variables of importance in valuing different asset payoffs.

From the perspective of the stochastic discount factor, the development of macro-finance may be seen as a process of grounding fluctuations in consumption and their valuation by a representative agent ever more rigorously in economic theory. Lucas (1978) was the first to focus attention on intertemporal tradeoffs in marginal utility via the consumption Euler equation, from which the stochastic discount factor may be derived. However it was Breeden (1979) who first modeled the intertemporal consumption decision in an explicit equilibrium model he called the consumption-based capital asset pricing model (CCAPM).

Breeden's CCAPM provides a powerful explanation of how risk premia are determined in equilibrium. Risk premia account for returns on investment in excess of the risk-free rate, and all risk premia may be decomposed into the product of a quantity of risk and a price for holding risk. In Breeden's model, the relevant quantities of risk are the variance of consumption and covariances of payoffs with respect to consumption. The price of risk, or the amount an investor is compensated to hold a unit of consumption risk in

equilibrium, is determined by the investor's degree of risk aversion, which is a feature of the utility function.

Breeden (1979) makes two assumptions that were decisive for subsequent work in consumption-based asset pricing. First, Breeden assumes that consumption is the only argument in agents' utility functions, while allowing that the representative agent in the CCAPM may aggregate heterogeneous preferences for consumption risk. By stipulating the sufficiency of consumption risk in determining asset prices, Breeden focused research on the specification of preferences over different consumption streams, while flagging aggregation and heterogeneity as important issues for theoretical and econometric analysis of asset prices. Second, Breeden's model leaves asset markets incomplete. Because the set of available securities does not provide exposure to all dimensions of consumption risk, not all risks to consumption are tradeable in asset markets.³ Thus a second stream of research motivated by the consumption-based theory focuses on expansions of risk premia prompted by agents' 'hedging demands' for those assets that let them offset non-tradeable risks to consumption, however incompletely.

Empirical testing of the CCAPM got off to a rocky start. Hansen and Singleton (1982, 1983, 1984) estimated and tested the consumption Euler equation, in some of the first studies applying Hansen's generalized method of moments (GMM) techniques. Their tests using aggregate consumption time series and constant relative risk aversion (CRRA) preferences rejected the consumption-based model.

Mehra and Prescott (1985) pioneered an alternative approach in which consumption time series are generated by simulating a calibrated equilibrium model. Unlike Hansen and Singleton, who estimate preference parameters from data, Mehra and Prescott set the parameters of their model with reference to a benchmark real business cycle (RBC) model and solve for the equity risk premium implied by the model for a range of values for the coefficient of relative risk aversion. Their efforts framed the famous 'equity premium puzzle,' based on the observation that the historical level of the equity risk

³Economists frequently employ the jargon of 'spanning' basis sets from linear algebra to say the above using the shorthand 'consumption risk is not spanned by asset markets.'

premium could not be reconciled with reasonable values for the degree of relative risk aversion in a model which, in their view, otherwise does a good job of capturing the dynamic behavior of the aggregate economy. We will have more to say about Mehra and Prescott's approach below.

Later developments in a 'first generation' of studies reviewed by Breen, Litzenberg and Jia (2015a) refined the measurement of consumption risk while working within the class of time-additive utility functions with CRRA preferences. Improvements in the measurement of consumption risk have been achieved by Hansen and Singleton (1983), Cox, Roll and Ross (1986), Litterman and Ludvigson (2001a, b), Parker and Julliard (2005), and Jagannathan and Wang (2007). At the same time a more disaggregated understanding of consumption risk led to refinements in the measurement of risk aversion. The studies of Mankiw and Zeldes (1991), Heaton and Lucas (1992), and Brav, Constantinides and Geczy (2002) show how disaggregating the consumption of subpopulations results in more volatile consumption streams for those with more exposure to equities, and therefore more reasonable measurements of risk aversion for individuals who actually hold equities. Consumption-based asset pricing research emphasizing market incompleteness is surveyed by Kocherlakota (1996) and Campbell (2003).

However focus would quickly shift away from market incompleteness to the specification of preferences over consumption streams. The 'second generation' of empirical work on the CCAPM surveyed by Breen, Litzenberger and Jia (2015b) takes a critical look at the first generation's use of time-additive, CRRA utility functions. By reformulating the utility function to incorporate habituated levels of consumption or recursive valuations of consumption researchers have posited preferences over the time path of consumption.

The models of Constantinides (1990) and Campbell and Cochrane (1999) posit that consumers value consumption streams relative to internal and external consumption habits, respectively.⁴ A typical specification of prefer-

⁴Habits are internal if they are personal to the agent, and external if they are based on aggregate consumption. The latter is more tractable mathematically and empirically.

ences with a habituated level of consumption is

$$U(C_t, C_{t-1}, \dots, C_{t-k}) = \frac{(C_t - C_H)^{1-\gamma}}{1-\gamma} \quad (1.2)$$

where C_H is the habituated level and γ is the local coefficient of relative risk aversion. The habituated level C_H is often defined by distributed lags of past consumption, i.e., $C_H = f(C_{t-1}, \dots, C_{t-k})$, eliminating time separability in the utility function. Its presence in the utility function makes the representative agent value consumption in the neighborhood of C_H like they would value near-zero consumption in a utility function without habits, allowing lower values of γ to be compatible with extreme levels of local risk aversion and the equity risk premium.

Following early work by Kreps and Porteus (1978), the recursive utility formulations of Epstein and Zin (1989) and Weil (1989) decouple risk aversion from the elasticity of intertemporal substitution. A typical specification of recursive preferences defines a constant-elasticity-of-substitution aggregate of current consumption and expected future consumption, while confining relative risk aversion to expected future consumption alone. Thus for a coefficient of relative risk aversion γ , a subjective discount factor β , and an elasticity of intertemporal substitution σ , define $\rho = \frac{\sigma-1}{\sigma}$ and write preferences recursively as

$$U_t = \left((1-\beta)C_t^\rho + \beta E [U_{t+1}^\gamma]^\frac{\rho}{\gamma} \right)^{1/\rho} \quad (1.3)$$

Recursive utility allows investors not only to prefer less uncertainty about consumption to more via the parameter γ , but also to prefer earlier resolutions of uncertainty to later resolutions via the elasticity of intertemporal substitution σ . Like utility functions with habituated levels of consumption, recursive utility functions also allow lower coefficients of relative risk aversion to be compatible with the equity risk premium for observed levels of consumption risk.

The more recent generation of empirical work has also led to further improvements in measurements of consumption risk. Lettau and Ludvigson

(2001a, b) condition consumption risks on the level of wealth, showing that aversion to consumption risk increases as wealth declines. In addition to applying Epstein-Zin/Weil preferences, Bansal and Yaron (2004) show the importance of persistence in the consumption process and time variation in the volatility of consumption. And Santos and Veronesi (2006) show time variation more generally in the risk characteristics of assets, leading to significant time variation in the equity risk premium.

Some interesting contributions by Eraker and Shaliastovich (2008) and Eraker (2008) connect a consumption-based foundation to the empirically-important class of exponential affine models for contingent claims prices. The class of exponential affine models potentially encompasses pricing implications for any asset with a term structure, which helps to push the boundaries of the consumption-based research program beyond the puzzle of the equity risk premium which had dominated earlier work.

Despite the empirical tenacity of the CCAPM and the intuitive role consumption risk plays in the determination of risk premia and asset prices, the CCAPM falls short as a deep, causal explanation of dynamic risk premia. A deep explanation must trace consumption risk to its origins in the productive capacity of the economy. To follow the consumption-based asset pricing research program in modeling production as an ‘endowment process’ gifting random levels of consumption to the economy is question-begging, and tweaking the preferences of a representative agent to dislike the output of the endowment process just enough to match the data is an unsatisfying strategy for ‘explaining’ asset prices. Inasmuch as the consumption-based theory hangs on the specification of the utility function, its empirical content is slim: a specification for the habituation process C_H or a value for the elasticity of intertemporal substitution σ that is not rejected by the data, given a tolerable constant of relative risk aversion γ . Perhaps more fundamentally, the consumption-based theory with time-dependent preferences gives us a subjective explanation of intertemporal substitution, rather than describing rates of intertemporal substitution as quantities determined by the behavior of optimizing agents in equilibrium.

1.2 Asset Prices and Production

Two alternative streams in the macro-finance literature aim to connect the determination of asset prices to the productive capacity of the economy. The first stream attempts to remove consumption decisions and preferences from the determination of asset prices all together, while the second stream models asset prices in a general equilibrium with endogenous production and consumption.

1.2.1 Production-based asset pricing models

Production-based models of asset pricing focus on the decisions of a representative profit-maximizing firm rather than a representative consumer. For the profit-maximizing firm the fundamental intertemporal tradeoff is that between producing output today and investing in partially-finished goods and productive capacity that will yield output in the future. The fundamental economic risk faced by firms is the risk of shortfalls in current output when resources are invested in expanding future production. Expected excess returns on physical investment compensate producers for risking current output. Returns on physical investment are linked to financial investment returns via Tobin's q theory of capital market equilibrium (Tobin 1969).

A stochastic discount factor may be derived for production-based models in analogy with (1.1). The stochastic discount factor of the consumption-based theory is the product of a subjective discount factor and a ratio of marginal utilities. The expected marginal utility of future consumption divided by the marginal utility of current consumption gives the marginal rate of substitution between current and future consumption. One can view the marginal rate of substitution as a ratio of prices for dated consumption. In the production-based theory the marginal rate of substitution is replaced by the marginal rate of transformation, the ratio of expected future and current marginal values of output to the firm. These marginal values may also be expressed as a ratio of prices. Cochrane (1991) accordingly formulates a

production-based stochastic discount factor of the form

$$m_{t+k} = \rho^k \frac{Q_{t+k}}{Q_t} \quad (1.4)$$

where $Q_{t+k} = P(s^{t+k})/\rho^{t+k}\pi(s^{t+k}|s^t)$ is the price of a unit of output delivered in a state s^{t+k} that follows s^t . The firm adjusts its production to equate its marginal cost of production to these contingent claims prices, which give the marginal benefit of production.⁵

The representative firm's first-order conditions imply an equilibrium return on physical investment which might be expected to explain returns on investments in equities. Cochrane (1991) develops a series of expected investment returns from aggregate data on gross investment and finds that ex post investment returns and stock returns are highly correlated, among other encouraging correspondences. If the investment decisions of a representative firm with a single production technology can reproduce aggregate stock returns, then we might surmise that multiple representative firms with multiple production technologies could be the basis of a multi-factor explanation for the cross-section of equity returns.⁶

Cochrane (1996) pursues a production-based explanation of the cross-section of equity returns using time series for gross residential and non-residential investment. He derives an expression for the expected return on physical investment with unknown sector-specific parameters for the marginal product of capital, the depreciation rate, and the cost of adjustment to new investment. The stochastic discount factor that prices the cross-section of equity returns is no longer expressed as a ratio of contingent output prices, but as a linear combination of the expected returns on physical investment. Using four selected decile portfolios, a proxy for the risk-free rate, and two instruments, Cochrane (1996) derives moment restrictions and jointly estimates the unknown parameters for the sector-specific production technologies and the

⁵Belo (2010), following Cochrane (1993), provides an alternative formulation in terms of productivity levels.

⁶Studies of 'the cross-section of equity returns' put asset pricing theories to the test by requiring them to explain returns on multiple portfolios simultaneously, as opposed to studies that seek to explain the movements of an aggregate index of equity prices.

stochastic discount factor by GMM.⁷ Cochrane's strategy thus wagers that time series on physical investment and well-chosen moment conditions can use stock returns to provide insight into the structure of production, while simultaneously leveraging producers' optimality conditions to parameterize a stochastic discount factor for the cross-section of equity returns.

Regrettably, Cochrane's empirical strategy did not live up to the ambitions that motivated it. The selected equity return series and instruments were not sufficient to identify even three unknown parameters in the specification of production technology, so Cochrane fixes depreciation rates and adjustment costs, leaving the marginal products of capital free. The latter are well-estimated, but a long-run average marginal product of capital for a sector does not tell us much about production. Indeed, Cochrane (1996) tests an alternative specification in which production technology is completely suppressed. Expected returns on physical investment are replaced with physical investment growth rates computed directly from the data, i.e., without reference to optimality conditions or a functional form for production. When production is suppressed, growth rates in investment actually do a better job of explaining the cross-section of equity returns. These results suggest that either (a) the linear form chosen for the stochastic discount factor makes poor use of the information contained in the physical investment return series, or (b) the sectoral investment return series absorb no information from the partial equilibrium theory of production beyond the information already present in the volume of investment.

The results of Cochrane (1996) boosted appreciation for GMM as an econometric methodology for testing asset pricing models, but dealt the production-based asset pricing theory a blow from which it hasn't really recovered. Li, Vassalou and Xing (2006) improve a bit on Cochrane's stochastic discount factor specification using a different sectoral disaggregation of investment growth rates, but further work in this direction isn't particularly informative; after all, why should residential investment by households (Cochrane) or investment by non-corporate business (Li, Vassalou and Xing) have anything to do with the returns on investment expected by the *corporate*

⁷The term premium and the dividend-price ratio are used as instruments.

business sector that actually issues publicly-traded equity? The careful study of Liu, Whited and Zhang (2009) shows that the cross-section of stock returns is consistent with the ‘investment Euler equation’ of a profit-maximizing firm and the q theory of investment. However they find, like Cochrane (1996), that equity returns provide little insight into the structure of production. Liu, Whited and Zhang use accounting data to determine the stock of capital, investment, output, depreciation, capital structure, effective tax rates and interest costs at the firm level, a strategy that is difficult to make operational at the macroeconomic level. Belo (2010) formulates a model with state-contingent production plans and estimates two parameters in an aggregate production function: the elasticity of substitution across states, and the sensitivity of production to a common productivity factor. His estimates of both parameters are often statistically indistinguishable from the neutral value of unity, and when differences from unity are statistically significant, they are not economically significant.

Under these circumstances it is not clear what equity returns can tell us about the structure of production, or what a better understanding of production would tell us about the structure of asset price risks. Nevertheless, dismissing production as an important source of risk would be unwise. I agree with Cochrane (2008: 290) that “If we want to link asset prices to macroeconomics, consumption seems like a weak link,” particularly when cyclical movements in output, investment, and employment dwarf the placid wanderings of aggregate consumption. Instead of the ‘either-or’ choice between consumption- and production-based theories, a ‘both-and’ approach employing consumer choice in a dynamic general equilibrium with explicit production might capture the best of both approaches.

The production-based theory also shows the limits of what econometric studies of asset prices can tell us about the structure of production and its role in determining risk premia. If production is to be specified in a model of asset prices, unknown parameters will have to be calibrated or fixed on the basis of prior empirical investigations. With the turn toward general equilibrium theories of asset pricing our primary tools of inquiry are calibrated macroeconomic models, solved numerically and simulated to

reveal the model's implications for asset pricing.

1.2.2 Asset pricing in general equilibrium

A paradigmatic example of a general equilibrium asset pricing model is Jermann (1998). In general equilibrium the first-order conditions of the producer and the consumer interact to determine equilibrium. Like Mehra and Prescott (1985), Jermann takes an RBC model as his starting point, based on the belief that the RBC model does reasonably well in accounting for the essential features of business cycles. In place of the Lucas (1978) model used by Mehra and Prescott, Jermann starts from the model of King, Plosser, and Rebelo (1988a, 1988b, 2002), which includes an explicit model of production.

Jermann was not the first to investigate whether the equity risk premium puzzle might arise from the absence of production in the Lucas (1978) and Mehra and Prescott (1985) models. Earlier efforts by Danthine, Donaldson and Mehra (1992) and Rouwenhorst (1995) found that introducing production only aggravated the equity premium puzzle, because the ability to control production through investment decisions gives risk-averse households another lever with which they might smooth their consumption. Accordingly, Jermann (1998) employs a concave adjustment cost for investment in a nod to the Q theory literature. The adjustment cost penalizes large changes in investment with output losses.

Jermann's modifications to Mehra and Prescott (1985) are not confined to the production side of the economy. In place of CRRA preferences he uses a simplified version of the habit formation preference specification of Campbell and Cochrane (1999).⁸ Jermann thus allows time-inseparable preferences and investment decisions to interact in reconciling the equity risk premium with low levels of aggregate consumption volatility. And consistent with the consumption-based asset pricing literature – but unlike King, Plosser and Rebelo (1988a, 1988b) – the representative agent in Jermann's model values only consumption and supplies his labor inelastically.⁹

⁸He cites the 1995 NBER working paper version.

⁹The elastic supply of labor in the King, Plosser and Rebelo model is arguably its most important feature. We will return to this point in Chapter 4.

Jermann (1998) sets the labor share in production, the rate of depreciation, the trend rate of growth, and the parameters of the process generating total factor productivity shocks with reference to the classic RBC literature, while fixing the local coefficient of relative risk aversion at 5. He then calibrates the remaining parameters of the model – the habit formation coefficient, capital adjustment costs, the subjective discount factor and shock persistence – to match the mean risk-free rate, the mean equity risk premium, and the ratios of consumption and investment growth volatilities to the volatility of output growth. Asset prices are calculated using the lognormal approximations of Campbell (1986). The log-linearized model is solved using the steady-state perturbation method of King, Plosser and Rebelo (1988a, 1988b, 2002).

Jermann finds he can match his target variables quite closely with his calibration; thus his model reproduces the mean risk-free rate and the mean equity risk premium almost exactly. Accordingly Jermann's model allows him to investigate the relative contributions made by preferences and investment frictions to determining the equity risk premium. He finds that neither is sufficient on its own. Habit formation alone fails because agents adjust production rapidly, as in Danthine, Donaldson and Mehra (1992) and Rouwenhorst (1995). Adjustment frictions alone produce an equity risk premium an order of magnitude smaller than the data. Jermann also uses impulse responses to gain insight into the causal origins of the equity risk premium. Decomposing the equity risk premium into payout uncertainty and valuation uncertainty, he finds that capital adjustment costs make dividends more procyclical, generating larger premia for payout uncertainty.

All in all, it must be said that Jermann (1998) makes a significant theoretical and methodological advance. He married key insights of the consumption- and production-based approaches in a general equilibrium model, showing that both aspects can interact to explain the equity risk premium and achieve a reasonable value for the risk-free rate. He also showed how asset pricing consequences could be read off of a calibrated model that had been solved with standard numerical techniques.

These achievements launched a new literature on asset pricing in gen-

eral equilibrium.¹⁰ Two especially interesting examples are Kung (2015) and Chen (2016). Kung (2015) studies the equilibrium term structure of interest rates and time-varying risk premia in a stochastic endogenous growth model. Chen (2016) incorporates time-to-build and time-to-produce delays into a general equilibrium model with recursive preferences. His model demonstrates the role of inventories in smoothing consumption risk, and reconciles the model's implications with the tendency for asset price movements to lead quantity movements in the data. Overall, Chen's model reproduces the level and volatility of the equity risk premium, and produces a low risk-free rate.

For all its virtues, however, Chen's paper suffers from the problems that plague much of the equilibrium asset pricing literature: an excessive focus on TFP shocks, a proliferation of 'frictions' to amplify minimal production risk, a dependence on non-standard utility to match asset pricing facts, and a reliance on essentially deterministic numerical solutions centered on a steady state.

1.3 A Preliminary Appraisal

Where do we stand after 50 years of effort to provide causal explanations of asset price movements in terms of macroeconomic aggregates?

There can be no doubt that the consumption-based approach to asset pricing dominates empirical macro-finance. The consumption Euler equation and the stochastic discount factor it implies are the fundamental lenses through which asset pricing theories are understood and tested. They focus attention on preferences and consumption, framing the empirical study of asset pricing in terms of consumption risk and a small set of unknown parameters. This narrow framing gives the consumption-based approach its power and accounts for its empirical successes. At the same time, the parsimony of the consumption-based approach restricts its potential empirical

¹⁰See Boldin, Christiano and Fisher (2001), Jermann (2010, 2013), Gomes, Kogan and Zhang (2003), Zhang (2005), Kaltenbrunner and Lochstoer (2010), and Croce (2014). Extensions of the general equilibrium approach that aim to explain the cross section of equity returns like Berk, Green and Naik (1999), Gomes, Kogan and Zhang (2003), Zhang (2005), and Gala (2006) introduce multiple production processes.

content and weakens it as a source of causal explanations for asset price movements. The behavioral component of such models has far more to do with agents' *valuations* of different outcomes than agents' *actions* to establish equilibrium. To say that holders of financial claims must be compensated more for certain risks to consumption because they particularly dislike those risks is deeply unsatisfying, and it strains credulity to think that all risks of interest in equilibrium are mediated exclusively through consumption.

Production-based approaches are an important corrective, highlighting the risk-taking and profit-maximizing behavior of firms and the uncertainty of the investment process as a source of risk. Yet the production-based approach has few empirical successes to its name. Though changes in gross physical investment have considerable power in explaining the cross-section of equity returns, the production-based approach has not delivered a clear story that establishes a causal link between investment flows and equity returns via the optimizing decisions of producers.

Unifying the two approaches in a general equilibrium theory of asset pricing has an almost-inevitable, Goldilocks-like feeling to it – a path that unites the empirical success and intuitive appeal of the consumption-based literature with the groundedness of the production-based approach. Surely an understanding of asset pricing in some kind of general equilibrium marks the way forward for macro-finance. But the current practice of general equilibrium modeling in macro-finance leaves much to be desired, and contains theoretical and methodological obstacles to further progress. These obstacles stem, in my view, from the strong framing of the problems of macro-finance in terms of the now-classical equity risk premium puzzle of Mehra and Prescott (1985).

The equity risk premium puzzle is classically a matter of reconciling a 'big' equity risk premium with surprisingly 'small' aggregate risk, where aggregate risk is defined as the volatility of consumption. The puzzle is not a recalcitrant feature that appears in a variety of macroeconomic models; rather, it arises in a very particular context: the model of Lucas (1978) and subsequent developments in real business cycle (RBC) theory. According to Mehra (2012: 396), the equity risk premium puzzle

arises because quantitative predictions of the [RBC] theory [for the equity risk premium] are an order of magnitude different from what has been historically documented. The puzzle cannot be dismissed lightly, given that much of our economic intuition is based on the very class of models that fall short so dramatically when confronted with financial data. It underscores the failure of paradigms central to financial and economic modeling to capture the characteristic that appears to make stocks comparatively so risky. Hence the viability of using [RBC] models for any quantitative assessment, say, to gauge the welfare implications of alternative stabilization policies, is thrown open to question.

Though the equity risk premium puzzle threatens a radical destabilization of the RBC theory, responses to the puzzle have been fairly conservative. The class of RBC models has most certainly not been dislodged by the puzzle as a preferred source of economic intuition.¹¹ Nor is economic intuition the only matter at stake. Economists have overwhelmingly built their protocols for solving and evaluating numerical models of the economy on RBC foundations. That the further elaboration of asset pricing theory in general equilibrium has been defined by comparatively slight modifications to a tenacious RBC paradigm should shock no one. A major modification would be a heavy lift.

Upon making the transition from partial equilibrium to general equilibrium modeling, the proving ground for a theory shifts from econometric modeling and specification testing to simulations and/or impulse responses obtained from a calibrated numerical model. The relationship between the model and empirical data becomes more difficult to establish, and leans heavily on the hypothesis of ‘rational expectations’ that dissolves the distinction between historical performance and expected future results. New methodological problems arise: how to obtain a numerical solution, how to parameterize the model, and how to study the properties of the model effectively. The handling of these methodological questions is so closely bound up with the development of RBC theory that we might well speak of an RBC

¹¹We accept Mehra’s assessment for purposes of argument; I do not mean to assert that RBC models command universal assent.

methodology that dominates general equilibrium research in macro-finance. The early work of Kydland and Prescott (1982) and King, Plosser and Rebelo (1988a, 1988b) laid the groundwork for how numerical macroeconomic models are calibrated, solved, validated against empirical data and studied every bit as much as they made signal advances in the development of RBC theory.

Before finding the numerical solution of a macroeconomic model one characterizes the equilibrium analytically through the first-order conditions and constraints of the model, and then expressing the equilibrium as an approximately linear function of log-deviations from the model's steady state (Uhlig 1999).¹² The need for an analytical starting point puts a premium on models for which such solutions may be found. RBC models offer analytical tractability, at the expense of several simplifications. In our view two of those simplifications are crucial. First, the use of a single aggregate production function eliminates questions of allocation from the analysis of equilibrium, preventing switches in processes and technologies from playing any role. Secondly, RBC models concentrate risk into total factor productivity fluctuations, a *deus ex machina* that forestalls further inquiry into phenomena that might actually *explain* fluctuations in output.

The conventional numerical methods employed to solve general equilibrium models further mask the risks that ought to be the primary object of study in macro-finance. Cochrane (2008: 300) writes:

I remain a bit worried about the approximations in general equilibrium model solutions. Most papers solve their models by making a linear-quadratic approximation about a non-stochastic steady state. But the central fact of life that makes financial economics interesting is that risk premia are not at all second order. The equity premium of 8 percent is much larger than the interest rate of 1 percent. Thinking of risk as a 'second-order' effect, expanding around a 1 percent interest rate in a perfect foresight model, seems very dangerous.

¹²I review standard numerical solution methods for general equilibrium models more fully in Chapter 3.

When a model is built to have globally stable dynamics and initialized in a steady state equilibrium, potential economic reactions to a shock are severely constrained. One suspects that the numerical methods used effectively bound the risks that can be generated and analyzed in general equilibrium models, as well as the range of behavioral responses that might be expected.

The rise of the RBC methodology as the central paradigm in macro-finance coincided with the submergence of the intertemporal portfolio theory of Merton (1973) that first got macro-finance going. Formulated in the language of continuous time stochastic processes and dynamic programming, the portfolio theory approach reached a highly developed form as a basis for macro-financial equilibrium in Breeden's (1979) CCAPM. Portfolio theory models make the choice of allocation a defining feature of equilibrium, and analyze aggregate risk in terms of a shifting investment opportunity set defined by any number of stochastic state variables. As we pointed out above, these aspects of the portfolio theory approach were not subsumed by the RBC approach so much as suppressed by it.

However the intertemporal portfolio theory approach to general equilibrium failed to develop robust solution methods. Few analytical solutions followed Merton's initial breakthroughs for a fixed investment opportunity set, and little progress was made with numerical solution methods. As a result, the field has generally left the portfolio approach behind. "Dynamic incomplete-market portfolio theory is hard," and "widely ignored in practice, though it has been around for half a century," writes Cochrane (2014: 4-5). "Institutions, endowments, wealthy individuals, and regulators struggle to use even the discipline of mean-variance analysis in place of name-based buckets, let alone to implement Mertonian state-variable hedging." Intertemporal portfolio theory would appear to have poor prospects as a basis for general equilibrium theory when even those with the greatest pecuniary interest in getting it right have given up trying to implement it.¹³

Yet it seems to me that the equity risk premium persists as a 'puzzle' for macro-finance mainly because the methodological choices that cause the

¹³We undertake a more detailed review of research on the intertemporal portfolio problem in Chapter 3.

puzzle to arise are so rarely called into question. Restricted to a steady state with a lonely, benignly stochastic total factor productivity shock as the sole source of uncertainty, there is really very little opportunity for risk to play a meaningful role in general equilibrium asset pricing models. We have swept all of the ‘state variables’ that determine investment opportunities and generate interesting patterns of risk in Merton (1973) and Breeden (1979) under the rug, and for no good reason.

Progress will not be made in macro-finance by adding new frictions to RBC models until the steady-state distributions of risk-free rates and equity risk premiums match those of the data. Instead, we should be asking whether an alternative to the RBC methodology might provide a more compelling basis for research in macro-finance.

1.4 An Alternative Paradigm

In my view the way forward is for macro-finance to return to its roots in dynamic portfolio theory, highlighting the choice of allocations as a central feature of equilibrium and explicitly modeling changes in the set of production possibilities. At the same time a general equilibrium model of asset prices should go beyond the analysis of investment embodied in dynamic portfolio theory to model the level and allocation of *labor* effort. In a model with labor and capital, changes in the state variables of Merton (1973) and Breeden (1979) affect not just returns on investment, but the set of multifactor production possibilities. Under such a construct, causal explanations of asset price fluctuations proceed from changes in the production opportunity set, to time- and state-dependent decisions made by optimizing agents about the level and allocation of productive resources among multiple production technologies, and ultimately to the configuration of risk-free rates and risk premia entailed by agents’ optimizing behavior.

The uncertain production yields and unnamed state variables of dynamic portfolio theory are ideal vehicles for re-introducing the aggregate economic risk that is so clearly missing in Mehra and Prescott (1985). Whereas the basic RBC framework concentrates aggregate risk into relatively small to-

tal factor productivity shocks, a framework grounded in dynamic portfolio theory models output as a fundamentally uncertain process and allows the distribution of output to depend on the stochastic state of the economy. We need a theory that names those state variables and connects them in intelligible ways to production and consumption. Only then will we have a useful framework in which to study the dynamics of economic risk and the means by which it is diversified, distributed, and held inside and outside of market institutions.¹⁴ We need to discover how risk fluctuates in response to intelligible shocks and what governs the dynamics of risk premia in general equilibrium. And then we may well find ourselves asking how it is possible that market expectations converge around a relatively narrow set of outcomes, rather than puzzling over why a civilization ostensibly gifted with virtually assured consumption demands so much return on equity investments. My goal is not to resolve the equity risk premium puzzle. I expect to turn it on its head.

The purpose and primary contributions of this thesis are to (a) propose, (b) elucidate, and (c) exemplify an alternative paradigm for macro-finance, motivated by the view that current general equilibrium approaches to macro-finance do not – and cannot – deal adequately with production risks or the actions taken by economic agents to manage those risks. The theoretical component of the alternative paradigm replaces the RBC framework with a generalization of the production-based stochastic control model of Cox, Ingersoll and Ross (1985a) (hereafter CIR85a). The methodological component of the alternative paradigm replaces standard DSGE numerical procedures with a novel solution procedure that finds time- and state-dependent optimal controls via deep learning, building on the methods of Han and E (2016) and Han, Jentzen and E (2018). The use of CIR85a as a point of departure serves to return the general equilibrium theory of asset pricing to

¹⁴Browning, Hansen, and Heckman (1999) criticize general equilibrium modeling on the grounds that it is insufficiently sensitive to heterogeneity and risk-sharing on the household or consumption side of the ledger. Campbell's (2018) research program picks up this gauntlet, showing how risk is managed by financial decision making within households. What I am suggesting is, in a sense, a version of this critique that emphasizes heterogeneity and risk management on the production side of the economy.

its roots in portfolio theory, and to focus attention on risk dynamics in a complete markets setting. The solution procedure based on deep learning allows us to obtain genuinely stochastic time- and state-dependent solutions in a high-dimensional setting, in stark contrast to the time-invariant and non-stochastic solution methods typically employed in general equilibrium analysis.

1.4.1 General equilibrium with risky production

We take CIR85a as our point of departure because we view it as the apotheosis of dynamic portfolio theory. Whereas Merton (1973) and Breeden (1979) made significant advances in describing investment allocations and consumption decisions in general equilibrium, they did so in an incomplete markets setting with an exogenous risk-free rate. CIR85a goes beyond Merton and Breeden by introducing a complete set of contingent claims and showing how the risk-free rate, the equity risk premium, contingent claims prices and expected excess returns on contingent claims are all determined in equilibrium. CIR85a unifies asset pricing theory in a general equilibrium setting, albeit at the cost of an explicit solution for optimizing behavior in equilibrium, which is left implicit in their model.

CIR85a models production as a disaggregated set of N stochastic processes whose outputs depend in an unspecified way on the evolution of K exogenous stochastic state variables. An endogenous state variable, wealth, accumulates returns on investment in the N processes. Agents choose how much to consume from each increment of wealth, and how to allocate their accumulated wealth among production processes. CIR85a thus models the allocation decision faced by a representative agent in the presence of fundamental uncertainty.¹⁵ The resulting state-contingent allocation of productive effort may be interpreted as a structural representation of the reduced-form aggregate marginal rate of transformation proposed by Cochrane (1993), Jermann (2010), and Belo (2010) for the production-based pricing kernel.

¹⁵The tension between methodological individualism and the representative agent and social planner interpretations of stochastic control problems is discussed in Chapter 2.

Dynamic portfolio problems are formulated over a finite time horizon to contrast the benefits of time-varying allocations with the static allocations implied by standard portfolio theory. CIR85a adopts this finite-horizon formulation, leading to a dynamic notion of equilibrium that may be contrasted with the familiar steady-state solutions favored by economists. Optimal consumption and allocation behaviors vary over time, as well as over different states of the economy. The state- and time-dependent CIR85a equilibrium entails state- and time-dependent asset prices à la Arrow and Debreu that are expected to prevail over the finite time horizon, based on the information available to agents at the beginning of the time horizon. Accordingly we interpret equilibrium in the CIR85a model as a set of forward-looking, state-dependent dynamic plans.

General equilibrium consequences for asset prices are summarized compactly in CIR85a by the equilibrium risk-free rate r , the equity risk premium ϕ_W/W , and a vector of risk premia ϕ_Y associated with the economic state variables.¹⁶ Genuine financial objects thus emerge from the equilibrium model, and may be combined in a modular way to compute values for contingent claims using formulas familiar to financial engineers and researchers in empirical asset pricing. A ready set of correspondences between model outputs and asset pricing formulas means that a variety of observable consequences can be generated for any candidate specification of the economy. In principle, the asset-pricing consequences of CIR85a have the potential to shift the proving-ground of macroeconomic models from aggregate time series to panels of contingent claim prices traded in financial markets.¹⁷

Risks associated with the economic state variables are assumed to be tradeable in complete markets for contingent claims. Some may see the assumption of complete markets as a step backwards from the ‘more general’

¹⁶The equity risk premium may be further disaggregated into multiple factors characterizing the cross-section of equity returns.

¹⁷Though a full pursuit of this point is beyond the scope of the present study, we speculate that enlarging the set of observable consequences in this way may eventually enable identification of macroeconomic model structures and parameters, a hitherto elusive goal for the discipline (see, e.g., David Romer 2016). We further note that it is exceedingly difficult to sustain the rational expectations hypothesis in a setting where agents’ plans are expected to change at every date.

setting of incomplete markets. We disagree. The *fiction* of complete markets for contingent claims provides us with a set of objects representing the economic risks faced by society. The *fact* of incomplete markets invites us to think hard about where those risks reside when they cannot be traded away.¹⁸

To make the CIR85a framework suitable for use as a paradigm theory for macro-finance, we generalize the model to incorporate labor supply and allocation decisions. The generalized model allows agents to control the supply and allocation of all productive resources to the set of available production technologies. Macro-finance models rarely consider elastic labor supply to be an important feature of equilibrium, though it has important consequences for general equilibrium and risk aversion. Using a stochastic generalization of a constant elasticity of substitution production function, we model output as a function of labor and capital inputs subject to factor-augmenting technical progress. Our generalized model of production thus incorporates uncertainties about input quality and technical progress that are mostly ignored in the RBC paradigm.

1.4.2 Numerical solutions by deep learning

If we plan to set up an analytically-intractable model as our paradigm case, we had better have a good numerical solution method. To this end we draw on deep learning methods.

Solving a dynamic economic model amounts to searching for functions that describe the responses of an optimizing agent to the state of the economy. Within the field of applied mathematics, deep learning is increasingly appreciated as a technology for general-purpose function approximation. The neural networks employed in deep learning have been proven to be universal function approximators (Hornik, et al. 1989, Cybenko 1989). Accordingly deep learning methods are seeing an expanding range of application in scientific computing wherever the solution of a problem is function, especially in the numerical analysis and solution of nonlinear partial differential equations

¹⁸See Staum (2008) for a (somewhat dated) survey of research on incomplete markets.

that define the solution of stochastic control problems.¹⁹

Our deep learning-based solution method permits a direct attack on the equilibrium of a complex economic model specified in terms of stochastic differential equations. It finds dynamic, state-dependent solutions to stochastic control problems over the entire state space. It does not require linearizations, approximations, or knowledge of a steady state. It treats risk by simulating the entire range of outcomes inherent in the specification of the model. It is also sufficiently scalable to solve problems with tens or hundreds of state variables. To my knowledge this thesis is the first application of such methods within the fields of economics and finance.

The application of a novel method to obtain a numerical solution of the CIR 1985a model represents a significant advance for the field of macro-finance. Standard solution methods truly pale in comparison. Instead of a deterministic steady-state solution rendered ‘stochastic’ by a small perturbation or a small-scale discrete approximation, numerical solution by deep learning methods uncovers a set of genuinely state-dependent functions that characterize equilibrium in the presence of actual risk. In the context of CIR85a, these capabilities permit exploration of the deep structure of a general equilibrium that includes financial markets. Whereas CIR85a analyzes the first-order conditions of an implicit solution, and CIR85b obtains an explicit solution under very restrictive specializations, our numerical methods allow explicit solutions to be obtained under quite general conditions. In particular we can see exactly how consumption and investment decisions depend on state variables, given a specification of the latter and of their influence on production.

We can also simulate equilibrium dynamics for the economy and asset prices in response to specific shocks. Our menu of shocks expands from the earlier, monolithic aggregate TFP shock to encompass all production processes and all state variables. Disaggregating risk in this way makes our modeling framework rich enough, in principle, to explain facts about

¹⁹See, for example, Sirignano and Spiliopoulos (2018), Al-Arifi et al (2018), Beck, E, and Jentzen (2019), Raissi, Perdikaris, and Karniadakis (2019), Hure (2019), and Wu and Xiu (2020).

the cross-section of returns. In addition this disaggregation of risk permits economists to be more specific about where disturbances to the economy originate and allows for heterogeneity in the exposure of production processes to ultimate sources of risk.

1.4.3 Benefits

In sum, this thesis makes an alternative methodology for studying ‘the relationship between asset prices and economic fluctuations’ ready for use, supplying a benchmark theoretical framework and tools for the numerical analysis of model solutions. We call this alternative paradigm ‘continuous-time macro-finance.’ Our alternative paradigm offers several benefits:

- Our model focuses attention on the risks of production. The process of production is disaggregated into multiple technologies. Agents control production through multiple decisions covering consumption, labor supply, and the allocation of capital and labor to alternative production technologies. Models are thus distinguished by how they specify aggregate production risk.
- We provide a general solution procedure that is scalable to large numbers of production processes and economic state variables. Indeed, on its own, the deep learning solution procedure is the most general method available to solve dynamic portfolio problems with state variable risks, a problem on which financial economists have long struggled to make progress (see Cochrane 2014 and references therein).
- Analytical work goes into formulating the model, but little is required to solve the model. One does not have to find first-order conditions, solve for a steady state, or linearize the model. Once the model is formulated as a system of stochastic differential equation, it can be solved.
- The lack of pre-solution analysis also means that the numerical solutions we provide are mostly free of artificial constraints, apart from

errors introduced by time discretization and Monte Carlo simulation. Errors from linear approximations and Taylor expansions in the neighborhood of the steady-state solution are not present.

I provide a proof of concept, a ‘sandbox’ for trying out different causal explanations that lead from production risks to asset prices in an environment of optimizing agents. Causal inferences come from trying out counterfactual scenarios in impulse response analyses, seeing how asset prices change in response to production risk. My overall theme is to bring a much richer specification of aggregate risk into the foreground of macro-finance. I see the results as contributions to a theory of risk-bearing, to echo the title of Arrow’s (1970) book.

Inevitably many questions will remain concerning the specification of the benchmark model. My overarching contribution is to provide an environment in which such questions may be asked and answered. I would like to see this thesis become the starting point for a continuous-time macro-finance research program.

1.5 Outline of the Thesis

We begin with a careful exposition, derivation and reinterpretation of the CIR85a model in Chapter 2. I situate the model in the context of the dynamic portfolio theory of Merton’s ICAPM and Breeden’s CCAPM, and show how an implicit solution is obtained for the stochastic control problem. Then I derive the implications of the model for asset pricing, including the equilibrium risk-free rate, the equity risk premium, and the equilibrium prices and risk premia for contingent claims. I emphasize the origins of aggregate risk in production, the determination of the risk-free rate from objective conditions, and the implications of the model for the cross-section of equity returns, which were largely passed over by the authors. In addition I will supply new proofs, probe some assumptions, and highlight a number of potential extensions.

Chapter 3 shows how to solve CIR models numerically. I present the

deep learning-based method and demonstrate its use in obtaining a solution of CIR85a. I specify the economic state variables, calibrate the model to historical levels for output and consumption, specify a set of production processes, and obtain a solution. I study the dynamic equilibrium of the CIR85a economy and dynamic equilibrium responses to exogenous shocks. The numerical solution and profiles of the optimal control decisions discovered by the deep neural network permit an explicit discussion of consumption and investment allocation decisions in CIR85a that were left implicit in CIR85a and CIR85b, and which do not otherwise exist in the literature. I compare and contrast my modeling approach and numerical procedure with standard numerical solution methods for DSGE models.

Chapter 4 generalizes the CIR85a model to include production with capital and elastically-supplied labor, making it ready for use as a paradigm theory comparable to RBC models in scope and complexity. The model incorporates multiple production processes with constant elasticity of substitution technology. Labor and capital are treated symmetrically by imagining labor as human capital that offers an uncertain yield in a particular production process. As a result of uncertainty concerning the yields of productive factors in each process, the output of each production process is uncertain. Additional uncertainty comes from the evolving state of technology, which has factor-augmenting and factor-neutral components.

I show in Chapter 4 that a calibrated continuous-time model with labor and capital does an excellent job of reproducing many characteristics of economic reality, including output, investment, consumption and labor force participation rates. Its consequences for asset prices are more exciting. My calibrated model produces reasonably-sized risk-free rates, equity risk premia, and option volatility smiles. Using an impulse response analysis, I show how movements in the equity risk premium depend on the technological characteristics of production and optimizing decisions undertaken by economic agents. In particular I show that movements in the equity risk premium can be decomposed into an ‘income effect’ driven by labor supply decisions and a ‘substitution effect’ driven by the need to reallocate productive resources among production processes.

The solution of my model in Chapter 4 suggests that fluctuations in TFP are insufficient to explain aggregate fluctuations in output and consumption. Instead, my model implies that aggregate fluctuations arise from breakdowns in the yield of capital within intensively-utilized production processes. Chapter 4 thus offers an alternative explanation of the business cycle, while showing that ‘idiosyncratic’ risks may play an underappreciated role in determining the equity risk premium.

Chapter 5 summarizes the characteristic features of the continuous-time macro-finance paradigm and collects a number of roads not traveled in the previous chapters as frontiers for future research. Following a discussion of the new paradigm’s limitations and weaknesses, I conclude by drawing out some implications of my modeling for economic policy.

Chapter 2

Production and Asset Prices in Continuous Time: A Portfolio-Theoretic Foundation for Macro-Finance

I believe Steve [Ross] understood, better than anybody else in his generation, the intuition of the financial market equilibrium...
– Jonathan Berk (2018: 71)

In this chapter we present the general equilibrium model of CIR85a. Our exposition makes the thesis self-contained and presents results we will refer to repeatedly in later chapters. However I do not merely rehearse the results of the original article. Most of the results in CIR85a are presented tersely. Patiently unpacking those results serves multiple purposes.

First, we obtain methodological guidance for the developments that follow. Stochastic calculus and control theory are not so familiar in macroeconomics that we can begin deriving results without further comment, and some familiarity with these methods helps to motivate our later chapters.

Second, the methods of proof employed by CIR85a help to build intuition about macro-financial equilibrium. One of the most remarkable aspects of CIR85a is the authors' ability to wring a steady stream of results from an

implicit solution. We revisit the main results, supplying more straightforward and explicit proofs. At the same time, our review highlights the gulf between the possibilities afforded by analytical and numerical solution methods, and just how much structure is missing from equilibrium for want of a full solution.

A third purpose is to re-interpret the central results. Interpretation highlights the original contributions of CIR85a, while putting them into dialogue with the literature surveyed in the Introduction. We give particular attention to the interpretation of the risk-free rate and the equity risk premium. The latter is hardly discussed in the original, which was primarily focused on deriving prices for contingent claims on the risk-free rate.

Finally, as we are elevating CIR85a to the status of a paradigmatic theory, we believe some amount of appreciation is also in order. To this end we situate the model in the context of earlier work on intertemporal equilibrium between consumption and portfolio investment by Samuelson and Merton, showing what CIR85a take from these formulations and where they go beyond these formidable first efforts. In addition we emphasize the aspects of CIR85a that make it a genuinely production-based theory, showing how all of the risks for which investors demand compensation may be traced to phenomena that generate fluctuations in output. In this sense we believe CIR85a supplies a deep causal foundation for asset pricing and macro-finance that is merely implicit in consumption-based theories. Our exposition of the model provides additional opportunities to develop this theme and others from the Introduction.

We begin with an exposition of the intertemporal portfolio problem, which provides historical and mathematical motivation for the CIR85a model. From there, we present the primitive assumptions of CIR85a and characterize equilibrium, giving particular attention to the risks of production and the ability to reallocate risk with the aid of contingent claims. Careful derivations of the equilibrium risk-free rate, equilibrium risk premia, and the master valuation equation for derivatives follow, with some signposting of junctures where closed-form results are necessarily replaced by numerical approximations in more complex versions of the model.

2.1 Intertemporal Equilibrium

The roots of CIR85a lie in the intertemporal portfolio theory of Robert C. Merton (1969, 1971, 1973a).¹ Merton formulated the problem of consuming and investing over time as a stochastic optimal control problem, bringing to bear the continuous-time mathematics of stochastic differential equations. In Merton (1971) he was able to claim that the intertemporal consumption-investment problem with fixed investment opportunities was essentially solved for the hyperbolic absolute risk aversion (HARA) class of utility functions, which includes constant absolute risk aversion (CARA) and constant relative risk aversion (CRRA) utility functions as special cases.

However it was in Merton (1973a) that Merton glimpsed the potential of his *intertemporal portfolio theory* to provide the basis of a *general equilibrium theory* connecting financial markets to the real economy. The problem of an individual trading current consumption for savings could be transposed into the problem of a representative agent trading current consumption for investment in a set of production technologies.² The menu of production technologies available to investors comprises the investment opportunity set, which evolves over time in response to exogenous movements in state variables. Accordingly the representative agent conditions his behavior on the state of the economy. Returns on investment add to the representative agent's wealth, which may then be allocated among production technologies according to the expected development of the investment opportunity set.

The intertemporal capital asset pricing model (ICAPM) Merton derived is not only a milestone for portfolio theory, but also a road not traveled for the development of dynamic general equilibrium models. This unexploited

¹Samuelson (1969) made an important early contribution as well, but acknowledged Merton had tackled the problem at a greater level of generality.

²The representative agent interpretation of control theory solutions has an equivalent interpretation of a social planner seeking to maximize social welfare. Wherever one says 'the agent chooses' one could also say something like 'society chooses.' We adopt the representative agent interpretation throughout, while noting the ambiguity. To avoid getting bogged down in methodological issues, we could adopt an attitude similar to potential theory where the agents under study act *as if* they are solving an optimization problem, just as particles in a system act as if they are trying to minimize an energy functional.

potential of portfolio theory as a basis for general equilibrium motivates our interest in CIR85a as the foundation of a core theory for macro-finance. The difficulty of obtaining solutions is the Achilles heel of intertemporal portfolio theory, however. In each of his papers Merton acknowledges the restrictions he must embrace in order to obtain an analytical solution to the intertemporal portfolio problem. Perhaps most importantly, Merton was not able to obtain a solution to the intertemporal equilibrium portfolio problem when the state variables are not constant. After Merton, researchers have struggled even to obtain numerical solutions to the intertemporal portfolio problem when the investment opportunity set changes over time in a predictable way.³

By contrast, CIR85a shows how much can be learned from an implicit solution of the intertemporal general equilibrium problem. CIR85a implicitly solves for an equilibrium with a full complement of state variables and changes in the investment opportunity set. The generality offered by CIR85a in defining the evolution of the economy's production possibilities is another unexploited source of portfolio theory's potential as core theory for macro-finance.

Then CIR85a adds two subtle but brilliant nuances to the picture of macro-financial equilibrium. First, they show that the dynamics of production and the investment opportunity set determine the risk-free rate endogenously. Though Merton (1973a) had recognized that the risk-free rate would be instantaneous in his model, Merton's risk-free rate is still exogenous. CIR85a further recognized that no net investment takes place at the risk-free rate in equilibrium.⁴ Thus the investment allocation problem faced by society concerns only the allocation of capital among risky assets. Merton's solutions assume a residual positive allocation to a risk-free asset, begging the question of net social investment in a risk-free asset and leaving an important feature of equilibrium mysterious.

Second, CIR85a integrated contingent claims as vehicles to transfer the risk of changes in the investment opportunity set. A set of contingent claims whose values fluctuate with changes in the investment opportunity set is

³Cochrane (2014) offers a brief review of the relevant literature.

⁴Both Merton (1973a) and CIR85a are moneyless economies.

therefore essential to a complete-markets characterization of general macro-financial equilibrium.⁵ Excess returns on contingent claims are determined in equilibrium, as are the risk-neutral values of a wide class of contingent claims encompassing futures, forwards, and options.⁶

Forgoing the possibility of a closed-form solution allowed CIR85a to adumbrate several crucial aspects of macro-financial equilibrium, as I will now show. Then in Chapter 3 I will solve the CIR85a model numerically, making the implicit solution explicit and demonstrating the power of the CIR85a framework with a concrete state process specification.

2.1.1 Risky production from capital

Following Merton (1973a), CIR85a embeds the portfolio problem in a model of general equilibrium by connecting consumption opportunities to capital allocations, investment returns and the economy's production technology. Returns on investment depend, in turn, on a set of variables describing the state of the economy. Uncertainty about production yields (output) and the state of the economy creates risks for current and future consumption. As a result, agents must be compensated for holding risky assets, and the risk premiums they demand are a function of production risk.

The specification given by CIR85a for the state variables defining the investment opportunity set is completely general, requiring only that the state variables follow a multidimensional Itô process. The state variables are not named or given a particular interpretation, leaving the researcher

⁵Assuming complete markets as a theoretical device is different from the empirical claim that markets are complete. We shall take pains to point out this distinction at other junctures below.

⁶A risk-neutral valuation of a contingent claim F is obtained by setting the drift of its stochastic differential dF equal to r , where r is the instantaneous risk-free rate. The transformation is justified on economic grounds by observing that contingent claims payoffs may be synthesized through continuous trading of a zero-cost replicating portfolio in a frictionless market. In the absence of arbitrage, a portfolio with no net investment should earn the risk-free rate of return. Mathematically, risk-neutral valuation changes the probability measure used by agents to evaluate the risks of uncertain outcomes. We shall see that this change of measure incorporates features of the agent's utility function, and that the difference between the physical and risk-neutral drifts is accounted for by equilibrium risk premiums.

free to specify the variables defining the state of the economy as desired. The specification of production processes will depend on the definition of the state variables. In the next chapter we will work through a concrete example which shows that the production and state variable processes may not be separately identifiable.

We begin by enumerating the sources of production risk in the economy. We define a Brownian motion dZ_t of length $N + K$, where N is the number of production processes employed in the economy, and K is the number of economic state variables. All of the stochastic differential equations in CIR85a are defined with respect to this Brownian motion and thus exposed to the same sources of uncertainty.⁷

Wealth, production, and state variables

The fundamental quantity in CIR85a is a vector η recording the amounts of capital invested in each of the economy's N production processes.⁸ Given total wealth W and a vector of allocations $a_i, i = 1, \dots, N$, with $0 \leq a_i \leq 1$ and $\sum_i a_i = 1$, the vector of investments η is defined by $\eta_i = a_i W$. The vector η therefore describes the allocation of society's wealth among available productive processes.

The capital stock is the sole factor of production in CIR85a. The rate of production is the yield of the capital stock, which is given in continuous time by the differential of η . Each individual η_i evolves according to a geometric Brownian motion

$$d\eta_i = \alpha_i(Y, t)\eta_i dt + \eta_i g_i(Y, t) dZ_t = \alpha_i(Y, t)a_i W dt + a_i W g_i(Y, t) dZ_t \quad (2.1)$$

where $\alpha_i(Y, t)$ is the expected rate of return on investment in process i and $g_i(Y, t)$ is a row vector of length $N + K$ describing the dependence of the return on process i on each of the $N + K$ sources of uncertainty in the

⁷However one should bear in mind that the *exposure* of each process to many of the $N + K$ sources of uncertainty will be zero.

⁸In the model, capital is anything that can yield output in the future and accumulate as wealth.

economy.⁹ Realized rates of return depend on particular realizations of dZ_t , which we might denote $dZ_t(\omega)$. The variance of $d\eta_i$ over multiple paths ω is therefore determined by $\eta_i g_i(Y, t)$.¹⁰

The expected rates of return on investment may be stacked in an N -vector $\alpha(Y, t)$ and the row vectors $g_i(Y, t)$ stacked to form the $N \times (N + K)$ matrix $G(Y, t)$. In this case the joint dynamics for returns on investment may be written

$$d\eta(t) = I_\eta \alpha(Y, t)dt + I_\eta G(Y, t)dZ_t \quad (2.2)$$

where I_η is an $N \times N$ diagonal matrix in which each entry is the amount reinvested in process i , $\eta_i = a_i W$.

Expected production yields $\alpha(Y, t)$ and their variability $G(Y, t)$ depend on a K -vector of exogenous variables Y that describe the state of the economy. As we pointed out above, the number and character of the state variables is undefined. Each state variable follows the diffusion process

$$dY_i(t) = \mu_i(Y, t)dt + s_i(Y, t)dZ_t \quad (2.3)$$

Again $\mu_i(Y, t)$ is the drift term, possibly time- and state-dependent, and $s_i(Y, t)$ is a row vector of length $N + K$. Stacking drift and diffusion terms, the K -dimensional state vector follows the process

$$dY(t) = \mu(Y, t)dt + S(Y, t)dZ_t \quad (2.4)$$

Though the state variables and the dependence of production processes on the state variables are as yet undefined, we may nevertheless interpret the

⁹Here and throughout, differential forms such as (2.1) are to be understood to mean the integral equation

$$\eta_i(t) = \eta_i(0) + \int_0^t \mu_i(Y(s), s)ds + \int_0^t \sigma_i(Y(s), s)dW_s$$

for a time interval $[0, t]$ where μ_i and σ_i are the drift and diffusion terms for the process followed by η_i , and the second integral is the Itô integral over a Brownian motion dW_t . The integral representation reminds us that the differential forms are really *accumulation* processes for the left-hand side variable.

¹⁰We omit the standard measure-theoretic incantations about ω .

model as generating a stochastic set of production (investment) opportunities characterized by $\alpha(Y, t)$ and $G(Y, t)$ from state variables following a general Itô process. While the processes for the states and production are both conditionally normal, the *unconditional* distribution of investment returns can assume quite general forms.¹¹

Aggregate risk

The outer product GG' is an $N \times N$ covariance matrix describing the risk and dependence structure of the investment returns, or the aggregate production risk that is inherent in the current production opportunity set. Similarly, a $K \times K$ covariance matrix SS' may be formed from the $K \times (N + K)$ matrix $S(Y, t)$ that summarizes the aggregate risk from uncertainty about the future state of the economy. Covariances between investment returns and state variables are given by the $N \times K$ matrix GS' , which may be understood as the risk to current output of adverse changes in the investment opportunity set.

Production risk is scaled by a chosen investment allocation a and the amount of wealth W . If we append the matrix S to the scaled matrix $a'GW$, aggregate economic risk is given by¹²

$$\begin{bmatrix} a'GG'aW^2 & a'GS'W \\ SG'aW & SS' \end{bmatrix} \quad (2.5)$$

We must bear in mind that the risk matrix above suppresses the dependence of G and S on the state variables Y .

The matrix (2.5) distinguishes the treatment of risk in portfolio theory models from the stochastic component of a typical real business cycle (RBC) model. The risk matrix provides a complete accounting of the economy's exposure to $N + K$ sources of uncertainty, whereas RBC models reduce uncertainty to technology and taste shocks. Its candid treatment of macroeconomic risk is another attractive feature of CIR85a as core theory. Merton

¹¹In fact, the risk profile of production in CIR85a may be viewed as a structural counterpart to Engle's (2008) dynamic conditional correlations approach.

¹²This is the definition of Σ in CIR85a: 380.

(1973a) and Breeden (1979) employ similarly general assumptions for aggregate risk, but do not consider whether the risk of changes in the investment opportunity set can be shared within the economy.¹³ The equilibria of Merton and Breeden are incomplete-markets equilibria with unspanned risks.¹⁴

Contingent claims

By introducing contingent claims markets, CIR85a are the first to define equilibrium with changing investment opportunities in a setting of complete markets. Contingent claims allow agents to adjust their risk exposures by exchanging claims on current output for state-contingent future claims. And in perhaps their most formidable achievement, CIR85a further show that the values of contingent claims are determined in equilibrium.¹⁵

A contingent claim i is defined by its payout profile $\delta_i(Y, t)$ and its economic risk exposure, given by the $N + K$ vector $h_i(Y, t)$. Its price F^i evolves according to the stochastic differential equation

$$dF^i = (F^i \beta_i(Y, t) - \delta_i(Y, t))dt + F^i h_i(Y, t) dZ_t \quad (2.6)$$

where $\beta_i(Y, t)$ is the expected rate of return. A contingent claim thus allows an agent to obtain an expected rate of return $\beta_i(Y, t)$ in exchange for holding the risk exposure defined by $h_i(Y, t)$ and payouts $\delta_i(Y, t)$.

Contingent claims are in zero net supply in equilibrium: there must be a long for every short. Therefore if we define an aggregate allocation vector b for contingent claims allocations, it must be the case that $b_i = 0$ for all i . However no obvious constraint exists on the number of contingent claims that may circulate in the economy, because payouts and risk exposures may

¹³Breeden (1979) appeared in print six years before CIR85a, but early versions of CIR85a and CIR85b were in circulation well before 1985. Breeden (1979) cites a 1977 draft.

¹⁴In other words, state-variable risks only enter the returns of tradeable securities through $a'GS'W$, which may be equal to zero if the state variables are static. If the state variables evolve over time and the basis of $a'GS'W$ is not also a basis for SS' , then some state-variable risks are unspanned by the universe of tradeable securities.

¹⁵“Although Arrow and Debreu first derived the concept of state prices, it was Steve [Ross] who first understood their importance in actually pricing financial assets.” (Berk 2018: 73)

be defined in any number of ways. If we were to stack the vectors $h_i(Y, t)$ to form a matrix H , the resulting matrix could have any number of rows. So rather than enumerate a specific set of contingent claims, CIR85a seeks a ‘basis set’ of contingent claims that spans whatever risks in the economy are not already spanned by other arrangements. If a basis set of claims can be characterized, then any other desired contingent payoffs can be synthesized by combining contingent claims in the basis set.

A bit of reflection on (2.5) helps to fix the dimension of the basis set and the definition of H . In equilibrium wealth is fully invested, so the economy must hold $a'GG'aW^2$ and $2a'GS'W$. Both G and S depend on Y . Accordingly all opportunities to hedge the risk of production must entail exposure to Y , the risk of which is captured in SS' , a square matrix of dimension K . Hence the basis set of contingent claims must also have dimension K . Thus without loss of generality we can stack K linearly independent row vectors $h_i(Y, t)$ to obtain a $K \times (N + K)$ matrix H , which may be identified with S .

Contingent claims pricing furnishes another window into the state variables Y . One is free to specify Y in a way that is most convenient to represent a given set of contingent claims. In applications, Y can capture spot prices and convenience yields for commodities, low- and high-frequency components of the term structure, or stochastic volatilities, among other examples.

2.1.2 Evolution of wealth and the control problem

As investment and consumption decisions are made over time, their consequences are accumulated in the level of wealth. Accordingly, wealth is an endogenous state variable for the economy.

The stochastic process followed by wealth is a geometric Brownian motion

given by

$$\begin{aligned}
dW &= W \left[\sum_{i=1}^N a_i(\alpha_i - r) + \sum_{i=1}^K b_i(\beta_i - r) + r - \frac{C}{W} \right] dt \\
&+ W \left[\sum_{i=1}^N a_i \left(\sum_{j=1}^{N+K} g_{ij} dZ_j \right) + \sum_{i=1}^K b_i \left(\sum_{j=1}^{N+K} h_{ij} dZ_j \right) \right] \\
&= W \mu(W) dt + W \sum_{i=1}^{N+K} q_i(W) dZ_{i,t}
\end{aligned} \tag{2.7}$$

The drift of the wealth accumulation process is comprised of four components: (1) the returns on invested capital in excess of the risk-free rate r from (2.1), (2) returns on contingent claims in excess of the risk-free rate per (2.6), (3) returns on risk-free lending, and (4) the rate at which gross investment returns are currently consumed, C . In equilibrium all wealth is invested in production so long as $\alpha_i(Y, t) > r$ for at least one i ; hence $\sum_i a_i = 1$. Contingent claims are in zero net supply, so $\sum_i b_i = 0$ and $b_i = 0$ for all i . Accordingly we can pull out r

$$-r1aW - r1bW + rW = 0$$

to see that r makes no net contribution to the drift. Evidently a more compact representation of the drift of (2.7) is

$$\left[\sum_{i=1}^N \alpha_i a_i W - C \right] dt \tag{2.8}$$

or the expected gross return on investment given current allocations, net of current consumption.

CIR85a interpret r as a risk-free rate at which agents may borrow and lend unlimited amounts. Under this interpretation, returns on capital α_i and returns on contingent claims β_i may be decomposed into the risk-free rate r and expected excess returns $\alpha_i - r$ and $\beta_i - r$ that compensate investors for holding risk. Incorporating r in the wealth evolution does not affect the dynamics of the economy, but it allows CIR85a to deduce expressions for the

risk-free rate and excess expected returns in equilibrium.

The diffusion term in (2.7) is given by two terms in dZ_j , reflecting how risk exposures from production and contingent claims, respectively, generate uncertainty about future levels of wealth. Once again, because $\sum_i b_i = 0$ in equilibrium, only the uncertainty generated by production will influence the evolution of aggregate wealth. Thus each production process contributes $a_i W g_i dZ$ to the diffusion term, and a more compact expression of the equilibrium wealth evolution process is simply

$$dW = \sum_{i=1}^N d\eta_i - C dt \quad (2.9)$$

or gross output minus consumption.

CIR85a formulate the determination of equilibrium as a stochastic control problem. We may interpret the problem as a computation undertaken by a representative agent or a social planner. Under the fairly strong assumption that agents agree about the structure and state of the economy, and therefore share common expectations concerning the evolution of the economy, the resulting equilibrium may be interpreted as a market outcome achieved by rational individuals without coordination, in which the wealth evolution equation is the budget constraint faced by a representative agent.¹⁶ For the remainder of the investigation we adopt the representative agent interpretation.

An optimizing representative agent controls the evolution of their wealth by choosing how much to consume (C), how to allocate investment among N opportunities (a), and which set of contingent claims (b) to hold. Choices are evaluated according to the cumulative utility they generate over the finite planning horizon. Define the class of value functions \mathcal{K} as

$$\mathcal{K}(\nu(t), W(t), Y(t), t) = E_{W,Y,t} \int_t^T U(\nu(s), Y(s), s) ds \quad (2.10)$$

¹⁶In particular, homogeneous expectations are needed to enforce the equilibrium condition $b_i = 0$ for all i , so heterogeneous expectations likely explain trade in contingent claims. We reserve questions about heterogeneous expectations for future work.

The class of value functions \mathcal{K} defines the utility of applying the control sequence $\nu(t)$ when the history of the economy is described by $W(t)$ and $Y(t)$. The function $\nu(t)$ that indexes the functions in the class \mathcal{K} is an admissible feedback control, meaning a choice of investment and consumption at time t that uses only information available at time t and maintains a non-negative level of wealth from t to T . The expectation operator $E_{W,Y,t}$ denotes the expectation conditional on the initial values of wealth W_t , the state variables Y_t and other information available at time t . Note that while the value functions \mathcal{K} are functions of the endogenous state variable $W(t)$, utility U is assumed not to depend on wealth. In principle, U may depend on any of the controls, the value of the state variables and time, but in what follows we assume, following CIR85a, that U depends on consumption alone and is time-separable.

The differential generator of $K \in \mathcal{K}$ associated with the control sequence $\nu(t)$, $\mathfrak{L}^\nu(t)K$, is

$$\begin{aligned} \mathfrak{L}^\nu(t)K &= \mu(W)WK_W + \sum_{i=1}^K \mu_i K_{Y_i} + \frac{1}{2}W^2K_{WW} \sum_{j=1}^{N+K} q_j(W) \\ &+ \sum_{i=1}^K WK_{WY_i} \sum_{j=1}^{N+K} q_j(W)s_{ij} + \frac{1}{2} \sum_{i=1}^K \sum_{j=1}^K K_{Y_i Y_j} \sum_{k=1}^{N+K} s_{ik}s_{jk} \end{aligned} \quad (2.11)$$

where μ_i is the drift of the i th state variable and the s_{ij} are the sensitivities of the i th state variable to source of uncertainty j , as given in (2.3), the terms $\mu(W)$ and $q_j(W)$ are defined by (2.7), and the subscripts on K denote partial derivatives with respect to the subscripted variables.¹⁷

If $J(W, Y, t)$ is the member of \mathcal{K} that solves the Hamilton=Jacobi-Bellman (HJB) equation

$$\max_{\nu \in V} [\mathfrak{L}^\nu(t)J + U(\nu, Y, t)] + J_t = 0 \quad (2.12)$$

¹⁷The differential generator may be seen as a device for applying Itô's lemma to a function of W and Y while holding t and $\nu(t)$ fixed.

with boundary conditions

$$\begin{aligned} J(0, Y, t) &= E_{Y,t} \int_t^T U(0, Y(s), s) ds \\ J(W, Y, T) &= 0 \end{aligned} \tag{2.13}$$

then the HJB equation is a martingale for an optimal control $\hat{\nu}$ and $J(W, Y, t) = K(\hat{\nu}, W, Y, t)$, subject to technical conditions given in CIR85a. Thus the optimal control $\hat{\nu}$ is the admissible feedback control that makes the drift of the HJB equation zero. Note also that the second boundary condition in (2.13) implies that wealth is depleted over the horizon $T - t$. In what follows, we refer to the optimal value function J as the indirect utility function for the representative agent.

2.1.3 First-order conditions and equilibrium

No constructive analytical procedure exists for finding the indirect utility function. An analytical solution may be found by making an inspired guess for the form of J , and applying an appropriate verification theorem to confirm the guess is correct.¹⁸ CIR85a do not follow this path and neither will I. Instead, CIR85a study the properties of a solution using the first-order conditions of the HJB equation. Though they do not obtain an explicit solution for a , b and C , CIR85a are able to derive expressions for the risk-free rate, risk premia and contingent claims prices in terms of *implicit* solutions for a , b , and C .

Expand (2.12) using (2.7) and (2.11) to write the HJB equation as

$$\begin{aligned} & \left[\sum_{i=1}^N a_i W (\alpha_i - r) + \sum_{i=1}^K b_i W (\beta_i - r) + rW - C \right] J_W + \mu' J_Y \\ & + \frac{1}{2} W^2 J_{WW} (a' G G' a + 2a' G H' b + b' H H' b) \\ & + (a' G S' + b' H S') W J_{WY} + \frac{1}{2} S S' J_{YY} + U(C, t) + J_t = 0 \end{aligned} \tag{2.14}$$

¹⁸See, for example, Rogers (2013) and Oksendal and Sulem (2019). Merton (1969, 1971) obtains closed-form solutions by this method.

The expressions J_{WY} and J_{YY} mean a K -vector of mixed partial derivatives $\frac{\partial^2 J}{\partial W \partial Y_i}$ and a $K \times K$ Hessian matrix $\frac{\partial^2 J}{\partial Y_i \partial Y_j}$, respectively. Similarly, J_Y is a K -vector of partial derivatives $\frac{\partial J}{\partial Y_i}$, while J_W and J_{WW} are scalar first and second derivatives of J with respect to wealth. The expansions in the J_{WW} , J_{WY} and J_{YY} terms follow from recognizing that $q_j(W)$ becomes $a'G + b'H$ under summation, while s_{ij} becomes S . Dependence of α , G , H , and S on Y and t has been suppressed, but should be borne in mind.

Differentiating the HJB equation with respect to the controls a , b and C yields the first-order conditions

$$\begin{aligned}\psi_a &= [\alpha - r1] W J_W + [GG'a + GH'b] W^2 J_{WW} + GS'W J_{WY} \leq 0 \\ \psi_b &= [\beta - r1] W J_W + [HG'a + HH'b] W^2 J_{WW} + HS'W J_{WY} = 0 \\ \psi_C &= U_C - J_W \leq 0\end{aligned}\quad (2.15)$$

In the first two expressions 1 is a conforming vector of ones. To ensure the non-negativity of consumption and investment allocations a_i , the complementary slackness conditions $C\psi_C = 0$ and $a'\psi_a = 0$ are also first-order conditions of optimality.¹⁹ The non-negativity constraint on investment reflects the impossibility of being ‘short’ a physical production process.²⁰ Contingent claims holdings are not subject to a sign restriction, which can be interpreted as allowing unlimited short sales.

Note that the terms of the Bellman equation that depend entirely on the state process, as well as the time derivative of the value function, drop out

¹⁹These conditions differentiate interior solutions with $a > 0$ and $\psi_a = 0$ (respectively, $C > 0$ and $\psi_C = 0$) from solutions on the boundary with one or more $a_i = 0$ ($C = 0$). Hardening the conditions to equalities rules out solutions on the boundary, requiring an interior solution.

²⁰If each element of η were to accumulate independently as $\eta_i(t + dt) = \eta_i(t) + a_i \sum_i d\eta_i(t + dt)$, the non-negativity constraint could be interpreted as an irreversibility condition, a possibility we take up in Chapter 4. As a further extension to such a model we might entertain the possibility of disinvestment from (liquidation of) a productive process, what Brunnermeier and Sannikov (2017) term ‘technological illiquidity’. In this case the values $-1 \leq a_i < 0$ would be admissible. Assuming a process-specific recovery rate ρ_i , disinvestment would increase the drift of the wealth equation by $\rho_i a_i W$ for all a_i with $-1 \leq a_i < 0$, subject to an additional constraint $a_{it} W_t + \eta_{i,t-1} \geq 0$, meaning that the amount of disinvestment may not exceed the amount originally invested in the process.

of the analysis. The state process is exogenous and therefore uncontrollable. In addition the time derivative J_t must be zero because J is a martingale for an optimal control sequence.

The first-order conditions must hold for all values of W , Y , and t . As a result, the derivatives of the indirect utility function J_W , J_{WW} and J_{WY} should also be read as functions of W , Y and t . Mathematically, the ‘local’ quality of the indirect utility function’s slope (J_W) and curvature (J_{WW} and J_{WY}) induces state-dependence in the derivatives that form the market prices of risk, as we shall see below.

In equilibrium, wealth is fully invested and contingent claims are in zero net supply. Accordingly $\sum_i a_i = 1$ and $b_i = 0$, $i = 1, \dots, K$ in equilibrium. The first-order conditions therefore simplify to

$$\begin{aligned}\psi_a &= [\alpha - r1] W J_W + GG' a W^2 J_{WW} + GS' W J_{WY} \leq 0 \\ \psi_b &= [\beta - r1] W J_W + HG' a W^2 J_{WW} + HS' W J_{WY} = 0 \\ \psi_C &= U_C - J_W \leq 0\end{aligned}\tag{2.16}$$

while the complementary slackness conditions remain unchanged.²¹

Our next task is to deduce the endogenous risk-free rate r and the endogenous rate of return on contingent claims β consistent with equilibrium. As noted above, CIR85a does not derive an explicit solution for the value function J or an optimal control policy. Instead, they proceed with the analysis of an implicit solution $a^*(W, Y, t)$, $b^*(W, Y, t)$, $C^*(W, Y, t)$.²² We now show how expressions r and β may be deduced from the implicit solution.

²¹It may not be obvious why $b_i = 0$ for the representative agent as well as the economy as a whole. If all agents hold the same allocation a they have the same exposure to production risk and the risk of changes in the investment opportunity set. Further, because they have identical utility functions, all agents are equally averse to the risks they bear in equilibrium. Under these conditions, any contingent claim desired would not be supplied, and any contingent claim offered would not be demanded. Thus equilibrium in the contingent claims market is a no-trade equilibrium. Demand for contingent claims arises either because agents are heterogeneous or because an obstacle to the equilibrium allocation exists, e.g., when an agent must hold a certain component of production risk. We believe this aspect of the complete markets equilibrium is worthy of further investigation.

²²To unclutter the notation we dispense with the stars in the sequel. It should be clear from the context that we intend optimal quantities.

2.2 The Equilibrium Risk-Free Rate

It is so common for stochastic control problems and dynamic economic models to include a discount factor in the value function that the absence of a discount factor from the CIR85a formulation may be surprising. However no discount factor is needed to ensure integrability in a finite-horizon problem, and in one of their central results, CIR85a show that intertemporal rates of substitution and the risk-free rate are determined in equilibrium even in the absence of a subjective discount factor. Time preference is an objective quantity – a price given in equilibrium – rather than a subjective measure of impatience. Accordingly the risk-free rate of CIR85a more closely resembles Wicksell’s natural rate of interest than Fisher’s rate of time preference or the Austrian agio.

Consider the first-order conditions on investment. Using (2.16), the complementary slackness condition on reinvestment reads

$$a'\psi_a = a'[\alpha - r1]WJ_W + a'GG'aW^2J_{WW} + a'GS'WJ_{WY} = 0$$

Rearranging and dividing through by WJ_W , we immediately have

$$\begin{aligned} r &= a'\alpha + a'GG'aW \left(\frac{J_{WW}}{J_W} \right) + a'GS' \left(\frac{J_{WY}}{J_W} \right) \\ &= a'\alpha - a'GG'aW \left(-\frac{J_{WW}}{J_W} \right) - a'GS' \left(-\frac{J_{WY}}{J_W} \right) \end{aligned} \quad (2.17)$$

because $a'1 = \sum_i a_i = 1$ in equilibrium. We assume $J_W > 0$ and $J_{WW} < 0$ by concavity. Thus $-\frac{J_{WW}}{J_W}$ is a positive number. The terms $-\frac{J_{WY}}{J_W}$ can have either sign.

When J is a function giving the utility of wealth, expressions of the form $-\frac{J_{WW}}{J_W}$ are known as Arrow-Pratt absolute risk aversion functions. (Pratt 1964, Arrow 1971) The ratios measure the extra compensation a risk-averse individual must be paid to take a risky gamble. Because we assume that our agent’s utility function is representative, all agents in the economy share the same degree of risk aversion, and $-\frac{J_{WW}}{J_W}$ and $-\frac{J_{WY}}{J_W}$ may be interpreted as the market prices of risk prevailing at a given level of social wealth. (Ingersoll

1987: 37-39)

Write the variance of a variable X as $\sigma^2(X)$ and the covariance of two variables X and Y as $\sigma(X, Y)$. Observing from (2.5) that $\sigma^2(W) = a'GG'aW^2$ and $\sigma(W, Y) = a'GS'W$, we obtain the alternative expression

$$r = a'\alpha - \frac{\sigma^2(W)}{W} \left(-\frac{J_{WW}}{J_W} \right) - \sum_{i=1}^K \frac{\sigma(W, Y_i)}{W} \left(-\frac{J_{WY_i}}{J_W} \right) \quad (2.18)$$

for the risk-free rate of interest.

2.2.1 Interpreting the risk-free rate

Thus we have an expression for r composed of three terms. The first, $a'\alpha$, is the expected rate of return on optimally-invested wealth. The second is the variance of wealth relative to its level, multiplied by the market price of risk for wealth. The third is a sum of state variable covariances with respect to wealth, relative to the level of wealth, and multiplied by a state variable-specific market price of risk. The variance of wealth must be non-negative, and its price is positive. The expected rate of return on optimally-invested wealth must also be positive, or investment would be wealth-destroying. The covariance terms and their prices are unrestricted in sign.

Analysis of (2.17) and (2.18) shows that the second and third terms defining r ‘de-risk’ the expected rate of return on investment over all sources of risk in the economy, at prices determined by the degree of risk aversion exhibited by the representative agent. The de-risking adjustments transform the drift of the investment process under the physical measure, defined by $\alpha(Y, t)$, to a drift under the risk-neutral measure defined by $r(W, Y, t)$. On these grounds, r is a rate of return that may be identified with the risk-free rate of interest. The risk-free rate defined by the model is not risk-free in the sense of being exogenous, constant, or invariant to the state of the economy, leading to a known return over a certain holding period. Rather, r is the risk-neutral drift that instantaneously defines equilibrium contingent claims prices, as we will show below.

The derivation provides some insight into the dynamics of the risk-free

rate. Holding the expected rate of return on the productive processes and the covariance terms fixed, risk-free rates fall when the volatility of wealth increases or when the market price of risk for wealth increases. Falling risk-free rates reflect a decline in the incentive required for investors to hold risk-free claims when they are less enthusiastic about their investment opportunities.

Now consider the covariance terms. State variables influence the risk-free rate according to their covariance with wealth. A state variable Y_i that increases when wealth decreases has a negative covariance with wealth and provides a kind of insurance against adverse outcomes. For such a state variable Y_i , $J_{WY_i} > 0$ and $-\frac{J_{WY_i}}{J_W} < 0$; that is, the market price of risk is negative because a risk-averse investor would pay for the opportunity to receive payouts connected to Y_i , rather than demanding additional compensation. Therefore the net contribution of Y_i to the risk-free rate, $\sigma(W, Y_i) \left(-\frac{J_{WY_i}}{J_W}\right)$ is to lower the rate. State variables that vary positively with wealth will have positive covariances and positive market prices of risk, so they, too, will tend to lower the risk-free rate. Accordingly the expected rate of return $a'\alpha$ is reduced by $K + 1$ positive increments defined by the endogenous and exogenous state variables W and Y , respectively, to reach the certainty-equivalent return r .

Building on an observation of Breeden (1979), changes in the risk-free rate may also be connected to changes in optimal consumption plans. By the first-order conditions (2.16) we have $J_W = U_C$ for an interior solution, an equality known as the envelope condition. For an optimal consumption plan $C = C(W, Y, t)$, application of the chain rule gives $J_{WW} = U_{CC}C_W$ and $J_{WY} = U_{CC}C_Y$. Substituting into (2.18), the expression for the risk-free rate becomes

$$r = a'\alpha - \frac{C_W \sigma^2(W)}{W} \left(-\frac{U_{CC}}{U_C}\right) - \sum_{i=1}^K \frac{C_{Y_i} \sigma(W, Y_i)}{W} \left(-\frac{U_{CC}}{U_C}\right) \quad (2.19)$$

The substitution shows that even if the numeraire is changed from the marginal indirect utility of wealth to the marginal utility of consumption, yielding market prices of risk $-U_{CC}/U_C$, quantities of risk are still defined

by the state variables, scaled by the derivatives of consumption with respect to W and Y . The quantities of risk will not, in general, be reducible to variances and covariances with respect to consumption, showing that the consumption-based theory does not encompass CIR85a.

Further insight into the risk-free rate comes from rewriting the complementary slackness condition $a'\psi_a = 0$ and applying duality theory from nonlinear programming. Observe in (2.16) that the term rWJ_W may be pulled out to the right-hand side of ψ_a , leaving all of the decision variables on the right-hand side. Accordingly we can consider the dual of (2.16) in which $\lambda = rWJ_W$ is varied, thereby relaxing the constraint ψ_a . Substituting λ into $a'\psi_a$ gives

$$a'\alpha WJ_W + a'GG'aW^2J_{WW} + a'GS'WJ_{WY} \leq \lambda$$

Once again, $a'\lambda 1 = \lambda$ when wealth is fully invested. Now put $\gamma = \alpha WJ_W + GS'WJ_{WY}$ and $D = \frac{1}{2}GG'aW^2J_{WW}$. The condition $a'\psi_a$ now has the quadratic form

$$a'\gamma + a'Da$$

When this expression is maximized subject to $a'1 = 1$ and $a \geq 0$, a solves a portfolio optimization problem in the form of a quadratic program, as in Markowitz (1952). The λ we have defined is the difference between $a'\gamma + a'Da$ and its optimum. Hence λ is a slack variable, the shadow price of an additional unit of investment when capital is optimally invested.

The risk-free rate of CIR85a is not the risk-free rate of Merton (1973a) or Breeden (1979). Though Merton and Breeden alike acknowledge the instantaneous quality of the risk-free rate, both authors treat it as exogenous, while CIR85a shows that it is endogenous. Merton and Breeden also tend to identify r with the return on a specific short-term investment in, e.g., Treasury bills, rather than recognizing it as a certainty-equivalent return that depends on economic risk.

The risk-free rate of CIR85a is not the risk-free rate of neoclassical or neo-Keynesian macroeconomics, either. In the neo-Keynesian view, the risk-free rate is undetermined in the economy and must be set by the authorities.²³

²³See, for example, Gali (2015)

CIR85a shows that the risk-free rate set by the authorities need not coincide with the equilibrium risk-free rate, which is determined by the state of the economy. In the neoclassical view the dynamics of the risk-free rate are determined by the subjective rate of time preference and the supply and demand for funds in a perfect capital market. CIR85a acknowledges the connection between the interest rate and the capital market by identifying the risk-free rate with the shadow price of capital, but goes further by tying the dynamics of the risk-free rate to quantities of risk, prices of risk, and the level of wealth—that is, to phenomena that can be read off from asset prices.²⁴

2.2.2 The market rate of return

Reversing the argument for the risk-free rate of return in (2.17) and (2.18) shows that the expected rate of return on the *market* is equal to the risk-free rate plus $K + 1$ *risk premia*. CIR85a is ultimately interested in the theory of the term structure elaborated in CIR85b, so CIR85a passes over discussion of the market rate.²⁵ Yet CIR85a includes a set of important consequences for equity markets. In order to bridge the gap between CIR85a and the macro-finance literature, we devote some space to CIR85a’s equity market implications here.

Since a describes the optimal allocation of capital, $a'\alpha$ may be interpreted as a capitalization-weighted index of individual process returns. Inverting the derivation of the risk-free rate in (2.18) gives a decomposition of the market rate of return into the risk-free rate and $K + 1$ premiums for risk:

$$a'\alpha = r + \frac{\sigma^2(W)}{W} \left(-\frac{J_{WW}}{J_W} \right) + \sum_{i=1}^K \frac{\sigma(W, Y_i)}{W} \left(-\frac{J_{WY_i}}{J_W} \right) \quad (2.20)$$

One can sum the $K + 1$ premiums to obtain a single ‘market risk premium.’ In this case—which we might term the CAPM case—the market risk premium compensates investors for the variance of wealth arising from current pro-

²⁴We intend to pursue the consequences of these observations for monetary policy in future work.

²⁵CIR85a: 363n1 indicates the paper is “an extended version of the first half of an earlier working paper titled “A Theory of the Term Structure of Interest Rates.””

duction *and* for changes in the investment opportunity set. Alternatively, one may leave the $K + 1$ premiums disaggregated and treat each one as a distinct risk factor, as in the arbitrage portfolio theory of Ross (1976).

Studies of the cross-section of equity returns rarely attempt to interpret risk premiums in terms of state variables per CIR85a. Fama and French (1996: 76-77) interpret their three-factor model in terms of the APT, and conjecture that their value factor may be related to “relative distress” while emphasizing that they “have not identified the two state variables of special hedging concern to investors that lead to three-factor asset pricing.” Their reference to two rather than three state variables suggests that they identify the market factor with the wealth variable. However Fama and French also allow for irrationality and data problems as potential alternative explanations for the empirical success of the three-factor model.

Further rearrangement of (2.20) shows precisely the sense in which CIR85a is a production-based theory of asset pricing. Recall from (2.2) that production (output) is given by $d\eta$, the yield on capital employed in all available production processes. We know from (2.5) that $GG' = \sigma^2(d\eta)W^2$ is the variance of output, and $GS' = \sigma(d\eta, dY)W$ captures the covariance of output with the state of the economy. In addition, $\eta = aW$. Hence the expression (2.20) may be rewritten

$$\alpha = r1 + \frac{\sigma^2(d\eta)\eta}{W^2} \left(-\frac{J_{WW}}{J_W} \right) + \frac{\sigma(d\eta, dY)}{W} \left(-\frac{J_{WY}}{J_W} \right)$$

For an investment equal to one unit of wealth, $W = 1$, this simplifies to

$$\alpha = r1 + \sigma^2(d\eta)\eta \left(-\frac{J_{WW}}{J_W} \right) + \sigma(d\eta, dY) \left(-\frac{J_{WY}}{J_W} \right) \quad (2.21)$$

In (2.21) it is clear that investors are compensated for the variance of output and its covariance with the state variables. The amount of compensation they demand per unit of production risk is determined by the indirect utility of consumption, defined over a space of wealth and investment opportunities.

In the case of a single process we have

$$\begin{aligned} \alpha_i = & r(W, Y, t) + \eta_i g_i(Y, t) g_i(Y, t)' \left(-\frac{J_{WW}}{J_W} \right) \\ & + \eta_i \sum_{j=1}^K g_i(Y, t) s_j(Y, t)' \left(-\frac{J_{WY_i}}{J_W} \right) \end{aligned} \quad (2.22)$$

which shows that, given a choice of basis for Y , investors' valuations of claims on each production process will be determined by the $N + K$ dimensional vector $g_i(Y, t)$, which specifies how output varies not only with respect to the state of the economy, but also the intrinsic uncertainty of production and covariances of this intrinsic uncertainty across production processes. These latter sources of risk tend to be overlooked in studies of the cross section of equity returns.

2.3 Contingent Claims Valuation

Now let us consider contingent claims valuation. Our proofs will take a somewhat different itinerary than CIR85a in order to focus attention on the risk premia ϕ_W and ϕ_Y . Our analysis shows that ϕ_W and ϕ_Y are sufficient statistics that summarize the asset-pricing consequences of the CIR85a model, when used in conjunction with the risk-free rate r , optimal controls a and C , and the system dynamics $\{\mu(Y), S(Y), \alpha(Y), G(Y)\}$.

2.3.1 Contingent claims risk premia

Because all risk in the economy is captured by variations in the endogenous state variable W and in the exogenous state variables Y , we can determine all of the risk premia for the economy by considering the joint dynamics of the processes for W and Y .

The partitioned covariance matrix has the form (2.5), where $\sigma^2(W) = a'GG'aW^2$ is a scalar, $a'GS'W = \sigma(W, Y) = (SG'aW)'$ is a $1 \times K$ row vector, and $SS' = \sigma^2(Y)$ is a $K \times K$ matrix. The market prices of risk associated with the disturbance matrix are collected in the $(1 + K)$ -vector

$\left[(-\frac{J_{WW}}{J_W}), (-\frac{J_{WY}}{J_W})'\right]$.²⁶ Multiplying the two gives our solution for the risk premia that prevail in equilibrium²⁷:

$$\begin{aligned} \begin{bmatrix} \phi_W \\ \phi_Y \end{bmatrix} &= \begin{bmatrix} a'GG'aW^2 & a'GS'W \\ SG'aW & SS' \end{bmatrix} \begin{bmatrix} -\frac{J_{WW}}{J_W} \\ -\frac{J_{WY}}{J_W} \end{bmatrix} \\ &= \begin{bmatrix} a'GG'aW^2(-\frac{J_{WW}}{J_W}) + a'GS'W(-\frac{J_{WY}}{J_W}) \\ SG'aW(-\frac{J_{WW}}{J_W}) + SS'(-\frac{J_{WY}}{J_W}) \end{bmatrix} \\ &= \begin{bmatrix} \sigma^2(W)(-\frac{J_{WW}}{J_W}) + \sigma(W, Y)(-\frac{J_{WY}}{J_W}) \\ \sigma(W, Y)'(-\frac{J_{WW}}{J_W}) + \sigma^2(Y)(-\frac{J_{WY}}{J_W}) \end{bmatrix} \end{aligned} \quad (2.23)$$

The scalar ϕ_W divided by W is the market risk premium that appears in (2.20). It reduces the rate of return on investment under the physical measure to the risk-neutral rate of return, as shown in (2.17) and (2.18). Thus the risk-free rate is equal to $a'\alpha - \phi_W/W$.

The K -vector ϕ_Y expresses risk premia for each of the K state variables in terms of the same $K + 1$ variables as ϕ_W . The vector ϕ_Y may be expanded as follows:

$$\begin{bmatrix} \sigma(W, Y_1) \\ \vdots \\ \sigma(W, Y_K) \end{bmatrix} \left(-\frac{J_{WW}}{J_W}\right) + \begin{bmatrix} \sigma^2(Y_1) & \cdots & \sigma(Y_1, Y_K) \\ & \ddots & \\ \sigma(Y_1, Y_K) & \cdots & \sigma^2(Y_K) \end{bmatrix} \begin{bmatrix} -\frac{J_{WY_1}}{J_W} \\ \vdots \\ -\frac{J_{WY_K}}{J_W} \end{bmatrix} \quad (2.24)$$

In this form, the state variable risk premiums are clearly determined by the covariances of each Y_i with W and Y .

Now return to the first-order conditions of (2.16). Rearranging the second equation to solve for β , we obtain

$$\beta = r1 + HG'aW \left(-\frac{J_{WW}}{J_W}\right) + HS' \left(-\frac{J_{WY}}{J_W}\right) \quad (2.25)$$

Comparing this expression with (2.23), it is clear that $\beta = r1 + \phi_Y$ when $H = S$, as we claimed above.

²⁶This is the definition of χ given at CIR85a: 380.

²⁷Together with Σ defined in (2.5), χ defines Arrow-Debreu state prices, a connection further discussed in Section 4 of CIR85a.

Now consider a specific contingent claim F . The sensitivity of F to W and Y is determined by the contractual terms of F and measured by the derivatives F_W and F_Y . The market risk premia for W and Y are given by ϕ_W and ϕ_Y . Hence the expected rate of return on F is $rF + \phi_W F_W + \phi_Y F_Y$. But we have already claimed in (2.6) that the rate of return on a contingent claim is βF . Thus:

$$\beta F = rF + \phi_W F_W + \phi_Y F_Y \quad (2.26)$$

This is Theorem 2 of CIR85a. It is also clear that β_i is the expected rate of return on an asset that varies one-for-one with state variable Y_i . Similarly, ϕ_{Y_i} is the expected excess rate of return on such an asset.²⁸

We have unearthed two important results. First, a set of contingent claims is complete if it spans the space of state-variable risks. As a result we can set $H = S$ without loss of generality. Second, we have derived equilibrium risk premia ϕ_Y for state-variable risks and constructed contingent claims for which β is the vector of expected excess returns.

The correspondence of ϕ_Y to observable quantities depends on (1) the particular contingent claims basis employed, as well as (2) the existence of markets for particular contingent claims. In the foregoing we chose the basis $H = S$ because it was the simplest way to span the space of state variable risks. However there are many ways to specify S . In addition, an arbitrary invariant transformation of H (call it H^*) will also span S .²⁹ For example, S may be measured in terms of productivity, while H^* might be measured in terms of commodity prices and economic indices. If H^* spans S , then contingent claims written on commodity prices and economic indices can serve to make state variable risks tradeable. In this case, we would need a mapping $S \rightarrow H^*$ to transform ϕ_Y to ϕ_{H^*} for the sake of observation.

The second condition on observability concerns the existence of markets, and may be seen as a restriction on potential transformations $S \rightarrow H^*$. One does well to choose a basis that aligns with contingent claims that are actually traded. At the same time, any H^* used for empirical testing will

²⁸CIR85a: 375 write “ ϕ_{Y_j} is the expected excess return on a security constructed so that its value is always equal to Y_j ”, which assumes the initial values are equal as well.

²⁹For example, multiplication by an invertible matrix and translation by a vector.

entail a residue of risks ϕ_{Y_i} that are partially spanned or unspanned, because markets are unlikely to be dynamically complete in practice.

Some care in interpreting ϕ_Y is therefore warranted. The risks priced by ϕ_Y are inherent in the structure of production and must be borne by someone. Some of the risks in ϕ_Y may be transferred in contingent claims markets defined in terms of arbitrary observable variables. The extent of the risk transfer will be defined by the projection of the observable variables on the economic state variables. Empirical studies of contingent claims prices (e.g., Dai and Singleton (2000) and Casassus and Collin-Dufresne (2005)) show how this indeterminacy of basis may be handled through the choice of canonical forms and careful specification analysis.

2.3.2 Master valuation equation

In equilibrium contingent claim values are functions of the state variables, $F(W, Y, t)$. Because $F(W, Y, t)$ is a stochastic variable, one can use Itô's formula to write its drift as

$$\begin{aligned} & \frac{1}{2}\sigma^2(W)F_{WW} + \sigma(W, Y)F_{WY} + \frac{1}{2}\sigma^2(Y)F_{YY} \\ & + [\mu_W - \phi_W] F_W + [\mu'_Y - \phi'_Y] F_Y + F_t \end{aligned}$$

The means multiplying F_W and F_Y have been decomposed into the physical drifts of the state variables μ_W and μ'_Y and adjustments for risk premia ϕ_W and ϕ'_Y . Accordingly we have translated the drifts of W and Y from the physical measure to the risk-neutral measure. Under the risk-neutral measure, the drift of F must be equal to the risk-free rate r , minus any payouts given by $\delta(W, Y, t)$. Thus we have the equality

$$\begin{aligned} & \frac{1}{2}\sigma^2(W)F_{WW} + \sigma(W, Y)F_{WY} + \frac{1}{2}\sigma^2(Y)F_{YY} + [\mu_W - \phi_W] F_W \\ & + [\mu'_Y - \phi'_Y] F_Y + F_t = r - \delta(W, Y, t) \end{aligned}$$

The terms $\sigma^2(W)$, $\sigma(W, Y)$, and $\sigma^2(Y)$ are defined by the joint covariance matrix (2.5). We have calculated the risk premia ϕ_W and ϕ_Y in (2.23) above.

Since μ_Y is undetermined, it remains to recall from (2.8) that

$$\mu_W = a'\alpha W - C$$

Substituting in for $\sigma^2(W)$, $\sigma(W, Y)$, $\sigma^2(Y)$, μ_W , ϕ_W and ϕ_Y results in the master valuation equation for contingent claims (CIR85a Theorem 3):

$$\begin{aligned} & \frac{1}{2}a'GG'aW^2F_{WW} + a'GS'WF_{WY} + \frac{1}{2}SS'F_{YY} \\ & + \left[a'\alpha W - a'GG'aW^2 \left(-\frac{J_{WW}}{J_W} \right) - a'GS'W \left(-\frac{J_{WY}}{J_W} \right) - C \right] F_W \\ & + \left[\mu'_Y - \left(SG'aW \left(-\frac{J_{WW}}{J_W} \right) + SS' \left(-\frac{J_{WY}}{J_W} \right) \right)' \right] F_Y \\ & + F_t - r(W, Y, t) + \delta(W, Y, t) = 0 \end{aligned} \quad (2.27)$$

In this form one recognizes the first three terms multiplying F_W as rW . The solution of the valuation equation depends on a boundary condition describing the terminal value of the contingent claim.

2.3.3 Specific contingent claims

The master valuation equation shows explicitly how risk-neutrality is obtained in equilibrium. When pricing contingent claims one replaces the drift of state variables under the physical measure μ with their risk-neutral counterparts. Theorem 3 shows that the adjustment to achieve the risk-neutral drift is a linear combination of risk aversion coefficients (a feature of the indirect utility function) and a set of sufficient statistics for risk (a feature of production risk) in equilibrium. The terms in this adjustment are exactly the terms defining the expected excess returns in Theorem 2.

The price $F(W, Y, t)$ of any contingent claim must satisfy the master valuation PDE. Specializing the equation to the price of a particular contingent claim is achieved through the following:

1. *Terms of the contingent claim contract.* Because contingent claim payoffs are usually defined with respect to a subset of the economic state variables, the derivatives F_W , F_{Y_i} , F_{WW} , F_{W, Y_i} , and F_{Y_i, Y_j} will pick out a

subset of terms, while setting many others to zero. Nonzero derivatives are specific to the contract, as is the time- and state-dependent payout policy $\delta(W, Y, t)$. The contract will also specify any initial, boundary, and terminal conditions to be satisfied by the partial differential equation.

2. *Equilibrium conditions.* Equilibrium instantaneously defines the risk-free rate $r(W, Y, t)$ and optimal consumption policy $C(W, Y, t)$. The variance $\sigma^2(W)$ and covariances $\sigma(W, Y_i)$ depend indirectly on optimal investment policy $a(W, Y, t)$, as indicated in the definitions above, as does the risk-free rate. Likewise the coefficients of risk aversion $-J_{WW}/J_W$ and $-J_{WY}/J_W$ are determined at the optimum.
3. *The structure of economic risk.* Finally the covariances $\sigma(Y_i, Y_j)$ are given by SS' , with $S = S(Y, t)$. Note that these terms were not present in the formula for the interest rate, where only the covariances of the states with wealth were priced.

It is clear that a number of familiar asset-pricing models are embedded in the master valuation equation, from the Black-Scholes-Merton PDE (Black and Scholes 1973, Merton 1973b) to the default-free term structure model of CIR85b to the class of exponential-affine models (Duffie and Kan 1996, Schwartz 1997, Duffie 2001).

2.4 Conclusions

We have shown how the model of CIR85a determines the risk-free rate and a complete set of risk premia in general equilibrium as a function of the production opportunity set, an optimal allocation of capital, and an optimal consumption plan. In particular we have emphasized how CIR85a obtains prices for all risks in the economy by completing the financial market structure with contingent claims.

Production is the origin of risk in general equilibrium. Production risk comes from uncertainty about output under current production arrange-

ments, as well as the risk that future production possibilities may be inferior to those available today. The price of risk is reckoned in terms an indirect utility function defined in production opportunity space; for a given set of production opportunities, risks may be computed in consumption space as well. Though prices of risk depend on the properties of investor risk aversion and the subjective disutility of changing consumption plans, the *objective, structural* risk faced by all investors originates in the structure of production.

CIR85a gives an accounting of risk premia in an environment of complete financial markets as the financial correlate of macroeconomic equilibrium. Accordingly it orients a search for economic risk that amounts to an unpacking of the equity risk premium and risk premia for contingent claims. Furthermore, the model allows economic risk to take a far more general form than the shocks to autoregressive total factor productivity used as a summary statistic for risk in real business cycle models. Regardless of whether the additional risks contemplated in CIR85a are transferable through empirical contingent claims markets, the existence of such risks cannot easily be ignored and we should be concerned with who holds such risks in equilibrium.

We may also say that CIR85a reveals that the investment opportunity set is treated as a static quantity in real business cycle models, reinforcing our suspicion that the equity risk premium ‘puzzle’ is an artifact of a model with too little structural risk.

In the next chapter we show how CIR85a may be solved numerically, affording a more complete characterization of macro-financial equilibrium in the model.

Chapter 3

A Deep-Learning Solution Procedure for Continuous-Time Macro-Finance Models

CIR85a formulates general equilibrium as a stochastic control problem. The Hamilton-Jacobi-Bellman (HJB) equation (2.14) with boundary conditions (2.13) summarizes the problem to be solved. Given an HJB equation, there are two established plans of attack to obtain an explicit solution to the stochastic control problem.

The first plan of attack yields an analytical solution: Guess the form of the value function J , and then verify that it is a solution by taking partial derivatives and checking other conditions. Merton (1969, 1971) solves the intertemporal consumption-investment problem with constant state variables and CRRA/HARA utility by this method. In some cases, the form of J can be determined up to one or more unknown functions of time. Using the method of undetermined coefficients, one can obtain a system of ordinary differential equations (ODEs) that must be satisfied simultaneously by the value function. The solution obtained is analytical, but numerical methods must be used nevertheless because the solution to the system of ODEs is usually inscrutable in closed form. Kim and Omberg (1996), Sangvinatsos and Wachter (2005), and Liu (2007) follow this latter approach.

When analytical solutions are not forthcoming, the second plan of attack is to treat the HJB equation as a nonlinear partial differential equation (PDE) to be solved with standard numerical PDE methods: that is, to find an approximation of J numerically. Upon discretizing W , Y , and t , the problem of finding J can be reduced to linear algebra using finite difference approximations for the unknown partial derivatives and a suitable time-stepping procedure, for example. Brennan, Schwartz and Lagnado (1997) and Barberis (2000) solve the intertemporal consumption-investment problem using numerical PDE methods.

Neither of the above approaches qualifies as a general-purpose solution method for the consumption-investment problem with time-varying investment opportunities. The most advanced analytical solution method of which we are aware, that of Liu (2007), is staked to a particular choice of utility function and return dynamics, and does not satisfy the CIR85a constraints on the investment allocation vector. Liu (2007) does not present a numerical implementation of his solution, and it is not clear how the properties of a solution might be analyzed. PDE methods, on the other hand, quickly run into a ‘curse of dimensionality’ as the state vector Y increases in size. Brennan, Schwartz and Lagnado (1997) deal with the $K = 1$ case, as does Barberis (2000). Even for $K = 2$ standard numerical PDE methods become challenging, to say nothing of problems with very large K .

In this chapter we introduce a general numerical procedure to calculate the solution to the intertemporal consumption-investment problem. The procedure, based on that of Han and E (2016), employs deep learning methods to approximate the functions that solve the optimal control problem. Deep learning methods have attracted a great deal of interest among applied mathematicians as a means for obtaining numerical solutions to nonlinear PDEs in very high-dimensional settings, and in other applications where the solution to a problem is a function.

Deep learning uses neural networks to approximate functions. Mathematically, a neural network is a linear transformation of the input passed through a nonlinear ‘activation function.’ Layering neural networks on top of one another and daisy-chaining the outputs makes the network ‘deep’ and

aids in finding the optimal approximation efficiently. Optimal parameters are found for all layers of the network via stochastic gradient descent. The gradient descent is ‘stochastic’ because gradients are calculated from a small sample of points in the support of the function. Each sample from the support allows the parameters of the network to adapt incrementally toward a solution. The optimizer takes a small step in the direction of the gradient for the sampled points, where the step size is controlled by a specified ‘learning rate’ parameter. Gradients are calculated for the multi-layer network via backpropagation, a fast numerical implementation of the chain rule.

By replacing a *linear algebra* problem with an *optimization* problem, deep learning subverts the curse of dimensionality that stands in the way of numerical PDE solutions. Han and E (2016) demonstrate the basic approach in a stochastic control setting, while Han, Jentzen and E (2018) show that deep learning can be used to solve a wide class of PDE problems in 100 dimensions or more. To our knowledge, the implementation we present in this chapter is the first application of Han and E’s deep learning-based method in the field of macroeconomics and finance.

We claim that the deep learning method presented here is superior to the DSGE methods used to obtain numerical solutions in macroeconomic models. It is quite general, allowing components of the CIR85a construct–utility functions, state dynamics, production specifications—to be swapped in and out at will. It characterizes equilibrium without requiring any knowledge about a steady-state solution, and permits the study of equilibrium through impulse responses. The model need not be linearized or localized. Perhaps most importantly, it does not require the economy to be in a steady state at any point of the analysis. Unchaining the solution from the steady state relaxes the notion of intertemporal equilibrium in an attractive way, while eliminating the numerical errors introduced by perturbation methods. While non-steady state solutions may also be computed for finite-horizon problems via dynamic programming methods, the deep learning approach presented here is both easier and more scalable. Backward recursions and discretizations of the state space are not necessary. One simulates the dynamics of the economy forward to create the data used by the neural network to learn the

optimal control laws.

Cochrane (2014: 5) derides the dynamic portfolio theory enterprise, writing “calculating partial derivatives of unknown value functions is hard and, more importantly, nebulous. People sensibly distrust model-dependent black boxes.” Though Cochrane has portfolio management applications in mind rather than general equilibrium analysis, I respectfully disagree. Deep learning makes it easy to differentiate the value function, and our implementation opens the ‘black box,’ making the specification of the economy, preferences, constraints and controls completely modular. We also show how optimal controls discovered by the neural networks can be profiled to facilitate interpretation. The advances in numerical methods presented here support a fresh look at the CIR85a portfolio-theoretic formulation of the consumption-investment problem as an equilibrium construct competitive with the real business cycle (RBC) and New Keynesian paradigms.

Our task is to compute an equilibrium of the CIR85a model. As noted in Chapter 2, CIR85a were able to derive many properties of a general equilibrium by working with the first-order conditions of an implicit solution. In this chapter we discover further properties of the CIR85a equilibrium by computing numerical solutions for $a(W, Y, t)$, the optimal dynamic allocation of capital, and $C(W, Y, t)$, the optimal consumption path. Similarly, we compute the risk-free rate $r(W, Y, t)$ and risk premia ϕ_W and ϕ_Y numerically to summarize the model’s implications for asset pricing. Given these quantities, we compute expected rates of return for equities and contingent claims.

To obtain a numerical solution we must specify the CIR85a economy more concretely than their abstract formulation of $N + K$ Ito processes for production and the state of the economy, respectively. CIR85b provides a simple specification of the economy with log utility and production returns defined by a constant scaling of the state vector. Under these restrictions, CIR85b derive equilibrium bond prices. Studies of the risk-free term structure employing the bond pricing model of CIR85b thus provide inference on the state processes of a CIR85a economy, subject to the restrictions of CIR85b. We use the estimates of Chen and Scott (2003) to parameterize the state processes with $K = 3$ and set production yields to be a constant multiple of the

state, while relaxing the assumption on utility to a power utility form. With our chosen parameters, we obtain reasonable values for macroeconomic and financial summary statistics.

The contributions of this chapter are substantial:

First, I introduce a novel numerical method to the fields of macroeconomics and finance, making it ready for applications and employing it to solve a problem of fundamental theoretical and practical importance.

Second, I provide a roadmap for specifying and calibrating a CIR85a economy. I develop a more extensive CIR85b economy with multiple production processes and state variables, and calibrate it using well-grounded empirical estimates. Though simple, the exemplar presented here provides a useful template for the specification and calibration of more elaborate models.

Third, I compute equilibrium numerically outside of a steady state and study its properties using impulse responses and numerical profiles of the optimal controls. Impulse responses are a preferred method for analyzing the properties of dynamic economic models and integrating them into our solution framework allows for continuity and dialogue with the literature. Control profiles illustrate the dependence of equilibrium policy responses on the economic state variables, making the function approximation found by the neural network less of a black box.

Fourth, the solution I compute permits exploration of the CIR85a equilibrium at a level that has so far been mathematically inaccessible. The original CIR85a model left equilibrium investment allocations and consumption as implicit features of equilibrium. Here they are brought out explicitly, showing how responses on the production side of the economy support an optimal consumption profile and generate asset price movements.

Fifth, I compare the results of the solved model to those generated by the incumbent RBC/DSGE paradigm. My alternative solution procedure offers several advantages. Most importantly, this chapter demonstrates the usefulness of our numerical solution method as a means of unburdening macroeconomic analysis of a particularly unfortunate implicit assumption about the economy—namely the assumption that the economy remains in the neighborhood of a steady state and, if disturbed, will return to the steady state

asymptotically. Our numerical tools provide another reason to banish the dogma of the economy's dynamic stability for good.

3.1 Computing Equilibrium in CIR Models

We seek to compute an equilibrium $(a, C, r, \phi_W, \phi_Y)$, where $a(W, Y, t)$ and $C(W, Y, t)$ are optimal control policies for investment and consumption, r is the equilibrium risk-free interest rate, and (ϕ_W, ϕ_Y) are the risk premia that determine the equilibrium rates of return for contingent claims given a contingent claims basis, as discussed in Chapter 2. For a finite-horizon problem, the optimal control policies will be time-dependent, in general, leading to non-trivial time paths for r , ϕ_W and ϕ_Y . Further, if we have more than $N = 1$ productive processes and $K = 1$ state variables, the optimal control policies will be high-dimensional from a computational perspective.

In cases where one cannot guess and verify a solution to the stochastic control problem, solutions must be found numerically. Standard numerical methods for the solution of partial differential equations (e.g., LeVeque 2007) may be used when the number of state variables is relatively small. And in infinite-horizon problems with time-invariant policies, approximate dynamic programming (Bertsekas 2017) may be brought to bear to obtain the value function and/or the policy function, which we have been calling the indirect utility function and controls, respectively.¹ But existing numerical methods tend to require a great deal of problem-specific customization, and for very high dimensional problems such methods quickly become intractable due to the familiar curse of dimensionality.

Our overall strategy for computing equilibrium is based on Han and E (2016), which combines stochastic simulation with deep learning to create a powerful numerical tool that can be used effectively in very high dimensional settings.² In addition to its applications in pattern recognition, predictive

¹Incidentally these are also the methods taught to economists by Ljungqvist and Sargent (2018).

²Han, Jentzen and E (2018) extends the method and provides code which has become the starting point for our own. However their reformulation of the control problem in terms of backward stochastic differential equations introduces additional complications

modeling, and natural language processing, deep learning is beginning to be appreciated within the field of applied mathematics as a technology for general-purpose function approximation. Deep learning is seeing an expanding range of application in scientific computing, especially in the numerical analysis and solution of partial differential equations.³

Deep learning algorithms generally require prodigious volumes of data to form a useful inference about a function. Here, the ‘data’ for training the neural network are obtained by simulating a discretized version of the state evolution process until the training process appears to achieve an optimum. Inference proceeds surprisingly quickly. Whereas applications of deep learning in data science learn a function mapping from a set of covariates to a target variable directly from the data with no other guidance, the structure imposed by the specification of the economy allows the neural network to converge to an optimum with relatively few simulations of the underlying processes. Achievement of an optimum is detected by monitoring the cumulative utility realized over the chosen time horizon. A lack of improvement in utility for a long sequence of training runs suggests that an optimum has been located.⁴

Time-dependence in the control functions is captured by approximating the control functions at each of the discretized time steps. The ability to solve for optimal policies at multiple time steps makes the imposition of a steady state unnecessary. On the other hand, the finite-horizon formulation of the CIR economy resembles an optimal resource depletion problem (e.g., Pindyck 1978). In the absence of a utility for final wealth, an optimal path for wealth will always terminate at zero. Mindful of this mathematical artifact,

which would take us too far afield. It is sufficient to think of the problem as a matter of learning the optimal response from repeated simulations of the system.

³See, for example, Sirignano and Spiliopoulos (2018), Al-Arifi et al (2018), Beck, E, and Jentzen (2019), Raissi, Perdikaris, and Karniadakis (2019), Hure (2019), and Wu and Xiu (2020).

⁴One can do no better than this criterion when solving a nonlinear optimization problem numerically. The optimum found may be local rather than global. We mitigate the risk of finding a local solution that is not globally optimal by using a convex objective, choosing initial values for the neural network parameters that imply an interior solution, and using a conservative learning rate to avoid overshooting and oscillatory behavior. Each of these choices is consistent with best practices in deep learning.

we solve the model for a long time horizon and analyze its consequences for a shorter initial interval.⁵

The method of Han and E (2016) has some resemblance to the techniques known as neuro-dynamic programming or deep reinforcement learning, which also employ the nonlinear function approximation powers of neural networks to solve stochastic control problems.⁶ The approach of Han and E (2016) is simpler, however. It allows the problem to be solved forward over the entire horizon rather than solving backward, step by step, using dynamic programming recursions, as in neuro-dynamic programming. In addition it makes direct handling of the value function unnecessary. State and control constraints are also far easier to specify and impose.⁷

3.1.1 Overview of the computational strategy

We are now ready to present an overview of the computational strategy. Before training begins, several features of the model are specified:

1. Choose a time horizon T and discretize it into increments Δt .⁸
2. Define the functions $\mu(Y)$ and $S(Y)$ for the state processes (2.3) and functions $\alpha(Y)$ and $G(Y)$ for the production processes (2.1).
3. Set initial values for the state variables W_0 , η_0 , and Y_0 .⁹

⁵One can tie the network weights across time steps to find a time-invariant policy, if desired. In this form our problem becomes a mean-field control problem (Bensoussan, Frehse and Yam 2013). Deep reinforcement learning has also been employed for infinite horizon problems in which time-independent controls are sought.

⁶See Bertsekas (2017) on neurodynamic programming and Sutton and Barto (2018) on deep reinforcement learning.

⁷The constraints imposed in the models studied here are enforced by choices of activation functions in the output layer of the neural networks, as we explain below. Other constraints may be handled by adding terms to the objective function that penalize violations of each constraint. The analytical approaches to constrained optimization of Zariphopoulou (1994) and Sethi (1997), for example, will promote appreciation for straightforward numerical methods.

⁸As with all methods using time-discretization, one must take care that the step size is not so large that convergence is lost, nor so small that computation is excessively expensive.

⁹Strictly speaking the states W and Y are sufficient to describe the system, as indicated by the arguments to the equilibrium variables. However for reasons explained below, we believe it is convenient to include the output vector $d\eta$ in the states.

4. Initialize the parameters defining the deep neural network.

Following standard practice in deep learning, we repeat our training procedure over a large number of ‘epochs’ until no further improvement of the objective function is evident. A preview of the procedure for a single epoch is as follows:

1. Simulate paths of the exogenous variables for the time horizon T .
 - (a) Draw M samples of the Brownian motion dZ_t .
 - (b) Simulate M realizations of Y_t from Y_{t-1} using the discretized SDE for the state evolution.
2. Learn the optimal controls for each time step from 1 to T :
 - (a) Using a deep feedforward neural network for each control, learn a mapping from the states to the controls $a(W, Y, t)$ and $C(W, Y, t)$, where M realizations of W_{t-1} , η_{t-1} , and Y_t play the role of training data.¹⁰
 - (b) Determine η_t and W_t using a and C .
3. Evaluate the solution found and perform feasibility checks. Compute $u(C_t)$ for the entire time horizon. Check for any violations of constraints on a , C , or W . If they are violated, compute appropriate penalty terms. Form the loss function from the utility function and the constraint penalties.
4. Update the parameters. Compute the gradient of the value function with respect to the unknown parameters defining the neural networks. Take a step in the direction of the gradient according to a chosen learning rate.

Running this optimization to completion will yield an approximation to the optimal control policies, i.e., functions $\hat{a}^*(W, Y, t_i)$ and $\hat{C}^*(W, Y, t_i)$ for $i = 1, \dots, T$ that obey all stated feasibility conditions and maximize the given utility function.¹¹

¹⁰There is no need to learn b because it is equal to zero as a condition of equilibrium.

¹¹We know $b^*(W, Y, t) = 0$ a priori.

Recalling the definition for $J^*(W, Y, t)$ given in (2.10), the indirect utility functions may be approximated using

$$E_{W,Y,t} \int_t^T U(a(s), C(s), Y(s), s) ds \approx \frac{1}{M} \sum_{j=1}^M \sum_{i=t}^T U_{ij}(a_i(t_i), C_i(t_i), W_j(t_i), Y_j(t_i)) \quad (3.1)$$

That is, for a given utility function, \hat{J}^* may be found by averaging total utility over many realizations of the state variables W and Y .

The calculations sketched above yield optimal controls and a numerical solution for the dynamic equilibrium of the system. Additional work is then needed to compute the equilibrium risk-free rate r and the risk premia ϕ_W and ϕ_Y . The remaining features of the equilibrium can then be found by numerically differentiating \hat{J}^* in a neighborhood of the optimum and performing calculations using the state and control processes:

1. State-contingent, time-dependent market prices of risk $-J_{WW}/J_W$ and $-J_{WY}/J_W$ may be calculated by numerically differentiating \hat{J}^* at each time step.
2. Ex ante variances and covariances may be computed from the system specification or, if necessary, estimated by sampling W and Y at the initial time step. Risk premia ϕ_W and ϕ_Y are the product of these risk statistics and the market prices of risk, as given in (2.23).
3. The equilibrium interest rate r can then be calculated using the closed-form expression (2.17), or numerically, using (2.18).

At this point computation of the equilibrium is complete.

We may then wish to undertake the further computations to enhance our understanding of the equilibrium. We consider impulse responses, profiles of the optimal controls learned by the neural network, and contingent claim prices consistent with equilibrium.

Dynamic economic models are often analyzed in terms of their responses to a shock, which is a straightforward exercise in a DSGE model or an esti-

mated vector autoregression. One can compute the equivalent of an impulse response in the present model by comparing the trajectory of the system with $dZ_t = 0$ for all t to the trajectory of the system for a pre-specified vector process dZ_t . For example, one could set $dZ_{0,i} = -1$ for some variable i and all other values to zero to examine the impact of a one-standard deviation adverse shock to the chosen variable. The result is not an expansion around a steady state but a comparison of two dynamically optimal responses.

The control function mappings learned by the deep neural network are not readily interpretable in terms of the network parameters. Nevertheless, we might want to understand the control response to different values of the states to evaluate the economic plausibility of the solutions learned by the neural networks. The controls are mappings from a multivariate state space to scalars. To visualize the mappings, we can compute graphical profiles of the controls, successively isolating each variable in the state space. Specifically, we allow one of the state variables to vary while fixing the other variables at their mean value. This yields a series of two-dimensional plots showing how each control responds to variation in each state variable, *ceteris paribus*.

Finally, we may wish to calculate the values of some fundamental contingent claims at each time step within our horizon. If some risk factors are identified as prices, interest rates, and volatilities, one can calculate equilibrium values of bond yields, forwards, futures and/or options. These calculations would show, in a fairly precise way, what sort of financial market developments accompany equilibrium in the real economy. In our calibration we compute long-run equilibrium yield curves as a check on our state process specification and as an aid in scaling output yields to empirically reasonable levels..

This completes our overview of the computational procedure. We now turn to the details of the stochastic simulation and deep learning procedures we glossed over previously.

3.1.2 Discretizing the state evolution

Fix a time step $\Delta t \approx dt$ that breaks the time horizon T into sub-intervals indexed by m . Using this time step, we wish to discretize the stochastic differential equations that define the state evolutions.

For a vector-valued stochastic differential equation of the form

$$dX = \mu(X, t)dt + \sigma(X, t)dZ$$

the Euler-Marayama discretization for the k -th element of X is given by (Kloeden and Platen 1999: 340-341)

$$X_{m+1}^k = X_m^k + \mu^k(X_m, m)\Delta t + \sum_j \sigma^{k,j}(X_m, m)\Delta Z^j$$

where ΔZ^j is the j -th component of a draw from the $\mathcal{N}(0, \Delta t)$ distribution¹² and σ^k is a vector of the k -th variable's exposures to all sources of uncertainty. Thus for a vector of means μ and a covariance matrix σ we have

$$X_{m+1} = X_m + \mu(X_m, m)\Delta t + \sigma(X_m, m)\Delta Z$$

The Euler-Marayama discretization provides strong convergence of order 0.5 and weak convergence of order 1. Strong convergence of order 0.5 means that the error of pathwise statistics will decline with the square root of M . Weak convergence of order 1 means that the error of averages will decline with M . For our purposes, we are interested in the order of weak convergence, since this will determine the Monte Carlo error in the computation of J^* . Other discretization schemes like the Milstein scheme and Runge-Kutta-like methods exist that achieve higher orders of strong and weak convergence (Higham 2001, Kloeden and Platen 1999). We leave an analysis of the benefits afforded by these alternatives to future research.

When computing the evolution of wealth we need not include terms for the risk-free rate and contingent claims. Recalling (2.7) and imposing the

¹²To be completely clear, the standard deviation is $\sqrt{\Delta t}$.

equilibrium conditions, β and r drop out of the wealth evolution equation:

$$dW = W \left[\sum_{i=1}^N a_i \alpha_i - \frac{C}{W} \right] dt + W \sum_{i=1}^N a_i \left(\sum_{j=1}^{N+K} g_{ij} dZ_j \right) \quad (3.2)$$

Indeed, the characterization of interest rates and contingent claim returns in CIR85a stems almost entirely from the authors' clever discovery of two ways to add zero to the wealth evolution equation. Under the definition of equilibrium, neither term affects the state evolution. Accordingly neither β nor r can be determined directly from solution of the control problem, as we explained previously.

3.1.3 Learning optimal controls

Deep learning refers to a set of methods developed for learning functions from data using neural networks (Goodfellow, Bengio, and Courville 2016). For an input vector $x \in \mathbb{R}^N$, a neural network computes $y = g(Ax + b)$, where A and b are an $M \times N$ matrix and an M -vector, respectively, and g is a nonlinear activation function such as the logistic sigmoid function or the rectified linear unit (ReLU). Deep neural networks compose several layers of such calculations, with the output of one layer becoming the input vector for the next layer. Writing

$$f(x; \theta) = g_N (A_N (\dots g_2 (A_2 g_1 (A_1 x + b_1) + b_2) \dots) + b_N) \quad (3.3)$$

for activation functions g_1, \dots, g_N exhibits the function approximation strategy of deep neural networks in a formal way. If we write the original input vector as x_0 and intermediate outputs as $x_i = g_i(A_i x_{i-1} + b_i)$, the calculations producing x_1 through x_{N-1} occur in the 'hidden' layers of the deep neural network. The calculation yielding x_N is done by the 'output' layer.

Neural networks have long been recognized as universal approximators for functions. Cybenko (1989) showed that a neural network can approximate any Borel-measurable function, while Hornik et al (1989) showed that deep feedforward neural networks have the same approximation property. The

latter result is of interest because deep neural networks are more readily optimized and admit more parsimonious representations than single-layer networks (Liang and Srikant 2017).

More recently mathematicians have begun to formalize the convergence of deep neural network approximations using tools in functional analysis. Rigorous definitions of the spaces of functions represented by specific neural network architectures and the norms in which those representations converge to target functions have enabled some theoretical convergence bounds to be established. For example, E, Ma and Wu (2019) define a compositional function space corresponding to the representation (3.3) and establish a theoretical rate at which the deep neural network converges to a target function. In addition, they derive an upper bound on the Rademacher complexity of the function space, suggesting that target functions may be efficiently estimated.

The theoretical results on convergence of deep neural networks for wide classes of practically-useful functions have been borne out by simulation studies in which the target function is known analytically. Han and E (2016) solve two stochastic optimal control problems using deep neural networks. They show that the control function approximation learned in the neural network solution to the optimal portfolio liquidation problem of Bertsimas and Lo (1998) converges quickly to its theoretical value, while the optimal liquidation cost simultaneously converges to its known value. Their neural network solution to the energy storage-allocation problem of Jiang and Powell (2015) improves on the solution found by approximate dynamic programming methods. Favorable results have also been demonstrated for very large-scale problems. Han, Jentzen and E (2018) solve a linear-quadratic-Gaussian stochastic control problem in 100 dimensions with a relative approximation error of only 0.17 percent.¹³

These simulation studies and theoretical results give us good reason to be confident that deep neural networks can efficiently approximate optimal controls in stochastic settings of very high dimension, and we can therefore expect that $a(W, \eta, Y, t)$ and $C(W, \eta, Y, t)$ can be well-represented by a deep neural network.

¹³As noted previously, we use their code as the starting point for our own.

In solving for the optimal controls we are seeking mappings

$$a(W, d\eta, Y, t) : \mathbb{R}^{N+K+1} \rightarrow \mathbb{R}^N$$

and

$$C(W, d\eta, Y, t) : \mathbb{R}^{N+K+1} \rightarrow \mathbb{R}$$

for each $t = 1, \dots, T$ that satisfy all feasibility constraints and maximize utility. Here we see that the advantage of including $d\eta$ in the state vector is to allow the dimension of the input vector for a to be larger than the dimension of the output vector. The extra input variables may ‘over-identify’ the mapping a , ensuring that the neural network matrices are full rank and introducing more information that can be useful in computing the optimal mapping.

I treat a and C as two separate mappings to N - and 1-dimensional outputs, respectively. The deep neural network is implemented with 2 hidden layers of dimension 60, with rectified linear unit (ReLU) activation functions in between the hidden layers.¹⁴ We implement batch normalization between network layers because it has been shown to improve the efficiency of training. (Ioffe and Szegedy 2015)

The ability to impose economically-motivated constraints on candidate optimal controls with well-chosen activation functions is a useful feature of our deep learning strategy. The ‘softmax’ activation function is used to impose the constraints on a . For an $N \times 1$ vector argument x ,

$$\text{softmax}(x)_i = \frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}}$$

By construction, the vector of outputs from the last hidden layer will be mapped to the interval $[0, 1]$, and the results will sum to one, satisfying the non-negativity and unit-sum constraints assumed of the investment allocations. For the consumption subnetwork I model the consumption-wealth ratio, C/W , rather than the level of consumption. Whereas consumption

¹⁴The ReLU activation is given by $g_i(A_i x_{i-1} + b_i) = \max[0, A_i x_{i-1} + b_i]$.

can take any positive value, the consumption-wealth ratio must fall in the interval $[0, 1]$. To impose this constraint on the values of C/W we apply a sigmoid (logit) activation function to the subnetwork associated with the control. The level of consumption is immediately recovered after multiplying by wealth. Thus at each time step t_i the neural network furnishes approximate optimal policies $\hat{a}(W, Y, t_i)$ and $\hat{C}(W, Y, t_i)$ that satisfy our non-negativity and summation constraints.

We now show how the neural network is embedded in the discretized problem for training. The following comprises the training steps for one batch of N simulated time series of length T . Let the initial levels of wealth W_0 , production η_0 , and the economic state variables Y_0 be given, and set $t = 1$:

1. Draw $\Delta Z \approx dZ$ from a $\mathcal{N}(0, \sqrt{\Delta t})$ distribution, where Δt is the time step, and compute Y_{t+1} using the discretized SDE:

$$Y_{i,t+1} = Y_{i,t} + \mu_i(Y_t, t)\Delta t + s_i(Y_t, t)\Delta Z$$

Because they are exogenous quantities, both ΔZ and Y may be calculated for all time steps at the outset.

2. Then for each time step $t = 1, \dots, T$:

- (a) Compute the return on investment for each process $d\eta_{i,t+1}$:

$$d\eta_{i,t+1} = \alpha_i(Y_{t+1}, y)\eta_{i,t}\Delta t + \eta_{i,t}g_i(Y_{t+1}, t)\Delta Z$$

- (b) Compute optimal policies \hat{a}_t and \hat{C}_t :

$$\hat{a}_t = N_a(Y_{t+1}, d\eta_{t+1}, W_t | \Theta_a)$$

$$\hat{C}_t = N_C(Y_{t+1}, d\eta_{t+1}, W_t | \Theta_C)$$

where N_j denotes the function approximation achieved by the neural network and $\Theta_j = \{A_{1j}, \dots, A_{lj}, b_{1j}, \dots, b_{lj}\}$ denotes the parameters of a neural network with l layers, per (3.3).

- (c) Capture the influence of the controls on future values of the state variables. Compute W_{t+1} from (3.2):

$$W_{t+1} = W_t + \sum_{i=1}^N d\eta_{i,t+1} - \hat{C}_t$$

and set η_{t+1} based on the investment decision:

$$\eta_{t+1} = \hat{a}_t W_t$$

3. Compute the loss function \mathcal{L} and the gradient \mathcal{L}_Θ , where $\Theta = \{\Theta_a, \Theta_C\}$.

We elaborate on the final step below.

3.1.4 Loss function and optimization

Training the neural network refers to the task of finding the parameter set Θ that minimizes a specified loss function. In other words, training a neural network involves solving a numerical optimization problem. An optimal parameter set Θ^* defines the approximately optimal mappings $\hat{a}(W, Y, t)$ and $\hat{C}(W, Y, t)$.

Deep neural networks find

$$\Theta^* = \arg \min_{\Theta} \mathcal{L}[y, f(x, \Theta)] \quad (3.4)$$

using an optimization algorithm, where x are the inputs, y are the targets, $f(x, \Theta)$ are the outputs and \mathcal{L} is a specified loss function. For example, a one-layer network with a linear activation function and loss function $\mathcal{L} = \frac{1}{N} \sum_i (y_i - f(x_i, \Theta))^2$ is equivalent to standard linear regression by the method of least squares.

Optimization algorithms for deep neural networks are based on the notion of stochastic gradient descent.¹⁵ The ‘stochastic’ part refers to using a

¹⁵See Bengio, Courville, and Goodfellow (2016) for a rigorous treatment.

random subset of the training data to evaluate the loss function.¹⁶ Gradients are then calculated for each of the parameters using backpropagation, an efficient algorithm for computing derivatives across all layers of the neural network using a numerical chain rule. As in deterministic gradient descent algorithms, parameters are updated in the direction of descent. The optimizer takes partial steps in the direction of the stochastic gradient controlled by a parameter known as the learning rate. Like Han, Jentzen and E (2018) I use the Adam optimization method implemented in TensorFlow, which allows the learning rate to adapt dynamically.

Given an optimization procedure, training is driven by the specification of the loss function, which the optimizer attempts to minimize. Since we are maximizing the indirect utility function, the negative of the indirect utility function appears in the loss function. We use a power utility specification

$$U(C_t) = \frac{(C_t + 1)^\gamma}{\gamma} \quad (3.5)$$

to ensure that outcomes in our model are driven by the production side of the economy. However nothing prevents us from using a more complicated specification that is not separable in time, such as a utility function with external habit formation. The indirect utility function is evaluated using the history $\hat{C}(W, Y, t)$, $t = 1, \dots, T$. For time-additive utility functions, utility could be calculated at each time step. Utility functions that are non-separable in time will take the entire history of consumption as input. In Chapter 4 we expand the utility function to include labor supply.

We can also include penalty terms in the loss function to enforce constraints on states and controls. Penalty terms should add large increments to the loss function when constraints are violated. Because we have already implemented constraints on the investment allocations and consumption via the activation functions in the neural network's output layer, explicit constraints are not needed for the present implementation.¹⁷

¹⁶One will get slightly different results from run to run for different draws of the data used in the stochastic simulation. I set a random seed for reproducibility.

¹⁷For appropriate initial values W_0 and reasonable parameter values the odds of wealth going to zero should be small in any event.

We note the importance of keeping utility functions and penalty functions separate in the implementation. Market prices of risk are derivatives of the utility function, while derivatives of the loss function are needed for optimization.

3.2 Specification and Calibration

In this section we develop a specification for a CIR economy, giving concrete forms to the abstract state and production processes (2.1) and (2.3). We also describe the parameterization of the model with a combination of estimates from the empirical literature and calibration to observable outcomes.

Economists typically parameterize their models in two subgroups (Dawkins, Srinivasan and Whalley 2001). It is not uncommon for the majority of parameters in the first subgroup to be set with reference to other literature. The parameter settings of a well-studied model may be taken over directly. In other cases one refers to empirical studies that estimate the parameters of interest. The estimates employed must have a reasonable chance of being invariant to the form of the model or the prevailing policy environment, which is to say, ‘structural.’ More practically, estimates are ideally drawn from an econometric model capturing most of the variables under study in the numerical model. When estimates are taken from microeconomic studies further reconciliations may be needed to convert the estimates into compatible macroeconomic parameters (Browning, Hansen and Heckman 1999).

The second subgroup of parameters that are not set on the basis of literature are then ‘calibrated.’ A strong interpretation of calibration chooses the remaining parameters so that the moments of time series simulated from the model match moments of empirical aggregate time series as closely as possible. Under this interpretation calibration closely resembles *estimation* by the simulated method of moments (see, e.g., Adda and Cooper 2003) or indirect inference more generally.

A weaker sense in which researchers calibrate their models is to choose the remaining subgroup of parameters so the predictions of their models match the levels of certain target variables, similar to the sense in which

one calibrates a Celsius thermometer to read 0 in freezing water and 100 in boiling water (Dawkins, Srinivasan and Whalley 2001). In this case one views calibration as an exercise in setting the model at a reasonable starting point from which the quantitative consequences of shocks in policy changes can be analyzed in an approximate sense. The model is not taken to have generated any observable data series. Instead, the model is set up in a way that its qualitative predictions may be analyzed on a relative scale.

In this study we adopt the weaker sense of calibration for two reasons. First, embedding the deep learning problem in a structural estimation of model parameters would significantly raise the level of computational complexity while raising questions about estimators, convergence and inference that would take us too far afield of our proof-of-concept exercise.¹⁸ Second, we are not sure which data ought to be used to estimate the model parameters. Our model does not imply steady-state distributions for its objects whose moments might be matched to the long-run moments of aggregate time series. Instead, the model is forward-looking, describing a dynamic plan that is time-consistent but subject to change as uncertainty about the state of the economy is resolved. Large panels of contingent claims prices may provide a suitable empirical basis for estimating our model, but the restrictions entailed for the specification of the state vector are non-trivial (Dai and Singleton 2000, Casassus and Collin-Dufresne 2005). Like the issues raised by estimation, we view the choice of a suitable empirical basis and the formulation of compatible state process specifications as hairy problems in need of separate study, and thus outside of our current scope.

In the present chapter we rely on the results of an econometric study of the term structure of risk-free interest rates based on a multifactor implementation of the CIR85b bond pricing model. In the next chapter we will draw estimates from multiple sources in the applied literature for an extended version of the present model. We calibrate the remaining free parameters in order to broadly match levels of output, consumption, savings, and interest rates consistent with recent experience.

¹⁸For some examples of the technical questions involved, see Chernozhukov et al. 2018 and Kaji, Manresa and Pouliot 2020.

3.2.1 Parameterization of the state process

Though we are not aware of calibrated examples of CIR economies, the literature does contain estimates of CIR state processes derived from studies of the term structure of risk-free interest rates in the United States. CIR85b derives equilibrium risk-free bond prices under the assumptions of logarithmic utility and production that is a scale multiple of the state vector. When econometricians estimate the parameters of the latent state process underlying the CIR85b pricing model, they are obtaining estimates of the state vector for the CIR85a economy under these restrictions. Under identification restrictions and the further assumption that the prices of US Treasury securities fully incorporate risks to production, the state processes that best fit the term structure are the best estimate of the state of the CIR85a economy for a chosen dimension K .

Reliable estimates of a multivariate state vector for the CIR economy are available from the careful study of Chen and Scott (2003) (hereafter, ‘CS’). CS use a Kalman filtering methodology to obtain quasi-maximum likelihood estimates of the latent CIR state vector from weekly and monthly panels of Treasury yields.¹⁹ CS assume the state vector consists of K independent Feller processes of the form

$$dY_i(t) = \kappa_i (\theta_i - Y_i(t)) dt + \sigma_i \sqrt{Y_i(t)} dZ_i \quad i = 1, \dots, K$$

The Feller process exhibits mean reversion in the drift term and a state-dependent scaling of the diffusion term. The speed of reversion to the unconditional mean of the state θ is controlled by the parameter κ . The state Y_i is prevented from taking non-positive values so long as $2\kappa\theta \geq \sigma^2$.

CS produce estimates for $K \in \{1, 2, 3\}$. Estimates of the state vector with $K = 3$ based on weekly data are reproduced in Table 3.1. One can solve for the half-life of deviations from the mean by solving $e^{-\kappa t} = 0.5$ for t . The CS

¹⁹CS use the canonical Shiller-McCullough dataset. Their estimates are quasi-maximum likelihood because the normal distributions used in the Kalman filter are only approximations to the state variable distributions, which are known (per CIR1985b) to be non-central chi-square distributions.

estimates suggest deviations of the state components from their means have half-lives of roughly 0.5, 40, and 20 years, respectively. While comparisons of the weekly estimates to the monthly CS estimates (not reproduced) and failures of the parameter restrictions indicate the parameters for the second and third components are not precisely estimated, it is clear that the state of the economy is driven by one relatively high-frequency process and two low-frequency processes.

We present the state evolution in system form to connect the parameters in Table 3.1 explicitly to the functions $\mu(Y)$ and $S(Y)$:

$$\begin{aligned}
\begin{bmatrix} dY_1 \\ dY_2 \\ dY_3 \end{bmatrix} &= \left(\begin{bmatrix} 0.062539 \\ 0.000043 \\ 0.000113 \end{bmatrix} - \begin{bmatrix} 1.4298 & 0 & 0 \\ 0 & 0.01694 & 0 \\ 0 & 0 & 0.04960 \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \end{bmatrix} \right) dt \\
&+ \left(\begin{bmatrix} 0 & 0 & 0.1604 & 0 & 0 \\ 0 & 0 & 0 & 0.1054 & 0 \\ 0 & 0 & 0 & 0 & 0.0496 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \sqrt{Y_1} & 0 & 0 \\ 0 & 0 & 0 & \sqrt{Y_2} & 0 \\ 0 & 0 & 0 & 0 & \sqrt{Y_3} \end{bmatrix} \right) dZ_t \\
&= (\mathcal{K}\Theta - \mathcal{K}Y_t) dt + S(Y)dZ_t \\
&= \mu(Y)dt + S(Y)dZ_t
\end{aligned} \tag{3.6}$$

In the numerical implementation we define an explicit reflecting boundary for the state vector to ensure positivity because the estimates for Y_2 and Y_3 do not satisfy the required parameter restriction.

The most complete estimate for the state vector with $K = 3$ is the estimate containing the maximum number of nonzero parameters that can be identified from the data. Dai and Singleton (2000) call such models ‘maximal’ models. The CS estimates for $K = 3$ are not maximal because they restrict \mathcal{K} to be diagonal in (3.6). Dai and Singleton (2000) show that the off-diagonal elements of \mathcal{K} may also be identified, in principle, and propose a simulated method of moments (SMM) estimator to estimate the additional

Parameter	$i = 1$	$i = 2$	$i = 3$
κ	1.4298	0.01694	0.03510
θ	0.04374	0.002530	0.003209
σ	0.16049	0.1054	0.04960

Table 3.1: Values for CIR state processes estimated by Chen and Scott (2003)

parameters consistently.²⁰ CS choose not to implement the SMM estimator because it is computationally expensive relative to quasi-maximum likelihood in the Kalman filter/state space model. Duffee and Stanton (2012) subsequently found that the small-sample performance of the SMM estimator is unacceptable even for simple term structure models, while the Kalman filter performs well in quasi-maximum likelihood settings. We are not aware of any multifactor CIR bond price estimates that have been obtained by SMM, or which otherwise estimate a maximal \mathcal{K} . The large likelihoods and small pricing errors achieved by CS suggest that the issue is not practically significant, in any case.

3.2.2 Parameterization of the production processes

Solving for the bond yields implied by the state processes using the multifactor pricing formula of CIR1985b shows that the three components measured by CS correspond to the familiar level ($i = 1$), slope ($i = 3$) and curvature ($i = 2$) components of the yield curve, respectively. One generally takes the level of the yield curve to signal the overall scarcity of capital in the economy, incentivizing agents to divert consumption to capital formation. Movements in the slope of the yield curve are associated with fluctuations in the business cycle, with an inverted yield curve serving as a hallmark of a recession. These interpretations of the state components are useful for designing the exposures of the production processes to the state vector.

When the state variables are at their unconditional means, summing the yield contributions from each of the state components results in a bond yield

²⁰Dai and Singleton (2000) do not furnish estimates of a three factor CIR model.

of roughly 5 percent. We can reproduce this yield by designing a production process with

$$\alpha_i(Y) = \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} Y_t$$

and $g_i(Y) = 0$. In other words, the bond pricing model establishes the scale of the production processes relative to the state process. Clearly any rescaling of the state process can be compensated by an offsetting rescaling of the production processes.

In CIR85b the expected returns on production $\alpha_i(Y, t)$ are scale multiples of the state process. What multiple is appropriate in this context? One is tempted to define $\alpha_i(Y, t)$ to be roughly 10 percent per annum for consistency with returns on equity.²¹ By the same rationale one might choose $g_i(Y, t)$ to match equity market volatilities.

Neither decision would be correct. The function $\alpha_i(Y, t)$ is not only the expected rate of return on capital in CIR85a, but also what determines the expected level of output relative to the capital stock. These statistics are equivalent in CIR85a because capital is the only factor of production, quite unlike production in reality. As a result, we can choose $\alpha_i(Y, t)$ to replicate either expected rates of return on capital or the expected level of output relative to wealth. We opt for the latter because expected rates of return in the capital-only economy of CIR85a should not be expected to match rates of return in the economy of our experience, in which labor makes the dominant contribution to production.

A wealth to GDP ratio of 4 corresponds to an output-wealth ratio of 0.25, significantly larger than the 10 percent rule of thumb for returns on equity. With bond yields at roughly 5 percent with $\alpha_i(Y, t) = 1$, we need a scaling factor of about 5 to set production in reasonable proportion to wealth.

To avoid over-complicating our proof of concept, we set $N = 2$ and define

²¹Let us stipulate for purposes of discussion that the equilibrium level of the risk-free rate is about 2 percent and the equity risk premium is about 8 percent, so the long-run average return on equities is about 10 percent.

the joint dynamics of the production processes as follows:

$$\begin{aligned} \begin{bmatrix} d\eta_1 \\ d\eta_2 \end{bmatrix} &= \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 5 & 5 & 5.1 \\ 5 & 5 & 4.9 \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \end{bmatrix} \right) \begin{bmatrix} \eta_1 & 0 \\ 0 & \eta_2 \end{bmatrix} dt \\ &+ \begin{bmatrix} \eta_1 & 0 \\ 0 & \eta_2 \end{bmatrix} \begin{bmatrix} 0.02 & 0 & 0 & 0 & 0 \\ 0 & 0.005 & 0 & 0 & 0 \end{bmatrix} dZ_t \quad (3.7) \\ &= \alpha(Y)dt + G(Y)dZ_t \end{aligned}$$

If both production processes had the same expected returns and volatilities, any allocation decision would be optimal. The processes have to be different, but not in a way that yields a corner solution to the allocation problem. Thus we define Process 1 to be more exposed to Y_3 than Process 2, increasing its expected return because $Y_3 > 0$. We handicap this superior expected return with a larger volatility of output to maintain a meaningful tradeoff to be resolved by the equilibrium allocation of capital between Process 1 and Process 2. The chosen volatilities keep the output-wealth ratio within reasonable bounds for any feasible allocation. The absence of covariance terms between the yield disturbances treats all shocks to expected production levels that do not originate in the state process as idiosyncratic.

3.2.3 Initialization

Initial values must be chosen for wealth, the allocation of capital, and the state variables. We start the state variables at their unconditional means given by θ and set the initial allocation of capital equally across the two production processes. We will see that this initial allocation is revised substantially. Though the economy will be in equilibrium during each time step, we do not expect the economy to converge to a steady state, and indeed there is nothing in the model or the solution method to force convergence.

Because the utility function exhibits constant relative risk aversion and the production processes have the constant-returns-to-scale property, the

level of wealth does not matter per se for the solution. Wealth simply needs to be set to a large enough level that consumption takes values away from the origin throughout the period. We set $W_0 = 4000$ and $\eta_1 = \eta_2 = 2000$.

As we mentioned previously, CIR bond prices entail a coefficient of relative risk aversion $\gamma \rightarrow 1$ because they are based on logarithmic utility. We set $\gamma = 0.5$. It will be clear that the choice of square-root utility is benign for our implementation. Again, differences between the CIR85a model structure and the structure of other macro models prevent easy comparisons of parameter values to other values used in the literature. Even ‘structural’ parameter estimates are not portable from one model structure to another.²²

3.3 Solution of the Model

Solving the model entails training the weights of a deep neural network to obtain an approximation to the state- and time-dependent functions $a(W, Y, t)$ and $C(W, Y, t)$ that define equilibrium. Training data is obtained by simulating the discretized stochastic differential equations (3.6) and (3.7) that define the model dynamics. We simulate 800 batches of 64 trials each to obtain a total of 51,200 training data paths. Each path extends for an horizon of 5 years divided into 60 time steps (monthly resolution). The number of paths and the length of the time step both appear adequate for the model to converge quickly to a well-behaved solution.²³

Neural networks are defined for each of the control functions at each time step. The neural networks consist of two hidden layers of 60 neurons each and an output layer, with batch normalizations preceding each layer. The output layer for the allocation function has dimension $N = 2$, while the

²²Much of the literature follows Mehra and Prescott (1985) in using the utility function $U(C_t) = C_t^{1-\gamma}/(1-\gamma)$ with values of γ between 0 and 10, which yield negative cardinal utility values for $\gamma > 1$. While the derivatives will have the required signs, we need our utility function to be positive for the optimization to be intelligible.

²³Solution of the model is reasonably fast. I solve the model on a laptop with 12 CPU cores and 32 GB memory. Initializing the computational graph entails some overhead and consumes about a minute of run time. Afterwards each batch run takes roughly one-quarter of one second. Total run time is about 4 minutes 15 seconds. Post-processing to compute equilibrium, control profiles, and impulse responses takes another 1-2 minutes.

output layer for the consumption function has dimension 1. There are no connections between network weights across time steps or controls. In all there are 530,820 trainable parameters available to approximate the optimal control functions.²⁴

In a finite-horizon model with consumable capital it will be optimal for agents to consume all capital by the end of the horizon. The CIR model is comparable to an optimal resource depletion model in this respect. Behaviors and outcomes in the early time steps are thus far more interesting than those toward the end of the horizon, where consumption and wealth both go to zero. I plot outcomes for the first 12 time steps, corresponding to one year in model time.

The indirect utility function is defined in (2.10) as the expectation, conditional on W , Y and t , of the integral of direct utility over the horizon $T - t$. Summing utilities for each simulated path and averaging the sums over a batch of simulations per (3.1) produces a Monte Carlo estimate of the required integral. Convergence of the model to optimal values is indicated by steady increase of the indirect utility function to a plateau. Increasing the number of runs beyond what is required to reach this plateau is not advised, because the model—like any neural network—is prone to overfitting to noise in the simulations beyond that point.

Once the optimum has been found in the training phase, I analyze equilibrium in the model by simulating a new batch of 500 paths. The final simulation provides data for which equilibrium quantities may be computed using the optimal control solutions found during the training phase. Market prices of risk are computed by differentiating the indirect utility function numerically. We use the TensorFlow gradient tape utility to calculate J_W , J_{WW} , and J_{WY} at the optimum.

We have three goals for our numerical analysis of the calibrated model. First, we wish to characterize the elements of the CIR85a equilibrium that remained implicit in the original article. We are interested in (a) the dimen-

²⁴Any econometrician will recoil in horror at the profligate over-parameterization of neural networks. This is the nature of the beast. Our task is not to *identify* each of these parameters, but to determine them enough to obtain a suitable approximation to a policy function—even if the parameter set producing that approximation is not unique.

sionless ratios of output and consumption to wealth; (b) how the path of optimal investment allocations combines heterogeneity, uncertainty, and risk aversion to generate state-dependent output profiles, and (c) the numerical values of r , ϕ_W , and ϕ_Y , as they summarize the asset pricing implications of the model.²⁵ Exposing the optimal control strategies and their dependence on the state of the economy also permits us to assess the economic plausibility of the solution learned by the neural networks.

Second, we want to undertake impulse response analyses to understand how equilibrium performance changes relative to a baseline system state. The impulse response analyses aid in developing causal explanations that lead from changes in the conditions of production to changes in asset prices.

Third and more generally, we want to show that the CIR model and our neural network-based solution procedure furnish a cogent alternative paradigm for macro-financial modeling and deliver an analysis that is competitive with that of the RBC/DSGE paradigm. We believe the best argument for the alternative paradigm is to put its results on display. Our proof of concept will show the time- and state-dependent results for equilibrium and comparative dynamics in impulse responses that are generally inaccessible in numerical models of the aggregate economy.

3.3.1 Equilibrium and control profiles

In this section we present and analyze the equilibrium solutions $a(W, Y, t)$ and $C(W, Y, t)$. Besides characterizing the equilibrium, we would like to convince ourselves that the policy functions learned by the neural networks are economically reasonable. To that end we present results on the levels of output and consumption relative to wealth, and profile the dependence of the policy functions on the state of the economy. Following our discussion of equilibrium behavior, we work out the consequences for asset prices.

We begin by looking at the equilibrium allocation between the two production processes. Because $N = 2$ and $a_2 = 1 - a_1$ it suffices to look at a_1 ,

²⁵Scale-free quantities are reported because the scale of wealth is arbitrary, as indicated above. Output is the sum of $d\eta_i$ for all i . Consumption is determined by the optimal control.

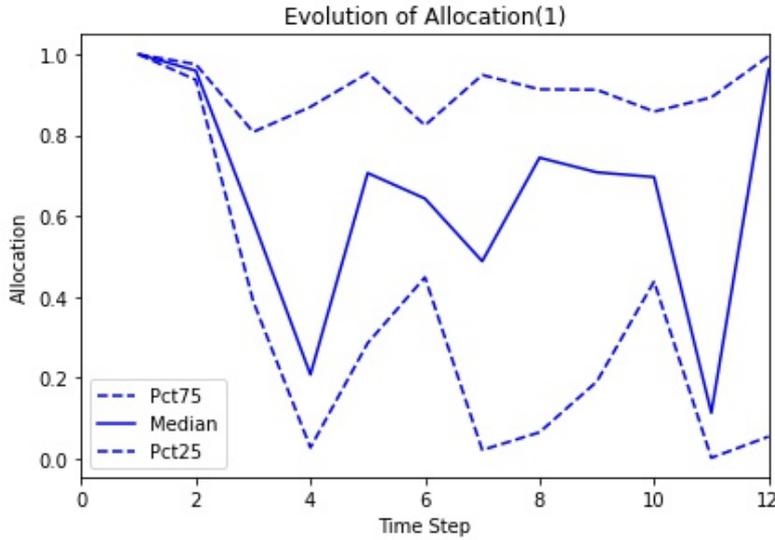


Figure 3.1: Allocation to production process 1 (a_1) in equilibrium.

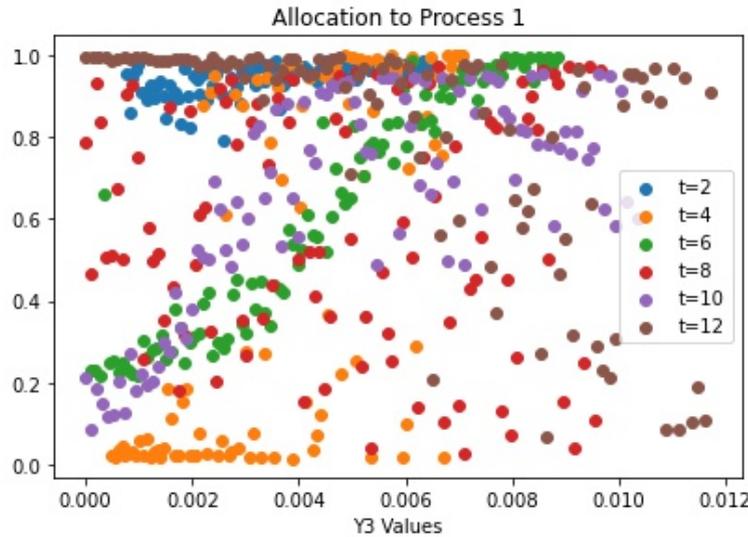


Figure 3.2: Sensitivity of a_1 to Y_3 .

the allocation to Process 1. Figure 3.1 shows that all capital is allocated to Process 1 at $t = 1$ and nearly all at $t = 2$. However from $t = 3$ forward the allocation of capital to Process 1 depends on the configuration of W and Y . The bounds given at the 25th and 75th percentiles of 3.1 suggest that

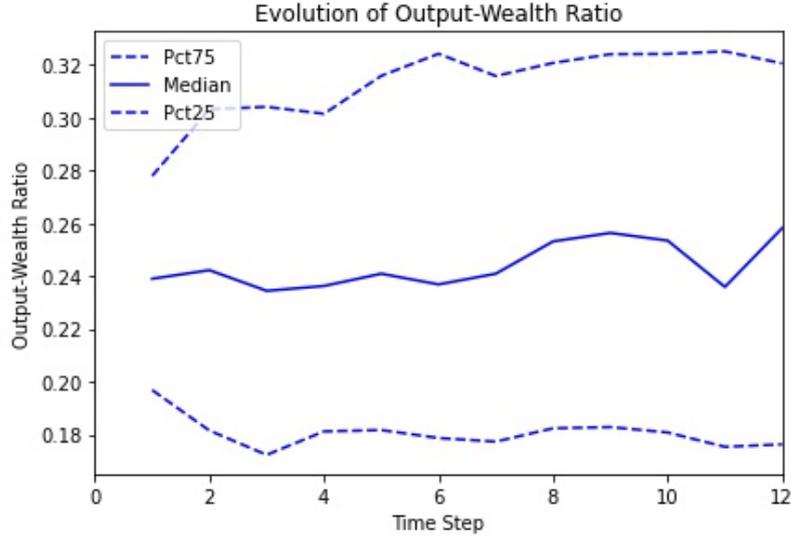


Figure 3.3: Equilibrium output-wealth ratio.

the optimal allocation can vary quite a bit with the state of the economy. Because we know that the production processes differ in their exposure to Y_3 , we expect that the sensitivity of a_1 to Y_3 will shed further light on the state-dependence of the allocation decision.

Figure 3.2 shows how the equilibrium allocation to process 1 reacts to Y_3 , with Y_3 varying on the horizontal axis and the allocation to process 1 on the vertical axis. For purposes of the computations, the values of W , Y_1 and Y_2 are fixed at their mean values from the simulation. The profile of the process 1 allocation is examined at six evenly-spaced time steps from $t = 2$ to $t = 12$, represented by the colors blue, orange, green, red, purple, and brown, respectively. We noted previously that the allocation goes almost completely to process 1 at $t = 2$; hence we see the blue dots clustered near 1.0 for all values of Y_3 . The orange dots at $t = 4$ show that the allocation to process 1 depends on more than the value of Y_3 , as no clear relationship between a_1 and Y_3 is evident. However at $t = 6$ and $t = 10$ a sigmoid relationship between the value of Y_3 and a_1 emerges. At these times Figure 3.1 shows that less uncertainty prevails concerning the optimal allocation.

The dynamic allocation of capital to the two production processes serves

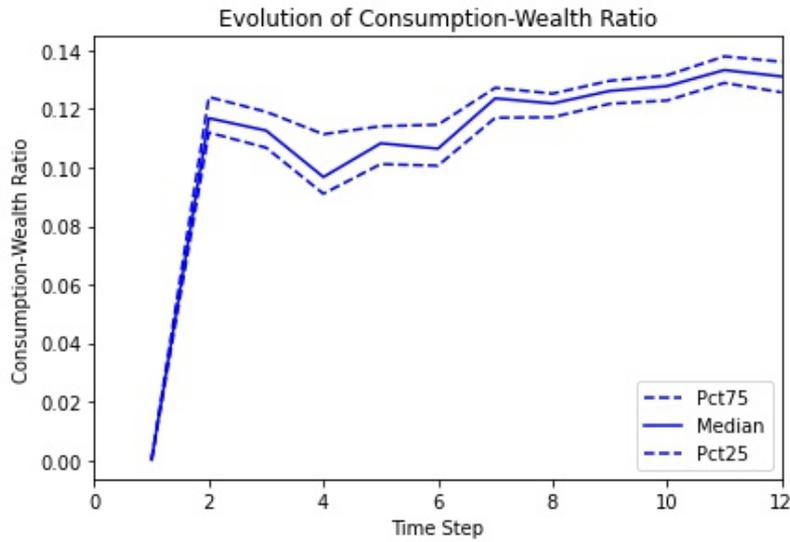


Figure 3.4: Equilibrium consumption-wealth ratio.

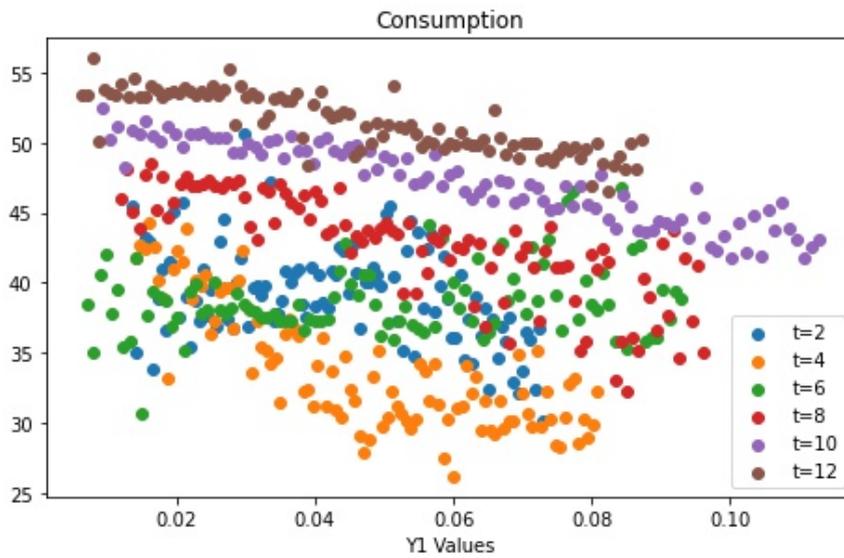


Figure 3.5: Sensitivity of equilibrium consumption to Y_1 .

to stabilize the level of output and lower the risk of consumption. Figure 3.3 shows that the average ratio of output to wealth remains close to the target value of 25 percent we used when calibrating the model, though the actual ratio varies a great deal between the first and third quartiles of the

simulated outcome distribution. Wealth to output ratios between 3 and 5 are consistent with the central tendency of our calibrated model.

The wide dispersion of outcomes for output should be contrasted with the tight distribution of the consumption-wealth ratio in Figure 3.4. Following an initial decision to reinvest all output, the consumption-wealth ratio rises and stabilizes within the 10 to 13 percent range. The difference between the consumption-wealth and output-wealth ratios implies that the reinvestment (savings) rate for the economy can reach one-half to two-thirds of output in the first year of the horizon studied. Such behavior is consistent with capital accumulation in early periods to support production and consumption through the end of the model horizon.

In Figure 3.5 we plot the level of consumption on the vertical axis against the value of Y_1 on the horizontal axis, similar to Figure 3.2. State variable Y_1 is the variable on which the level of output most strongly depends because its unconditional mean is an order of magnitude larger than those of Y_2 and Y_3 . Based on the bond pricing formulas of CIR85b we interpreted Y_1 as the risk factor that determines the level of the yield curve. Figure 3.5 shows the representative agent's willingness to reduce consumption when the level of the yield curve is higher, as one would expect. However we would like to understand how consumption reacts to the conditions of production which underlie the determination of the yield curve and the optimal consumption decision alike.

The dependence of consumption on Y_1 is actually more subtle than it appears. Larger realizations of Y_1 coincide with more current-period output for any choice of allocation. If larger realizations of Y_1 were transitory, the representative agent would do well to consume the output windfall in the current period, because the rise in current output would imply nothing about output in subsequent periods. But because Y_1 follows a Feller process, an elevation of Y_1 above its unconditional mean value will persist; for the parameter values in our calibration, the half-life of the decay is roughly 6 months. The representative agent knows the structure of the economy and takes the forecastability of Y_1 into account when making optimal plans. Thus the representative agent recognizes from the above-average realization of Y_1

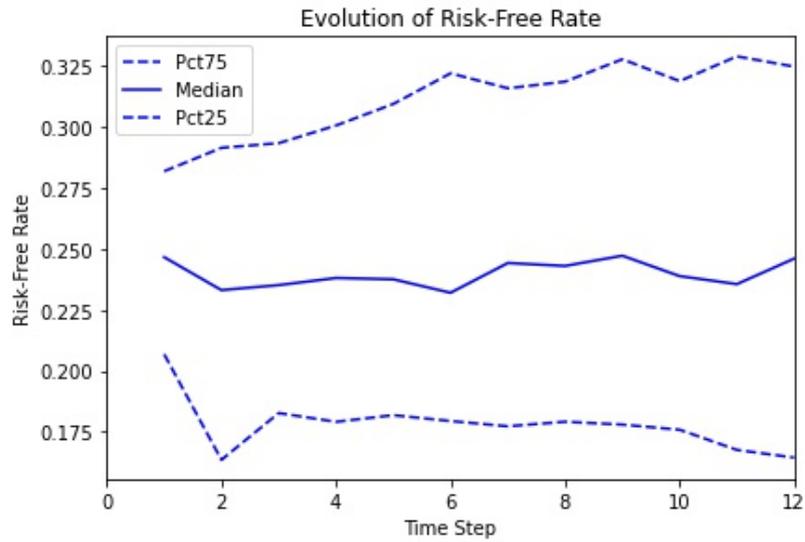
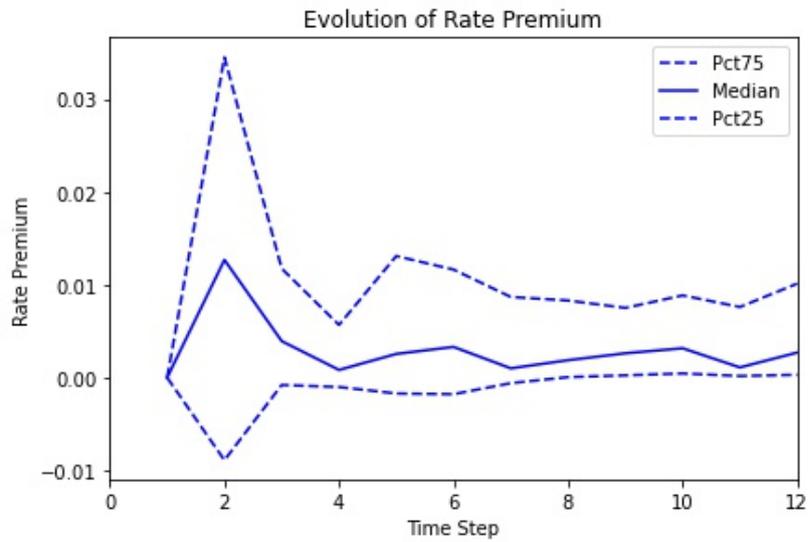
that output will be elevated for several future periods as well. The representative agent knows he will be better off if he reinvests more output in production and consumes less in the current period because higher levels of consumption will be possible in the future. The profiles in Figure 3.5 show that the representative agent reinvests (saves) more for larger realizations of Y_1 . Given that process 1 has a higher expected rate of return than process 2, it is unsurprising that the profile of a_1 versus Y_1 (not shown) reveals that allocations to process 1 depend positively on Y_1 .²⁶

Transitioning from a discussion of consumption sensitivities in terms of the yield curve to a discussion in terms of output illustrates the grounding of CIR85a and CIR85b as a production-based model of financial market equilibrium. Though Y_1 is interpretable in terms of the yield curve, its influence stems from its impact on production. Refining our explanation of the optimal consumption policy to a genuinely causal explanation in terms of production demonstrates that interest rates are not the prime mover in general equilibrium. Higher expected returns on production are reflected in higher interest rates, which serve as an elegant signal to agents that reducing current consumption is optimal.

Having characterized the optimal policies $a(W, Y, t)$ and $C(W, Y, t)$, we now proceed to analyze the consequences of equilibrium behavior for asset prices. The risk-free rate r and risk premiums ϕ_W/W and ϕ_Y that prevail in equilibrium may be computed as functions of a and C . Figure 3.6 shows that the risk-free rate closely tracks the output-wealth ratio. Because we saw in (2.20) that the risk-free rate is the expected return on production minus the equity risk premium, Figure 3.6 suggests that the equity risk premium ϕ_W/W must be quite small for this economy and the calibrated parameter values.

Figure 3.7 confirms our intuition. The equilibrium equity risk premium ϕ_W peaks at a median of 1 percent in period $t = 2$ and remains between 0 and 1 percent in most simulated scenarios thereafter. A very low equity risk premium indicates the volatility of wealth is small, and varies little with the state vector. The conditional volatilities of the state variables in (3.6) that

²⁶Our framework thus reproduces the intuition of Campbell and Viceira (1999).

Figure 3.6: Equilibrium risk-free rate r .Figure 3.7: Equilibrium equity risk premium ϕ_W/W .

define S and the volatilities of output in (3.7) that define G are both small, so the minimal risk premium is unsurprising. If we were to increase them in an effort to calibrate to the equity risk premium, the volatility of output would expand further beyond the already wide range seen in Figure 3.3, and

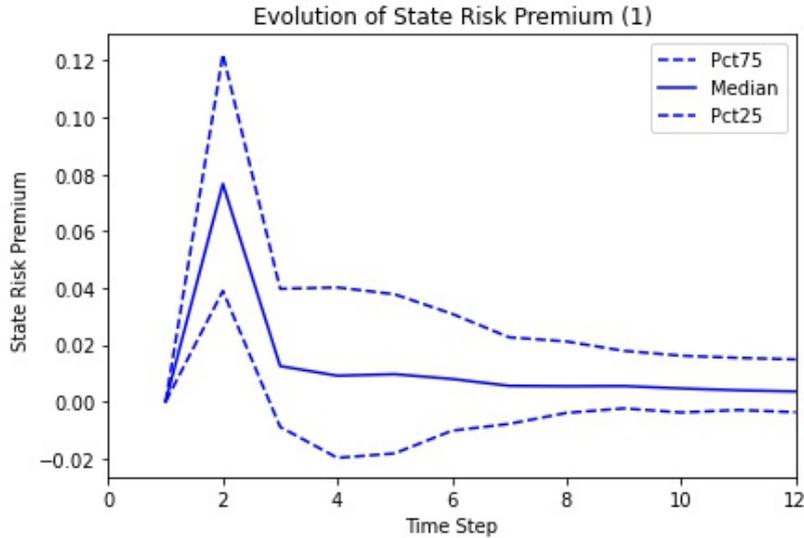


Figure 3.8: Equilibrium contingent-claim risk premium ϕ_{Y_1} .

we would still be left with a risk-free rate in the mid- to high double digits.

The spike in the equity risk premium at $t = 2$ corresponds to the surge of initial investment implied by Figure 3.4. The spike recedes afterwards as investment and consumption return to stable levels. We see the same pattern of a quickly-reversed spike at $t = 2$ in Figure 3.8, which plots ϕ_{Y_1} . Based on (2.26), we reasoned that ϕ_{Y_1} was the equilibrium compensation demanded for taking on an extra unit of exposure to a claim on future consumption that varies one-for-one with Y_1 . The required compensation is small except for $t = 2$. Output already depends heavily on the value of Y_1 , and the value of Y_1 at $t = 2$ will largely determine the levels of output and consumption enjoyed after $t = 2$. Because the representative agent is already so exposed to the value of Y_1 at $t = 2$, additional exposure to Y_1 at this time would further increase the variance of consumption outcomes. That a representative agent would require significant compensation for additional Y_1 risk at this time is unsurprising.

Note that the risk premium attached to Y_1 is 6-7x larger than the equity risk premium. The equity risk premium is so much smaller because it reflects the ability of the representative agent to mitigate his exposure by adjusting



Figure 3.9: Response of C/W to a negative one-standard deviation Y_1 shock.

consumption and changing investment allocations to the less-volatile Process 2. By design, there is no way to mitigate the risk of a contingent claim on Y_1 .

3.3.2 Impulse responses

Impulse response analysis generates additional insight into the dynamics of the model and helps to trace causal pathways from the $N + K$ economic risks in the model to output and asset prices via changes in equilibrium allocations and consumption. In this section we analyze the CIR85a equilibrium by studying impulse responses for Y_1 and Y_3 . We shock Y_1 because it is the state variable to which output is most sensitive. We use a one-standard deviation negative shock, where the standard deviation is given in Table 3.1.

It is obvious that a negative shock to Y_1 at $t = 0$ entails a decline in output relative to the baseline. Optimizing agents may compensate the decline in production by changing production process allocations and their rate of consumption. As we saw in the previous section, changes in allocation help to dampen the effect of changes in output for consumption. However Figure 3.9

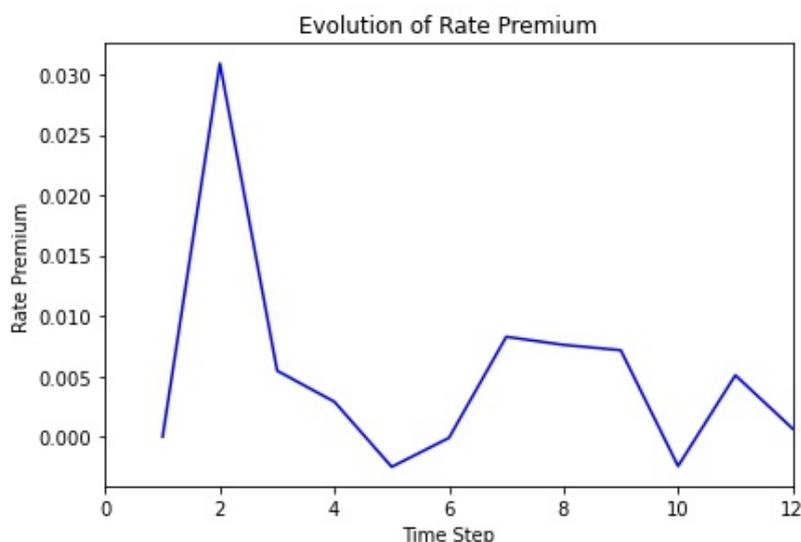


Figure 3.10: Response of ϕ_W to a negative one-standard deviation Y_1 shock.

shows that revisions to allocations cannot completely offset the consumption risk created by a shock to Y_1 . The consumption-wealth ratio fluctuates by -0.08 percent to 0.04 percent relative to the baseline ratio shown in Figure 3.4.²⁷ Is this a big change?

Figure 3.10 shows that the additional volatility in consumption is economically meaningful. The equity risk premium ϕ_W increases by three percentage points at $t = 2$ and nearly one percentage point in $t = 7 - 9$. These are large changes relative to the baseline values shown in Figure 3.7. The shock to Y_1 means that expected output at $t = 0$ is $5 \times 0.16049 \approx 0.8$ percentage points lower relative to wealth than in the baseline scenario. This is not a small difference. A decline of 0.8 percentage points is about 3.2 percent of total output.

A negative 3.2 percent deviation of output from trend usually indicates a recession. In recessions risk premiums increase, causing asset prices to fall until their expected rates of return are sufficient to induce investors to hold them. The increase of three percentage points seen in Figure 3.10 amounts

²⁷Technically the difference is relative to the mean rather than the median shown in Figure 3.4, but the practical difference between the two is small.

to a quadrupling of the baseline risk premium from one to four percent per annum. Because κ_1 is large we can expect the majority of the effect to dissipate after 12 months as shown in the figure.

Our impulse responses do not follow the usual pattern of smooth decay from an initial peak. The second and third rises in the equity risk premium shown in Figure 3.10 are not preceded by symmetric declines that would suggest an oscillatory return to equilibrium.²⁸ In fact these delayed responses correspond to significant increases in the equilibrium allocation to Process 1 in $t = 6 - 9$ and $t = 11$ relative to the baseline scenario (not shown). A negative shock to Y_1 has the deferred effect of motivating agents to take on more risk in the future. After the effect of the shock dissipates, the representative agent increases investment in Process 1, which is riskier, but has a higher expected return.

This simple example shows how our impulse response analysis differs from the typical analysis obtained by expansion around a steady state. Our analysis shows the comparative dynamics of two scenarios where equilibrium behavior is allowed to adjust to the complete dynamic profile of the shock. The response we model is not constrained by assumptions about dynamic stability. Thus the reaction to a shock we model is far richer than what a simple eigenvalue analysis would permit.

3.3.3 Comparison to DSGE methods

We have provided a stochastic simulation of a calibrated benchmark model and analyzed the consequences of an impulse response for consumption, investment allocations, and risk premiums. The model and its numerical solution allow for some useful initial comparisons to be made to the canonical RBC/DSGE approach. What do we gain from employing this new methodology?

DSGE models are typically solved by perturbation methods or stochastic dynamic programming methods. Per Judd (1998: 447), “The basic idea of

²⁸In a dynamically-stable model such behavior is indicative of a complex-conjugate pair of eigenvalues in the steady-state transition matrix.

[perturbation or] asymptotic methods is to formulate a general problem, find a particular case that has a known solution, and then use that particular case and its solution as a starting point for computing approximate solutions to nearby problems. ... Economists have often used special version of perturbation techniques, such as ... linearizing around a steady state.”

The popularity of perturbation methods is attributable in large part to the emergence of Dynare as the numerical software platform of choice for solving DSGE models.²⁹ Dynare has found wide application in the solution of RBC and New Keynesian models alike (Miao 2014, Torres 2015, Costa 2016). The latest version of the Dynare model documentation (November 2020) indicates “the main algorithm for solving stochastic models relies on a Taylor approximation, up to third order, of the expectation functions.” (Adjemian et al, 2020: 51) Upon referring to the technical documentation, we find the solution method is “...essentially a variation on the methods presented by ... Uhlig (1999).” (Villemot 2011: 1)

Uhlig (1999) describes a “general procedure” for solving a DSGE model:

1. “Find the necessary equations characterizing the equilibrium, i.e., constraints, first-order conditions, etc. ...
2. Pick parameters and find the steady state(s)...
3. Log-linearize the necessary equations characterizing the equilibrium of the system to make the equations approximately linear in the log-deviations from the steady state...
4. Solve for the recursive equilibrium law of motion...
5. Analyze the solution via impulse-response analysis....”

²⁹Miao (2014: xvii-xviii) writes, “In earlier years it was quite cumbersome to numerically solve dynamic stochastic general equilibrium (DSGE) models. Students and researchers found it hard to replicate numerical results in published papers. This changed in the 1990s. Researchers finally developed efficient numerical methods to solve medium- to large-scale DSGE models.... These methods were made popular with the launch of Dynare in the late 1990s. Dynare is a software platform for handling a wide class of economic models, and in particular, DSGE models and overlapping generations (OLG) models.”

Uhlig (1999) emphasizes the continuity of his approach with methods previously used to solve RBC models, such as those of King, Plosser and Rebelo (1988a, 1988b, 2002).³⁰

RBC/DSGE methods are widely used in macro-finance. Jermann (1998) uses a perturbation method combined with lognormal pricing approximations. More recently, perturbation methods have become the tool of choice for the solution of consumption-based models with long-run risks (Chen, Cosimano and Himonas 2014, Pohl, Schmedders and Willms 2018). The log-linearization and perturbation technique of Campbell and Shiller (1988) has become so ubiquitous that “it is difficult to find studies that do not rely on the Campbell-Shiller approach—it has become the standard method for solving asset pricing models with long-run risk.” (Pohl, Schmedders and Willms 2018: 1062).

Perturbation methods do not model stochastic processes per se. Instead, perturbation methods analyze stochastic shocks using a Taylor series expansion around a deterministic solution, relying on an early result in control theory demonstrating the equivalence between regular perturbations and stochastic perturbations (Judd 1998: 471-74). In RBC models the stochastic variable of primary interest is total factor productivity, which is taken to be autoregressive with little volatility.³¹

Dynamic stochastic programming is a popular alternative to perturbation methods (Adda and Cooper 2003). Such methods start from the dynamic programming approach, where a time-invariant optimal policy function is iteratively (Ljungqvist and Sargent 2018). A stochastic solution is found by discretizing the state space and modeling state transitions with an approximating matrix (Kushner and Dupuis 2002). Upon embedding the state transitions in the iterative search for the policy function, one finds an optimal state-dependent policy function.

Neither perturbation methods nor dynamic stochastic programming methods admit a dynamic equilibrium. Perturbation methods begin from a steady

³⁰Though King, Plosser and Rebelo (2002) looks like an anachronism, the technical appendix had long been in circulation among researchers in working paper form prior to publication (2002: 111).

³¹See, for example, Mehra and Prescott (1985).

state and generate a time-invariant optimal policy. Dynamic stochastic programming finds a time-invariant optimal policy because the notion of convergence applied to dynamic programming solutions is the convergence of an operator that maintains a fixed point in equilibrium. Both methods are thus premised on a ‘tight’ notion of equilibrium in which changes in the state of the economy are short-lived.

Further, both methods face severe limitations on their ability to represent the stochastic elements of an economic model. Perturbation methods do not characterize state-dependent optimal policies. They model stochastic disturbances as local displacements around the equilibrium of a linearized model. In order for the perturbation representation to be valid, the disturbances must be small and immaterial: small enough to remain within the radius of convergence of the Taylor expansion and the neighborhood for which the linearization is a good approximation; and immaterial in that the disturbance does not change the optimal policy in a discontinuous way. Stochastic dynamic programming methods find state-dependent controls, but require a discretization of the state space and representation of the state dynamics as a matrix. The discretization introduces a source of numerical error, while the matrix representation requires the state dynamics to be stationary. It is questionable whether either numerical approach can represent risk adequately for studies of asset pricing in equilibrium.

Our numerical solution method compares favorably with DSGE methods by representing multi-dimensional risks faithfully and solving for time- and state-dependent optimal policies. The method does not require a steady state as a basis for the solution of the model and no linearizations are necessary. We have devised a framework in which the behavior of the economy recognizes and responds to risk in a genuine way.

In our example, we have been able to specify the risks at the heart of the production process in a fairly intricate way. We have many factors at our disposal relative to the monolithic total factor productivity at the heart of the RBC model, and we have the freedom to link the factors according to theoretical motivations. Risk is present in the model due to the variability of the investment opportunity set, the impact of changes in the investment opportu-

nity set on expected returns from production, and the independent variability of output in each production process. The portfolio-theoretic framework of CIR85a, combined with our solution method, opens the door to a careful disaggregation of economic risk.

By introducing a simple form of heterogeneity into our model's treatment of production, we have been able to appreciate the economic impact of time- and state-dependent control. We see that investment allocations change in response to shocks, dynamically responding to the impact of the shock. We also see that dynamic consumption decisions reflect the predictable component of long-run equilibrium investment opportunities: agents save more when output is currently high, because elevated returns on investment can be expected over the near term. Such adaptations are inconceivable within the RBC framework due to its rudimentary description of technology, homogeneous aggregate production and time-invariant decision rules.

Thus the key innovation of our methodology is its ability to relax the intertemporal equilibrium concept from the steady state concept at the heart of the RBC/DSGE methodology to a true dynamic equilibrium concept. The ability to analyze the dynamic response of the economy and financial asset prices to disturbances outside of a steady state is, in my view, a huge advantage. The weaker sense of intertemporal equilibrium allows for adaptive dynamic behavior that has no feasible way of entering the RBC/DSGE or dynamic programming models typically used in macro-finance. Moreover, by solving our model over a finite horizon with time-dependent controls, we have been able to set aside questions about dynamic stability. Indeed, so long as non-stationary processes do not cause numerical values to explode within the horizon of the model, an absence of stationarity poses no problems for our model.

In our model setting, intertemporal equilibrium entails a planned sequence of optimal behaviors over a given planning horizon. Agents formulate their plans based on their knowledge of the laws governing the system and the initial values of wealth, the distribution of capital among productive processes, and the state variables. The stochastic nature of the problem means that agents formulate their plans pathwise: For each trajectory that the state

variables can take, there is an optimal plan for investment and consumption corresponding to that trajectory. But because uncertainty about state variables cannot be resolved at the outset, intertemporal equilibrium does not imply a single production and consumption plan. Accordingly our notion of equilibrium is one in which risk is essential and irreducible. It is hard to see why any other construct would be preferred for analyzing financial phenomena.

3.4 Conclusions

In this chapter we have implemented a flexible, scalable methodology based on deep neural networks to solve for general equilibrium in a fully stochastic setting. We regard our results as a first proof of concept for a solved macroeconomic model with genuine financial objects. This is a major contribution and advance for macro-finance.

The methodology provides a fast, reliable solution for equilibrium including time-dependent control. Time-dependent control allows for equilibrium solutions outside of a steady state, relaxing the notion of equilibrium needed to implement dynamic economic models. Delivering a methodology to solve dynamic models outside of a steady state makes a dynamic equilibrium analysis accessible numerically, and this, too, amounts to a significant contribution.

We are able to obtain insight into the time-dependent control solution of a CIR economy and analyze responses to shocks from any position in the evolution of the economy. We saw that the approximate solutions found by the neural networks are economically reasonable and quite flexible. Control decisions were also seen to reflect predictability in returns: consumption declined in response to a positive output shock because it implied elevated returns to investment in future periods. In an impulse response analysis, we saw that adjustments on the production side of the economy serve to mitigate output, wealth, and consumption volatility, while driving fluctuations in the equity risk premium. Moreover, we saw that dynamic adjustments in behavior continued well beyond the period of the initial shock. Each of these features

is likely to be misinterpreted or missed all together by consumption-based asset pricing models because endowment economies cannot generate such fluctuations, and by DSGE methods because dynamic responses to shocks are bounded apriori.

Our proof of concept already provides compelling reasons to bury the DSGE methodology. Optimal behavior is constrained in the DSGE framework because it cannot vary over time or depend on the state of the economy. And dynamic responses to shocks in the DSGE methodology are constrained by an apriori imposition of dynamic stability on equilibrium. Under such onerous limitations on behavior, it is difficult for risks to materialize on the scale observed in financial markets, and there are no means by which the economy may transform itself over time through investment.

Our enthusiasm for the potential of the numerical method does not yet extend to our specification of the economy, however. The stylized character of the CIR85a model presented challenges for calibration and the interpretation of our numerical results. The all-capital economy of CIR85a achieves a defensible wealth to output ratio of 4 when the expected rate of return on capital is 25 percent. Our simulations divide this return into a risk-free rate of 24 percent and an equity risk premium of 1 percent. The small equity risk premium reflects the stabilization of output, wealth, and consumption that may be achieved through changes in consumption and the equilibrium allocation, while the large risk-free rate reflects capital's status as the sole limiting factor in production. Though the baseline equity risk premium is small, a negative shock to the economy's primary risk factor produces a significant increase in the equity risk premium, consistent with the behavior of asset prices in a recession, but likely far too large in scale. The present model's limited ability to reproduce stylized facts is a drawback.

At the same time it is not obvious that the results of the stylized CIR economy *ought* to match empirical risk-free rates and equity risk premiums. We need to integrate labor into the process of production to decouple output from returns on capital. Thus we reserve judgment on CIR85a's predictions for asset prices until labor can be incorporated in the model. A labor-augmented CIR economy will furnish a more reliable basis for comparison to the facts.

Defining such an economy is the task of our next chapter.

Chapter 4

A Continuous-Time Macro-Finance Model with Labor and Capital

4.1 Introduction

A primary reason we chose CIR85a as our point of departure for an equilibrium asset pricing model was its conception of production as the ultimate source of risk determining asset prices. In contrast to consumption-based asset pricing models that treat production as an exogenous ‘endowment process’, CIR85a makes production an endogenous, controlled process within the economy.

Even so, the treatment of production in CIR85a is highly stylized and, as we saw in Chapter 3, largely unspecified. As the model stands, the allocation of capital to production processes is the sole production decision available to agents. Consumption and capital accumulation decisions are completely coupled, and all productive resources are continuously deployed to maximize output. And because capital is the only factor of production, seemingly extraordinary returns on capital are required to generate output- and consumption-wealth ratios of an empirically reasonable size. Such an economy generates summary statistics for asset pricing that are difficult to

interpret and bear little resemblance to their empirical counterparts.

One then concludes that the CIR85a model is not yet fit for purpose as the hard core of a theoretical macro-finance paradigm. We would go further and suggest that the absence of labor from production stands in the way of comparing the predictions of CIR85a with macroeconomic data. Moreover the absence of labor from production makes comparisons with competing general equilibrium models difficult. Accordingly, the purpose of this chapter is to expand and enrich the model of production in CIR85a to include labor as a second factor of production. The result is a rich general equilibrium model with risky multifactor production that accounts for the main stylized facts of economics and asset pricing in a realistic and dynamic macroeconomic setting.

4.1.1 Multifactor production with risk

The production side of an RBC model is defined by an aggregate production function combining known quantities of capital and labor. Risk arises in the form of shocks to total factor productivity (TFP) that depress output from its steady-state level. We believe that the treatment of production in RBC models is deficient in several ways, particularly with regard to its handling of risk, and suspect that the equity risk premium puzzle that arises in RBC models may be attributed to its inattention to uncertainty in production. Our model therefore departs from the treatment of production in the RBC paradigm in several ways.

First, we treat labor and capital symmetrically as known ‘stock’ inputs that yield uncertain flows when employed in a production process. The incongruous combination of a stock of capital with a flow of labor effort in aggregate production functions precipitated the famous Cambridge controversies about the measurement of capital and the existence of aggregate production functions.¹ We propose to avoid this incongruity by treating labor as a stock of human capital. The representative agent devotes *known*

¹For a retrospective, see Cohen and Harcourt (2003) and the subsequent comments by Pasinetti, Fisher, Felipe, McCombie, and Greenfield (with responses by Cohen and Harcourt) in the Fall 2003 issue of the *Journal of Economic Perspectives*.

quantities of human and non-human capital to production. Their *yields* in a process of production are uncertain, however, because these productive factors are imperfectly matched to any one production technology, and because their usefulness in any process depends on the ability of entrepreneurs and managers to deploy them effectively in changing economic conditions.² We model the uncertain yields of labor and capital as Itô processes, in a direct generalization of the CIR85a capital yield process. The factor yield processes enter a deterministic production function, whose output is found using the multivariate version of Itô's lemma. The result is an Itô process representing an uncertain flow of output from production.

Second, we abandon the standard Cobb-Douglas production function in favor of production with a constant elasticity of substitution (CES) between capital and labor. The Cobb-Douglas specification assumes a unit elasticity of substitution between labor and capital, which forces marginal rates of substitution between capital and labor to evolve linearly as the capital-labor ratio changes. Empirical evidence suggests that elasticities of substitution in actual production processes are significantly less than unity, so marginal rates of substitution are convex with respect to the capital-labor ratio. In any case, it seems unwise to assume away limitations on the substitutability of labor and capital in a theory concerned with the risks of constructing, allocating and owning capital.

Third, we consider factor-augmenting technical progress in addition to the total factor productivity (TFP) contemplated by the Cobb-Douglas production function. Technological improvements may be embodied in factors, in contrast to disembodied total factor productivity operating at the level of the production process. Our handling of technical progress allows technology to affect factor endowments by making a unit of capital or labor function as if it were more than one unit.³ Uncertainty about the rate of factor-augmenting

²For a classic treatment of the management of uncertain factor yields in production by entrepreneurs and professional managers, see Knight (1921). Entrepreneurial and managerial behavior are not modeled separately and may be regarded as part of the human capital aggregate at the disposal of the representative agent. Future work might profitably revisit this simplification.

³The augmented factors can be viewed as factors in efficiency-equivalent units.

technical progress creates uncertainty about the effective stock of labor and capital being allocated to production. Agents decide on allocations of physical stocks of human and non-human capital, but the *effective* stocks used in production are determined by an exogenous stochastic process describing the state of technology. The factor *supply* uncertainty created by technology is then amplified by factor *yield* uncertainty. We model technical progress as a system of three independent state variables corresponding to TFP, capital productivity, and labor productivity.⁴ The three state variables scale output, capital input and labor input, respectively, in each production process. Shocks to the state variables are propagated to all production processes in proportion to their relative rates of technology adoption.

Fourth, we disaggregate production into multiple process technologies that each have the properties enumerated above. Disaggregating production allows agents to make decisions about the allocation of labor and capital among production technologies.⁵ The other non-standard elements of our model of production ensure that agents take uncertain factor yields, limitations on factor substitutability and uneven changes in technical progress into account when making allocation decisions. An aggregate production function, by contrast, leaves no room for decisions about allocation to impact asset prices.

The model of production at the core of our model thus disaggregates production into multiple stochastic processes distinguished by rates of technology adoption, capital shares, elasticities of factor substitution, and distributions of yields on labor and capital inputs. Our specification of the production process is more flexible than the standard RBC specification, which sets the elasticity of substitution to unity and restricts technical progress to improvements in total factor productivity (TFP). More generally, our model captures risks inherent in making *inputs* ready for production, where RBC

⁴In the growth theory literature the variables correspond to Hicks-neutral, Solow-neutral, and Harrod-neutral technological change.

⁵Allocations among productive process may be interpreted as movements of resources between sectors characterized by sectoral-level production functions, or as changes in the intensity of use for different aggregate production technologies. Though our discussion below inclines to the interpretation in terms of sectors, the technology-switching interpretation is equivalent and may be preferred by some readers.

models only consider risks to the *output* from production. At the same time, our approach remains true to and expands on the central insights of CIR85a by accounting for uncertainty in production and highlighting factor allocation decisions as a primary feature of equilibrium behavior.

4.1.2 Elastic labor supply

Equilibrium asset pricing models in the RBC tradition either ignore labor (e.g., Mehra and Prescott 1985) or incorporate labor, but assume that it is supplied inelastically (e.g., Jermann 1998). Under either treatment of labor, the utility of the representative agent is defined over consumption alone. When consumption is the only argument in the utility function, very high levels of global or local relative risk aversion are needed to reconcile the low volatility of aggregate consumption with the level of the equity risk premium.⁶

Elsewhere in the RBC tradition – and in economic theory more generally – it is commonly recognized that individuals value consumption *and* leisure, leading to a disutility for labor effort and an elastic supply of labor. In a canonical result, King, Plosser and Rebelo (1988a, 1988b) find that variations in labor supply are necessary to generate realistic business cycles in concert with serially-correlated technology shocks and the dynamics of capital accumulation.

Elastic labor supply impacts the analysis of risk aversion. Agents who are free to vary their labor supply are less risk-averse because they can increase their labor effort and ‘work their way out of’ an adverse outcome. Chetty (2006) shows that empirical evidence on the responsiveness of work effort to changes in wages can be used to set an absolute upper bound of 2 for the coefficient of relative risk aversion in a time-separable utility function, while a tighter bound of 1 is defensible under reasonable assumptions about the complementarity of consumption and leisure. Such results are manifestly at odds with double- and even triple-digit estimates from asset pricing stud-

⁶Mehra and Prescott (1985) set risk aversion globally in a CRRA utility function. Jermann (1998) uses a habit-formation specification that makes risk aversion a local feature vis-a-vis the habituated level of consumption.

ies, as well as “the 1 to 10 range” deemed plausible for the coefficient of relative risk aversion in consumption-based asset pricing models (Breedon, Litzenberger and Jia 2019a: 69).

Thus it is surprising that elastic labor supply is largely absent from studies of the equity risk premium. An early survey of the equity risk premium ‘puzzle’ by Kocherlakota (1996) reviews many innovations in functional forms for the utility of consumption and devotes ample attention to market incompleteness, but does not once mention labor supply. A more recent 100-page survey of research on the CCAPM by Breedon, Litzenberger and Jia (2019a, b) contains but a single mention of labor supply in a quotation from Parker and Julliard (2005), in which the authors suggest that changes in labor supply will be absorbed into ‘ultimate’ consumption risk, eliminating the need to model them separately.

Nevertheless, empirical data and RBC models agree that variations in labor supply are an essential part of the economy’s dynamic response to shocks. A model with elastic labor supply allows agents to respond to shocks by choosing among feasible configurations of consumption and labor effort. The models of Mehra and Prescott (1985) and Jermann (1998), by contrast, only allow agents to change their level of investment. For want of additional decision variables, models like Jermann (1998) and Chen (2016) must introduce frictions for the investment process and nonstandard utility functions in order to reconcile production risk to asset prices.

In our model we endow our representative agent with a utility function in which consumption and leisure are non-separable. A complete dynamic equilibrium thus consists of a set of paths for the supply of labor, the allocation of labor to multiple processes, the allocation of capital to multiple processes, and the level of consumption. Elastic labor supply works in concert with risky production to generate a rich set of asset pricing consequences that have been largely ignored.⁷ Labor supply decisions are another mechanism that couples equilibrium asset prices to the conditions of production, and studying the connection sheds new light on the relationship between labor and asset prices.

⁷An honorable exception is Kung (2015).

In the literature on asset pricing, fluctuations in labor income have long been treated as a partial equilibrium problem to be solved through transactions in asset markets. The permanent income model of Friedman (1957) posits that individuals will smooth disturbances to their income through capital market transactions, introducing a ‘hedging motive’ for equity market participation. Subsequent literature has emphasized the consequences of non-insurable labor income for asset prices. Heaton and Lucas (1992) argue that interruptions in labor income increase individual consumption volatility enough relative to aggregate consumption volatility to rationalize the equity risk premium, while Lettau and Ludvigson (2001a) note that negative deviations of aggregate labor income from their long-run values predict increases in the equity risk premium.

In contrast to this line of research, our model treats labor income as an endogenous feature of a general equilibrium with complete markets. Rather than choosing investment allocations *in response to* the characteristics of a stream of labor income, agents choose their labor effort and investment allocations *jointly* in response to the state and structure of the economy. Labor income is not just a problem for the representative agent to solve with capital market transactions, but also a means for the representative agent to solve problems with the structure of capital. One can just as easily attribute a hedging motive to the labor supply decision in a general equilibrium model. Thus in our model, neither labor income hedging demands nor incomplete markets play a role in explaining the equity risk premium. Both are removable artifacts of modeling decisions made in earlier literature.

4.1.3 Challenges to the asset-pricing literature

The consumption-based asset pricing literature has long contended that it is not necessary to model production when studying asset prices. For example, Breeden (1979: 269) claims

it is not necessary to explicitly examine firms’ production decisions and the supply of asset shares, provided that the assumptions made are consistent with optimal behavior of firms

in a general equilibrium model. To be consistent with general equilibrium, prices must be recognized to be endogenously determined. ... The model presented is consistent with endogenously-determined prices if, as assumed, all random shocks to the economy are captured as elements of the state vector.

When all random shocks to the economy are included in the state vector, incomes, output and technology may all be represented by Itô processes that depend on the state of the economy, and indeed we will work out these Itô processes below. Furthermore, we have seen in (2.19) that the risks relevant to asset pricing can be represented in terms of consumption risk and the agent's degree of aversion to consumption risk, so long as consumption risk is defined with respect to all state variables.

But for a model to have non-trivial empirical content, the processes for output, technology and incomes must be determined more precisely than saying they are Itô processes dependent on a suitably rich state. As Hansen (1987: 239) writes,

Competitive equilibrium models of asset prices and aggregate fluctuations have little empirical content without a set of auxiliary restrictions on tastes, technology, and the interaction between observable and unobservable forcing variables. For this reason, explicit modeling of the hypothesized estimation environment is essential. Only in this way can we hope to interpret the time series evidence. ... Analyzing the empirical implications of [suitably-restricted] models also will be facilitated by calculating the stochastic equilibria even if such calculations are numerically intensive. For this reason, the study of stochastic models of asset prices with computable equilibria will continue to be an important avenue of research.

We could not agree more.⁸ In our reading, the consumption-based asset pricing literature has long pursued restrictions on tastes at the expense of

⁸Admittedly our purpose is different from Hansen's. Whereas Hansen (1987) is motivating research into different forms of market incompleteness, we propose to study the consequences of restricting production processes.

restrictions on observable and unobservable forcing variables that determine production possibilities and output. While restrictions on tastes may be informative about the psychology of a representative consumer, equilibrium models of asset prices only have empirical *economic* content to the extent that the risks inherent in capital formation are specified, and decisions about the deployment and allocation of resources in production are part of the definition of equilibrium.

We confront the ‘production-skepticism’ of the asset pricing literature in two ways. First, our choice of CES production functions defines a non-linear technology that is not easily subsumed into the kind of linear endowment process routinely employed in the consumption-based asset pricing literature. There is no good theoretical or empirical reason to insist that asset prices are determined as if technology is linear. At the same time, the CES production function allows us to consider the consequences of factor-augmenting technical change, which enjoys pride of place in the growth literature but plays no discernible role in asset pricing.

Second, we condition the representative agent’s level of risk aversion on his (variable) labor effort. Our decision to incorporate elastic labor supply distances us from the asset pricing literature, while—perhaps paradoxically—bringing us closer to the RBC tradition. As in RBC models, agents vary their equilibrium labor supply in response to external shocks, but we go beyond RBC models by disaggregating production and expanding the set of shocks to which agents may react. We divide shocks into factor yield shocks and technology shocks. Shocks to factor yields represent unexpected outcomes for the yield of labor or capital in a particular production process, while shocks to technology affect output, capital productivity, or labor productivity in all processes. The entire class of factor yield shocks is absent from the asset pricing literature, while the study of technological shocks is often limited to TFP.

On the one hand, our model is constructed to create an economy in which production is too important to be ignored. On the other hand, our model illustrates just how many features of the aggregate economy must be suppressed in order to generate an ‘equity risk premium puzzle.’ The Mehra-

Prescott model is an oversimplification if one wishes to explain asset prices. Upon recognizing the risks of production and endowing our representative agent with micro-founded levels of risk aversion and labor supply elasticity, we obtain an economy with low risk-free rates and a reasonably-sized equity risk premium. In our calibrated economy there is no equity risk premium puzzle. At the same time, we do not proffer our model as a ‘solution’ to the equity risk premium puzzle, because we are not concerned with shoring up the foundations of the RBC paradigm. We believe a more radical theoretical and methodological response to the puzzle is long overdue, and offer our continuous-time model as an example of the direction in which future model-building might go.

We do not contend that our model is the final word on the aggregate economy merely because it has plausible asset pricing consequences. Nor do we claim that the equity risk premium is definitively explained by our model. While our model incorporates a wider-than-usual menu of risks, we have chosen parameter values that make some risks small, while employing simplifications that eliminate other risks.

Instead, we propose an inverted puzzle. The classical equity risk premium puzzle asks why the equity risk premium is so big given that the volatility of consumption is so small. Our model throws an entirely different question into focus: *In a dynamic economy characterized by pervasive uncertainty, why is the equity risk premium so small?*

4.1.4 Plan of the chapter

The next section modifies the baseline CIR85a model of Chapter 2 to incorporate labor and multifactor production. The labor supply decision is introduced along with a new utility function in which labor supply and consumption decisions are non-separable. We define stochastic yield processes for capital and labor, and then derive the stochastic process followed by output for a general production function. After showing how factor-neutral and factor-augmenting technical progress enter the general output process, we specify the production function and calculate its derivatives. The model is

completed by defining the distribution of income and modifying the CIR85a asset pricing equations to incorporate the risks of multifactor production.

Section 3 implements a two-step parameterization procedure for the model. We begin with a review of the empirical literature which permits us to fix parameters for the production processes, the evolution of the state of technology, and the utility function. The undetermined factor yield parameters are then calibrated with the goal of matching several dimensionless quantities describing output, consumption, labor effort and asset prices.

The model is solved in Section 4 using the deep learning technique of Chapter 3. We verify that the character of equilibrium is consistent with the targets used in calibration and study the dynamic behaviors producing the equilibrium. We then study the impulse responses of the model to understand the relative importance of different shocks and agents' dynamic responses to them. In particular, we examine changes in the equity risk premium in response to shocks as a means of studying the origins of aggregate fluctuations. Our model suggests that aggregate fluctuations arise not from technological shocks to aggregate output, but from shocks to the yield of *inputs* in particular production processes. Thus we offer a new explanation of the business cycle that is consistent with the stylized facts of asset pricing. Further, we show that dynamic responses to shocks are far more nuanced than the smooth returns to equilibrium fabricated by perturbation methods.

Section 5 concludes.

4.2 Production with Capital and Labor

Before we transplant a multifactor production function into our model we must consider where to make the incision in CIR85a. We take the scalpel to the wealth accumulation process, which endogenously determines the state variable W for any admissible control sequence.

We begin by rewriting the accumulation process (2.7) slightly to motivate

our development:

$$\begin{aligned}
dW = & W \left[\sum_{i=1}^N a_i \left(\alpha_i dt + \sum_{j=1}^{N+K} g_{ij} dZ_j \right) \right] \\
& + W \left[\sum_{i=1}^N b_i \left(\beta_i dt + \sum_{j=1}^{N+K} h_{ij} dZ_j \right) \right] - C dt
\end{aligned} \tag{4.1}$$

In CIR85a, the first term on the right-hand side of (4.1) is the sum of the yields on capital, weighted by the capital allocation, while the second term is the return on a portfolio of contingent claims. We can rewrite the wealth evolution as follows:

$$dW = \sum_{i=1}^N d\eta_i + W \left[\sum_{i=1}^N b_i \left(\beta_i dt + \sum_{j=1}^{N+K} h_{ij} dZ_j \right) \right] - C dt \tag{4.2}$$

where we recall that $d\eta_i = \alpha_i(a_i W)dt + (a_i W)g_i dZ_t = \alpha_i(\eta_i)dt + (\eta_i)g_i dZ_t$ for all processes i . In CIR85a, $d\eta_i$ is the output of process i . The key to the further development of our model is to determine an expression for $d\eta_i$ when production involves capital and labor.

4.2.1 Labor supply and factor input processes

Our first task is to develop a model for labor input. We define labor input as the yield on human capital and construct the yield on labor in direct analogy with the yield on physical capital.

Let \bar{H} be a fixed endowment of human capital. We introduce a new control function $S_L(W, Y, t)$ mapping to the interval $[0, 1]$. The supply of labor is controlled by S_L , which may be interpreted in our representative agent model as a combination of the relative number of hours worked (the intensive margin) and the labor force participation rate (the extensive margin). Therefore the amount of human capital available for production is $H = S_L(W, Y, t)\bar{H}$.

Labor is supplied elastically and its supply is determined jointly with consumption decisions. Following Cantone et al. (2015) we use a time-

separable specification of the utility function⁹ in which consumption and labor supply are non-separable arguments:

$$U(C, S_L) = \frac{1}{\gamma} \left[(C + 1)^{1-\psi} (1 - S_L)^\psi \right]^\gamma \quad (4.3)$$

Here, γ is the coefficient of relative risk aversion and ψ is the Frisch elasticity of labor supply.

The total amount of labor supplied to the market H is allocated to N production processes by the control vector a^L , with $\sum_{i=1}^N a_i^L = 1$ and $0 \leq a_i^L \leq 1$ for all i .¹⁰ The amount of human capital supplied to process i is $\eta_i^L = a_i^L H = a_i^L S_L \bar{H}$. Thus the amount of labor supplied to each production process is the result of two decisions by the representative agent: how much to work (S_L), and how to allocate effort among multiple activities (a^L).

Extending the analogy with non-human capital, we define the yield of human capital supplied to process i as

$$d\eta_i^L = \alpha_i^L(Y, t)\eta_i^L dt + \eta_i^L g_i^L(Y, t)dZ_t \quad (4.4)$$

The yield of non-human capital supplied to process i follows the law (2.1) as before and is now denoted with superscripts K

$$d\eta_i^K = \alpha_i^K(Y, t)\eta_i^K dt + \eta_i^K g_i^K(Y, t)dZ_t \quad (4.5)$$

where $\eta_i^K = a_i^K W$ as in CIR85a.

We have defined $2N$ factor yield processes that determine the output yielded by N production processes. To allow disturbances to labor input yields $d\eta_i^L$ to vary independently of capital input yields $d\eta_i^K$ for each i , we extend dZ_t to become a $2N + K$ -dimensional Brownian motion. Conforming g_i^K and g_i^L to this definition makes them row vectors of length $2N + K$.¹¹

⁹Cantone et al. (2015) incorporate preference shocks and external habit formation in consumption, while treating S_t as a fixed parameter. By contrast, we suppress habit formation, dispense with preference shocks, endogenize S_L , and reparameterize slightly.

¹⁰The summability condition fixes the scale of S_L , while the positivity condition reflects natural constraints or, if you like, the inability to borrow and lend labor effort.

¹¹Describing the inputs to production as stochastic factor yields is sufficiently hetero-

4.2.2 Output processes

Production is defined by a function F that decomposes into N production functions of similar form:

$$F(d\eta_1^K, \dots, d\eta_N^K, d\eta_1^L, \dots, d\eta_N^L) = \sum_{i=1}^N F^i(d\eta_i^K, d\eta_i^L) \quad (4.6)$$

Because η_i^K and η_i^L are defined by diffusions, the production yield of F is given by its stochastic differential, which may be found using the multivariate form of Itô's lemma (Baldi 2017: 236-7). Let the derivatives of F^i with respect to labor and capital be written F_L^i and F_K^i , respectively. Per Itô's lemma, the stochastic differential of F is

$$\begin{aligned} dF &= \left[F_t + \frac{1}{2} \sum_{i=1}^N F_{KK}^i g_i^{KK} + \sum_{i=1}^N F_{KL}^i g_i^{KL} + \frac{1}{2} \sum_{i=1}^N F_{LL}^i g_i^{LL} \right] dt \\ &+ \sum_{i=1}^N F_K^i d\eta_i^K + \sum_{i=1}^N F_L^i d\eta_i^L \\ &= \left[F_t + \sum_{i=1}^N \left(\frac{1}{2} F_{KK}^i g_i^{KK} + F_{KL}^i g_i^{KL} + \frac{1}{2} F_{LL}^i g_i^{LL} \right) \right] dt \\ &+ \sum_{i=1}^N (F_K^i d\eta_i^K + F_L^i d\eta_i^L) \end{aligned} \quad (4.7)$$

where $g_i^{KK} = (\eta_i^K)^2 g_i^K g_i^{K'}$, $g_i^{KL} = (\eta_i^K \eta_i^L) g_i^K g_i^{L'}$, and $g_i^{LL} = (\eta_i^L)^2 g_i^L g_i^{L'}$. The second equality emphasizes that the processes remain separable.

dox to warrant further comment. In configuring a production process, an entrepreneur contracts with a known number of workers and finances a supply of capital. However the entrepreneur does not know how much labor the workers will yield until they are hired and enter the workplace. In the workplace, workers labor more or less purposefully and coordinate their efforts with their colleagues and the available capital more or less effectively. Similarly, the entrepreneur does not entirely realize what his capital resources will yield until they are employed in a larger design. Accordingly we believe it is enlightening to treat inputs to production as uncertain quantities. As in CIR85a, we allow these uncertainties to be partly inherent in the production process (the first $2N$ elements of dZ_t) and partly determined by the state of the economy (the last K elements of dZ_t).

Thus let us concentrate on the contribution from a single process i :

$$\left(\frac{1}{2} F_{KK}^i g_i^{KK} + F_{KL}^i g_i^{KL} + \frac{1}{2} F_{LL}^i g_i^{LL} \right) dt + F_K^i d\eta_i^K + F_L^i d\eta_i^L \quad (4.8)$$

Expanding $d\eta_i^K$ and $d\eta_i^L$ we obtain

$$\begin{aligned} F_K^i d\eta_i^K + F_L^i d\eta_i^L = & \\ & (F_K^i \alpha_i^K(Y, t) \eta_i^K + F_L^i \alpha_i^L(Y, t) \eta_i^L) dt \\ & + (F_K^i \eta_i^K g_i^K(Y, t) + F_L^i \eta_i^L g_i^L(Y, t)) dZ_t \end{aligned} \quad (4.9)$$

Substituting (4.9) into (4.8), the contribution from each process i to dF is

$$\begin{aligned} & \left(F_K^i \alpha_i^K \eta_i^K + F_L^i \alpha_i^L \eta_i^L + \frac{1}{2} F_{KK}^i g_i^{KK} + F_{KL}^i g_i^{KL} + \frac{1}{2} F_{LL}^i g_i^{LL} \right) dt \\ & + (F_K^i \eta_i^K g_i^K + F_L^i \eta_i^L g_i^L) dZ_t \end{aligned} \quad (4.10)$$

which shows that production with Itô process inputs is again an Itô process.

We can simplify (4.10) by imposing some additional structure on g_i^K and g_i^L . Let dZ_t be ordered so that the first N elements are shocks to capital yields, the next N elements are shocks to labor yields and the last K elements are shocks to the state variables. If we assume that shocks to input factor yields are independent of shocks to other factor yields as well as independent of shocks to the state variables, the row vectors g_i^K and g_i^L may be summarized by single non-zero entries σ_i^K and σ_i^L , respectively. In economic terms, we assume that variations in the yield of a factor in a production process are specific to that input and process, so that a failure of labor to yield its expected product in process i does not spill over to the yield of capital in process i or to the yields of labor in processes $j \neq i$.¹²

Under these simplifying assumptions, the matrix obtained by stacking the row vectors g_i^K and g_i^L reduces to a diagonal matrix of dimension $2N \times 2N$ with a $2N \times K$ matrix of zeros appended on the right. Now $g_i^{KK} = (\eta_i^K \sigma_i^K)^2$,

¹²In our view the first of these assumptions is potentially the most controversial. If a firm marshals its human resources poorly its capital resources are likely to underperform, and vice versa. We proceed under the assumption of independence in light of our ignorance concerning more fundamental quantities such as the conditional means of the factor yields.

$g_i^{KL} = 0$ and $g_i^{LL} = (\eta_i^L \sigma_i^L)^2$ for each process i , and the contribution from process i to dF is

$$\begin{aligned} & \left(F_K^i \alpha_i^K \eta_i^K + F_L^i \alpha_i^L \eta_i^L + \frac{1}{2} F_{KK}^i (\eta_i^K \sigma_i^K)^2 + \frac{1}{2} F_{LL}^i (\eta_i^L \sigma_i^L)^2 \right) dt \\ & + F_K^i \eta_i^K \sigma_i^K dZ_{it}^K + F_L^i \eta_i^L \sigma_i^L dZ_{N+i,t}^L \end{aligned} \quad (4.11)$$

where the extra scripting on the diffusions emphasizes that the relevant elements of dZ_t are to be found at positions i and $N + i$.

4.2.3 Technological change

Recall that the expressions α_i , σ_i^K and σ_i^L that define the factor yield processes (4.4) and (4.5) are functions of the state vector Y .¹³ The functions $\alpha_i(Y)$ and $\sigma_i(Y)$ are the points at which the state of the economy influences production possibilities.

We define Y as the state of technology. Setting $K = 3$, we define Y_1 as the state of total factor productivity, Y_2 as the state of capital-augmenting technical change, and Y_3 as the state of labor-augmenting technical change. In the literature on economic growth, the three dimensions of Y are known as Hicks-neutral, Solow-neutral, and Harrod-neutral technical progress.

Producers will vary in their ability to incorporate technical progress into their methods of production; some will adopt technologies at an above-average pace, while others will lag in adoption or find new developments ill-suited to their operations. We define a process-specific scaling coefficient $\delta^i = \begin{bmatrix} \delta_1^i & \delta_2^i & \delta_3^i \end{bmatrix}$ which adapts the state of technology in the economy to the state of technology in a particular process.

Let the capital and labor yield processes have base yields of π_i^K and π_i^L ,

¹³Per the previous section, the vector-valued function $g_i(Y)$ simplifies to the scalar-valued functions $\sigma_i^K(Y)$ and $\sigma_i^L(Y)$.

respectively, and define¹⁴

$$\begin{aligned}\alpha_i^K(Y) &= \pi_i^K \delta_2^i Y_2 \\ \alpha_i^L(Y) &= \pi_i^L \delta_3^i Y_3\end{aligned}\tag{4.12}$$

The drifts of the factor yield processes are

$$\alpha_i^K(Y)\eta_i^K = \pi_i^K \delta_2^i Y_2 \eta_i^K = \pi_i^K (\delta_2^i Y_2 \eta_i^K) = \pi_i^K K_i\tag{4.13}$$

and

$$\alpha_i^L(Y)\eta_i^L = \pi_i^L \delta_3^i Y_3 \eta_i^L = \pi_i^L (\delta_3^i Y_3 \eta_i^L) = \pi_i^L L_i\tag{4.14}$$

where $K_i = \delta_2^i Y_2 \eta_i^K$ and $L_i = \delta_3^i Y_3 \eta_i^L$ are the *efficiency-equivalent* inputs of capital and labor used in process i . Technical progress captured by Y_2 and Y_3 is factor-augmenting: η_i^K and η_i^L are scaled in the production function to behave like more capital and labor input, respectively. We assume yield volatilities are unaffected by the state of technology, so uncertainty about technical progress affects only expected factor yields, and the functions σ_i^K and σ_i^L reduce to constants.

Applying (4.13) and (4.14), the contribution of process i to output is now given by

$$\begin{aligned}\delta_1^i Y_1 \left(F_K^i \pi_i^K K_i + F_L^i \pi_i^L L_i + \frac{1}{2} F_{KK}^i (K_i \sigma_i^K)^2 + \frac{1}{2} F_{LL}^i (L_i \sigma_i^L)^2 \right) dt \\ + F_K^i K_i \sigma_i^K dZ_{it}^K + F_L^i L_i \sigma_i^L dZ_{N+i,t}^L\end{aligned}\tag{4.15}$$

Total factor productivity Y_1 now appears in the drift term, raising the output of the process proportionately according to the process-specific scaling factor δ_1^i .

¹⁴We use expressions $\delta_j^i Y_j$ to unclutter the notation. More precisely, the adjustments are given by $Y_{jt} \exp((\delta_j^i - \delta_j)t)$, which scale the compounded growth rate δ_j of Y_j to give a process-specific growth rate.

4.2.4 The production function

It remains to specify F^i and its derivatives in (4.15). Any twice-differentiable deterministic function of capital and labor will do. We use the Constant Elasticity of Substitution form (Arrow et al, 1961):

$$F^i(K_i, L_i; \beta_i, \sigma_i) = [\beta K_i^\rho + (1 - \beta)L_i^\rho]^{1/\rho} \quad (4.16)$$

The capital share β controls the relative weights of physical and human capital yield used in the production process.

The parameter $\rho \in (-\infty, 1]$ is calculated from the elasticity of substitution parameter $\sigma \in [0, \infty)$ as

$$\rho = \frac{\sigma - 1}{\sigma}$$

In the completely inelastic case, σ goes to zero, sending ρ to negative infinity. The CES production function reduces to the Leontief form in which K and L are perfect complements. Conversely, in the completely elastic case, σ goes to infinity, ρ approaches 1, and K and L are perfect substitutes. In the Cobb-Douglas case, $\sigma = 1$.

Non-constant returns to scale may be entertained by replacing the outer exponent with ν/ρ , where $\nu > 1$ implies increasing returns to scale and $\nu < 1$ gives decreasing returns. We leave this possibility aside since we are not currently concerned with matters of industrial organization, competition or the life cycle of individual firms.

The parameters β_i and σ_i are a convenient way to introduce heterogeneity into our model, along with the process-specific technology adoption coefficients δ^i . Production processes may be distinguished by their capital shares, their elasticities of substitution, and their relative rate of technology adoption.

The CES production function makes output a genuinely nonlinear function of the inputs. The cost of nonlinearity is a loss of tidy analytical results for the first-order conditions of the HJB equation.

However we like nonlinearity for four reasons. First, a nonlinear pro-

duction function is not easily subsumed into a mere endowment process. If production were a linear transformation of state variables following a linear process, we would again have a linear process that could be treated as primitive.

Second, a nonlinear production process showcases the power and flexibility of our solution procedure. We don't need tidy first-order conditions to solve the model. And whereas standard methods would solve a log-linearized approximation to the nonlinear model in the neighborhood of a steady state, we can let the nonlinearity stand while solving for a dynamic equilibrium.

A third benefit is the flexibility of the CES specification with factor-augmenting technical change. The CES production function can account for imperfect substitution between factors, while factor-specific technology shocks provide another pathway besides TFP through which changes in the state of technology can affect output.

Finally, the CES specification has proved its mettle in macroeconomic models. Cantore et al (2015) show that the choice of a CES aggregate production function has a material impact on business cycle analysis in DSGE models. They estimate a variant of the Smets and Wouters (2007) DSGE model that replaces Cobb-Douglas production with CES production. Their results indicate that the CES variant with an aggregate elasticity of substitution well below 1 outperforms the Cobb-Douglas version significantly. Cantore et al (2015: 149) strikingly conclude "we should dismiss once and for all the use of [Cobb-Douglas production functions] for business cycle analysis."

For the multivariate Itô's lemma we need the derivatives of the CES production function (Sato 1967). The first derivatives of the CES production function are:

$$F_K = \beta \left(\frac{F}{K} \right)^{1-\rho} \quad F_L = (1 - \beta) \left(\frac{F}{L} \right)^{1-\rho} \quad (4.17)$$

and the second derivatives are:

$$F_{KK} = \frac{1-\rho}{F} F_K^2 \left(1 - \frac{F}{K F_K} \right) \quad F_{LL} = \frac{1-\rho}{F} F_L^2 \left(1 - \frac{F}{L F_L} \right) \quad (4.18)$$

The cross-derivative $F_{KL} = \frac{1-\rho}{F} F_K F_L$ does not appear in our specification but may be needed in a more general model that allows factor yield shocks to be correlated within processes.

For specified $\{\beta^i, \sigma^i\}$ the flow of output may now be found by substituting F^i and its derivatives into the evolution of dF^i in (4.15).

4.2.5 Incomes, expected returns and risk premia

In CIR85a the distribution of income is trivial because capital is the only factor of production. Capital ‘earns’ the expected gross rate of return on all production processes, weighted by the allocation of capital, $a'\alpha$. The gross return is comprised of the risk-free rate and a risk premium that we have been calling the equity risk premium. The equity risk premium is the inner product of $K + 1$ market prices of risk and $K + 1$ covariances of wealth with $\{W, Y\}$. Netting the equity risk premium from the gross rate of return leaves the risk-free rate as a residual.

Introducing labor requires defining how the proceeds of production are divided among labor and capital. In competitive equilibrium each factor earns its marginal product.¹⁵ Hence the expected gross rate of return on capital conditional on the state of technology is

$$r_K = \frac{1}{W} \sum_{i=1}^N \delta_1^i Y_1 \left(F_K^i \pi_K^i K_i + \frac{1}{2} F_{KK}^i (K_i \sigma_i^K)^2 \right) \quad (4.19)$$

while the conditional expected gross rate of return on human capital is

$$r_H = \frac{1}{H} \sum_{i=1}^N \delta_1^i Y_1 \left(F_L^i \pi_L^i L_i + \frac{1}{2} F_{LL}^i (L_i \sigma_i^L)^2 \right) \quad (4.20)$$

Actual rates of return include the terms $F_i^K K_i \sigma_i^K dZ_{i,t}$ and $F_i^L L_i \sigma_i^L dZ_{N+i,t}$,

¹⁵We could depart from competitive equilibrium by defining process-specific ‘appropriation shares’ that transfer the product of one factor to the income of another, or by integrating Y out to obtain unconditional expectations and introducing a margin of safety for the income paid to labor. For now we explore the conditional competitive benchmark for simplicity and comparability with the literature, and because the income distribution is not our primary concern.

respectively. Note that the rate of return on human capital depends (through L_i) on the optimal allocation of human capital to productive processes, in exact analogy with the rate of return on non-human capital. We normalize by H because wages are only paid to those who are in the labor force.

Aggregate economic risk is still given by (2.5), and the equilibrium market risk premium and state variable risk premia maintain the form

$$\begin{bmatrix} \phi_W \\ \phi_Y \end{bmatrix} = \begin{bmatrix} \sigma^2(W) & \sigma(W, Y) \\ \sigma(W, Y)' & \sigma^2(Y) \end{bmatrix} \begin{bmatrix} -\frac{J_{WW}}{J_W} \\ -\frac{J_{WY}}{J_W} \end{bmatrix} \quad (4.21)$$

as in (2.23). Because we hold \bar{H} fixed, W is still the only endogenous state variable, and the accounting for the equity risk premium in our model with labor remains the same as in CIR85a. If we had defined an accumulation process for \bar{H} and an aggregate of social wealth $W_0 = \bar{H} + W$, we would have a new numeraire for risk and our system of risk premia would become

$$\begin{bmatrix} \phi_H \\ \phi_W \\ \phi_Y \end{bmatrix} = \begin{bmatrix} \sigma^2(\bar{H}) & \sigma(\bar{H}, W) & \sigma(\bar{H}, Y) \\ \sigma(\bar{H}, W) & \sigma^2(W) & \sigma(W, Y) \\ \sigma(\bar{H}, Y)' & \sigma(W, Y)' & \sigma^2(Y) \end{bmatrix} \begin{bmatrix} -\frac{J_{HW_0}}{J_{W_0}} \\ -\frac{J_{WW_0}}{J_{W_0}} \\ -\frac{J_{YW_0}}{J_{W_0}} \end{bmatrix} \quad (4.22)$$

where ϕ_H is the excess rate of return in equilibrium to holding an additional increment of human capital. In this case, one would arguably net ϕ_H/H from the gross rate of return on labor to obtain a risk-free rate of return on human capital, analogous to an unskilled wage rate. The equity risk premium would also pick up a contribution from the covariance of aggregate wealth with aggregate human capital. However we do not pursue this possibility further here.

Returning to the system (4.21), recall that computation of the equity risk premium in CIR85a was simplified by having closed form expressions for $\sigma^2(W) = GG'$ and $\sigma(W, Y) = GS'$, given in (2.17). We can obtain closed forms once again by breaking a and G into labor- and capital-specific pieces. The state covariance matrix $\sigma^2(Y) = SS'$ remains unchanged.

Based on the simplifications we made in (4.11), define the $N \times (2N + K)$

capital yield covariance matrix G_K as

$$\begin{bmatrix} F_K^1 \sigma_1^K & & 0 & 0 & 0 & 0 & \cdots & 0 \\ & \ddots & & & \ddots & & \vdots & \vdots \\ 0 & & F_K^N \sigma_N^K & 0 & 0 & 0 & \cdots & 0 \end{bmatrix} \quad (4.23)$$

and define the labor yield covariance matrix G_L analogously as

$$\begin{bmatrix} 0 & 0 & F_L^1 \sigma_1^L & & 0 & 0 & \cdots & 0 \\ & \ddots & & & \ddots & & \vdots & \vdots \\ 0 & 0 & 0 & & F_L^1 \sigma_N^L & 0 & \cdots & 0 \end{bmatrix} \quad (4.24)$$

The terms $a'GG'aW^2$ and $a'GS'W$ in (2.17) become

$$\sigma^2(W) = a'_K G_K G'_K a_K W^2 + a'_L G_L G'_L a_L S_L^2 \bar{H}^2 \quad (4.25)$$

and

$$\sigma(W, Y) = a'_K G_K S'W + a'_L G_L S' S_L \bar{H} \quad (4.26)$$

respectively. Accordingly we obtain the following expression for the equity risk premium:

$$\begin{aligned} \phi_W/W &= [a'_K G_K G'_K a_K W^2 + a'_L G_L G'_L a_L S_L^2 \bar{H}^2] \left(-\frac{J_{WW}}{W J_W} \right) \\ &+ [a'_K G_K S'W + a'_L G_L S' S_L \bar{H}] \left(-\frac{J_{WY_i}}{W J_W} \right) \end{aligned} \quad (4.27)$$

Subtracting the equity risk premium from the expected gross rate of return on capital yields the risk-free rate as a residual once again. The computation of $\sigma^2(Y)$ is unchanged, but contingent claims risk premia will be affected by the modification to $\sigma(W, Y)$ in (4.26).

In the dynamic equilibrium of our multifactor model, the equity risk premium (4.27) evolves in response to the endogenous accumulation of wealth, endogenous labor supply decisions, and endogenous factor allocation decisions.¹⁶ The levels of relative risk aversion that determine the market prices

¹⁶A potential extension of our model in which the stock of human capital accumulates

of risk are based on a utility function in which consumption and labor supply are non-separable. Our model has far more levers than consumption-based models to explain the level and variability of the equity risk premium, and our study of the model's impulse responses below will aid in understanding how each feature of dynamic equilibrium mediates the response of the equity risk premium to factor yield and technology shocks.

4.3 Parameterization and Calibration

The model contains several free parameters, which I determine in two steps. (Dawkins, Srinivasan and Whalley 2001) In the first step, I collect empirical estimates for the unknown parameters of the production functions, the state evolution equation, and the utility function. In the second step, I calibrate the few parameters that remain undetermined with reference to a set of target values based on data for the US economy. The parameterization obtained grounds the model in empirical evidence and 'centers' the outputs of the model to create a reasonable baseline from which the properties of the solved model may be studied.

4.3.1 Parameters determined by empirical studies

Some recent studies have developed estimates of production functions and rates of technical change at the sector and industry levels. Estimated elasticities of substitution are available from Young (2013), estimated factor shares from Valentinyi and Herrendorf (2008), and calculated rates of technical progress from Jorgenson, Ho and Stiroh (2005, hereafter JHS).

Young (2013) estimates constant elasticities of substitution at the industry level. Young's industry-level production function estimates show that $\sigma^i < 1$ for most of the 35 industries he studies. His estimates of σ^i become more stable and robust when aggregated into sectors. We adopt Young's (2013: Table 9) generalized instrumental variables estimates for the elastic-

and economic risk is given by (4.22) is left as a possibility for further research.

ity of substitution.¹⁷

Valentinyi and Herrendorf (2008: Table 1) measure capital income shares at the sector level, using data from government statistical agencies. Using duality theory, income shares correspond to the shares β_i in optimal production plans. The authors find important differences in capital shares across sectors that are decomposable into intensities of use for land, structures, and equipment.

Rates of technical change δ_i^j are measured at the industry level by JHS (2005: Table 7.2), who construct constant quality aggregates of capital, labor, and intermediate inputs to study the impact of technological change on gross output at the industry level. JHS find that most technical progress for most industries is embodied in intermediate inputs and improvements in the quality of labor and capital. Disembodied improvements in TFP contribute comparatively little to growth, though important exceptions exist in certain industries like software and computing equipment.¹⁸ The quantity of capital in our model properly encompasses capital and intermediate input, owing to the consumability of the capital good. Hence we combine the contribution from intermediate input with improvements in capital quality to obtain an estimate of technical progress embodied in capital. Similarly, the productive processes we model are best understood as producing gross output rather than value added.

Production process specifications

Due to the mix of sector- and industry-level estimates available from the above studies, we define each of our production processes at the industry level, where the industry is a member of a selected sector. We use industry-

¹⁷Young's GMM estimates produce the same ordering of elasticities across sectors. Young also provides industry-level estimates of factor augmentation (2013: Tables 8A, 8B, 8C) that are conditioned on the CES functional form. They are poorly estimated, however, and not easily aggregated to the sector level.

¹⁸Other studies decomposing industry-level gross output (Bartelsman and Beaulieu 2004, Gullickson and Harper 2002, and Bosworth and Triplett 2003) highlight the importance of the disaggregated approach and find significant heterogeneity in TFP across industries. Further comparisons are difficult due to differences in the underlying data and variations in the measurement of labor input.

Parameter	Value	Sector/Industry/Concept	Source
β_1	0.5400	Agriculture	VH (2008)
σ_1	0.6000	Agriculture	Young (2013)
δ_1^1	0.0190	Agriculture	JHS (2005)
δ_2^1	0.0045	Agriculture	JHS (2005)
δ_3^1	0.0014	Agriculture	JHS (2005)
β_2	0.4000	Manufactured consumption	VH (2008)
σ_2	0.4600	Manufacturing	Young (2013)
δ_1^2	0.0045	Fabricated metal	JHS (2005)
δ_2^2	0.0124	Fabricated metal	JHS (2005)
δ_3^2	0.0011	Fabricated metal	JHS (2005)
β_3	0.3400	Services	VH (2008)
σ_3	0.6800	Services	Young (2013)
δ_1^3	-0.0026	Professional services	JHS (2005)
δ_2^3	0.0188	Professional services	JHS (2005)
δ_3^3	0.0027	Professional services	JHS (2005)
β_4	0.3400	Services	VH (2008)
σ_4	0.6800	Services	Young (2013)
δ_1^4	-0.0001	Other services	JHS (2005)
δ_2^4	0.0176	Other services	JHS (2005)
δ_3^4	0.0022	Other services	JHS (2005)
δ_1	0.0021	41-industry median	JHS (2005)
δ_2	0.0130	41-industry median	JHS (2005)
δ_3	0.0011	41-industry median	JHS (2005)
$\sigma(\delta_1)$	0.0010	41-industry median	Based on JHS (2005)
$\sigma(\delta_2)$	0.0060	41-industry median	Based on JHS (2005)
$\sigma(\delta_3)$	0.0005	41-industry median	Based on JHS (2005)
γ	0.7100	Relative risk aversion	Chetty (2006)
ψ	0.3500	Elasticity of labor supply	Cantone et al (2015)

Table 4.1: Parameters determined with reference to empirical studies. The first four blocks determine the production functions. The fifth block determines the state process. The final block determines the utility function. Sources: VH (2008) = Valentinyi and Herrendorf (2008: Table 1), Young (2013) = Table 9, GIV estimates (rounded), JHS (2005) = Jorgenson, Ho and Stiroh (2005: Table 7.2).

level estimates for productivity growth rates and sector-level estimates for capital shares and elasticities of substitution. Our goal is to define a set of production processes with heterogeneous frontiers and exposures to the

evolution of technology.

We set $N = 4$. Process 1 is modeled on the agriculture industry, which coincides with the agriculture sector. Agricultural production is characterized by a relatively large capital share due to its intensive use of land, elevated rates of total factor productivity growth and relatively small contributions to growth from intermediate inputs.

Process 2 is modeled on the capital-intensive fabricated metal industry, an example selected from the manufacturing sector. Due to the specialized and energy-intensive tasks performed by capital equipment, the elasticity of substitution within the manufacturing sector is the lowest of the four production processes we define.

Processes 3 and 4 correspond to the professional and non-professional service industries, respectively, within an overarching services sector. The services sector is characterized by a larger labor share (smaller capital share) and high rates of factor-augmenting productivity growth, partially offset by slowly-declining total factor productivity. Somewhat surprisingly, the substitutability of capital for labor in the services sector is greatest among the examples we study.

Our parameter choices for the production functions are consistent with other literature. Leon-Ledesma et al (2010: 1344) use an aggregate capital share $\beta = 0.4$, but find that values of 0.3 and 0.6 are also consistent with their qualitative conclusions. Cantone et al. (2015) estimate σ far below 1 for an aggregate CES production function, while Koesler and Schymura (2015) estimate substitution elasticities between 0 and 1 for most sectors, leading them to reject Leontief and Cobb-Douglas production specifications alike.

Our choices for technical progress parameters are more difficult to assess by comparison with other studies. We are not aware of studies that match JHS in scope and granularity. The JHS results are broadly consistent with the earlier study of Jorgenson, Gollop and Fraumeni (1987), in that both support the view that technical progress is predominantly embodied in factors. Though we are convinced the findings of JHS are sound, we would not expect that all economists would agree that TFP plays a relatively minor role in the evolution of technical progress.

Specification for the state of technology

The rates of technical progress that define the state of technology in our model are re-scaled within each production process. Indeed, if we were prepared to treat technical progress as a deterministic time trend, we could dispense all together with specifying a state process for the economy. However technical progress is not assured, so we would like to treat it as a random variable. Fixing the rates of technical progress for the economy at large is useful because it lets us choose their volatilities intelligently. We choose volatilities to assign a small probability (less than 5 percent) to negative rates of technical progress.

Accordingly we use the JHS median estimates to define the process followed by the state of technology as follows:

$$\begin{bmatrix} dY_1 \\ dY_2 \\ dY_3 \end{bmatrix} = \begin{bmatrix} 0.0021Y_1 \\ 0.0130Y_2 \\ 0.0011Y_3 \end{bmatrix} dt + \begin{bmatrix} 0.001 & 0 & 0 \\ 0 & 0.006 & 0 \\ 0 & 0 & 0.0005 \end{bmatrix} dZ_t^Y \quad (4.28)$$

We assume that the processes are independent, and set the volatilities large enough to permit the process increments to be negative with meaningful probability.¹⁹

Note that productivity shocks will not be persistent under our specification, in a departure from the benchmark RBC specification. In our model the dynamic impact of productivity shocks will be determined entirely by changes in time path of wealth accumulation. Like the RBC specification, however, shocks to any of the productivity variables will be perfectly correlated in the cross-section – that is, across production processes.

Utility function

Our utility function (4.3) contains two unknown parameters. Following Cantone et al. (2015) we set the elasticity of labor supply $\psi = 0.35$. Our chosen value is well-supported by other research. Chetty (2012) establishes a range

¹⁹We are not aware of reliable time series for factor augmenting technical change that would permit more precise specification of the state covariance matrix.

of 0.33 to 0.47 for the Frisch elasticity based on microdata estimates, with the low end of the range occurring when the elasticity of intertemporal substitution is zero, a useful benchmark for our time-separable specification. Reichling and Whalen (2012) cull Frisch elasticities ranging from 0.27 to 0.53 from a review of the microeconomic literature, with a central estimate of 0.40.

As discussed in the Introduction, relatively small coefficients of relative risk aversion are compatible with elastically-supplied labor. We set $\gamma = 0.71$ based on the average of labor supply-based estimates presented in Chetty (2006). It bears repeating that our choice of γ is *two orders of magnitude smaller* than the equivalent levels of local risk aversion implied by other equilibrium asset pricing models.

The parameters in the model determined with reference to the above empirical studies are collected in Table 4.1. The sources and concepts pertaining to each are listed alongside for ease of reference.

4.3.2 Calibrated parameters and targets

We have conceptualized the inputs to production as the stochastic yields of labor or capital. As economists do not typically think of productive factors in this way, the empirical literature provides little guidance on reasonable values. Thus we make the simplifying assumption that $\{\pi_i^K, \pi_i^L, \sigma_i^K, \sigma_i^L\}$ are equal for all four production processes, and calibrate their values by targeting selected aggregate values which we now describe.

In Chapter 3 we analyzed the macroeconomic consequences of the CIR85a model in terms of dimensionless output/wealth and consumption/wealth ratios. We would like to do the same here. In order to bound these ratios with empirical data, we refer to data on fixed assets, gross output, and personal consumption expenditure for the United States compiled by the Bureau of Economic Analysis (BEA).

We measure wealth using the current-cost stock of fixed assets and consumer durable goods.²⁰ At the end of 2019 the total value of fixed assets

²⁰Available at <https://fred.stlouisfed.org/series/K1WTOTL1ES000>.

Concept	Lower bound	Upper bound	Source
Fixed assets and consumer durables	7.00×10^{16}	8.49×10^{16}	US BEA
Gross output, private industry	3.36×10^{16}	3.40×10^{16}	US BEA
Personal consumption expenditure	1.45×10^{16}	1.48×10^{16}	US BEA
Output/wealth ratio	0.396	0.486	US BEA
Consumption/wealth ratio	0.171	0.211	US BEA
Labor force participation rate	0.620	0.680	US BLS
Risk-free rate	0.000	0.030	Consensus
Equity risk premium	0.040	0.080	Consensus

Table 4.2: Targets for macroeconomic aggregate values.

and durable goods was 70.74 trillion USD. We round this down to set our lower bound estimate of wealth at 70 trillion. The upper bound allows for wealth to be 20 percent greater than the measured stock of fixed assets and consumer durable goods, as the value of intangible assets is not included in the BEA data.

Flows are measured at the middle and the end of 2019 for comparability with the year-end total of wealth. The gross output of private industry was 33.632 trillion USD in 2019Q2 and 34.052 trillion USD in 2019Q4.²¹ For the same quarters personal consumption expenditure was 14.5 trillion USD and 14.8 trillion USD, respectively.²² Based on these data we establish a range of 39.6 percent to 48.6 percent for the output-wealth ratio and 17.1 percent to 21.1 percent for the consumption-wealth ratio.

The United States Bureau of Labor Statistics (BLS) tracks labor force participation rates for the working-age population.²³ Since 1970 the labor force participation rate has ranged between 59.8 percent and 67.3 percent. The participation rate climbed steadily upward from 1970 to 1990. During the two decades from 1990 to 2010 it hovered between 65.5 and 67.5 percent before declining over the next decade to roughly 63 percent. Thus we adopt 62 to 68 percent as a reasonable range for the share of labor supplied.

As the risk-free rate and the equity risk premium are not observable

²¹ Available at <https://fred.stlouisfed.org/series/GOPI>.

²² Available at <https://fred.stlouisfed.org/series/PCE>.

²³ Available at <https://fred.stlouisfed.org/series/CIVPART>.

Parameter	Value
π_i^K	0.40
π_i^L	0.30
σ_i^K	0.05
σ_i^L	0.02

Table 4.3: Calibrated values for factor yield process parameters.

per se we rely on other scholars' judgment about reasonable levels for the equilibrium risk-free rate and the equity risk premium. Bansal and Yaron (2004) set out to explain an equity risk premium of 6 percent and a "low" risk-free rate. Barro (2007) ballparks the equity risk premium at 4 to 6 percent and the risk-free rate at 1 to 2 percent. Cochrane (2008) pegs the risk-free rate at 1 percent and the equity risk premium at 8 percent. We adopt a range of 0 to 3 percent for the targeted risk-free rate, and 4 to 8 percent for the equity risk premium.²⁴

Our target values are summarized in Table 4.2. We calibrated the remaining parameters in Table 4.3 with the goal of reproducing values within the targeted ranges in Table 4.2. The parameters in Table 4.3 are chosen heuristically rather than through a systematic search to minimize distances between targeted values and actual outcomes. Given the symmetry between labor and capital in our model, the values for labor and capital could well be swapped. However we believe that a higher but more uncertain yield for capital reflects labor's relatively flexible use in production, as well as the tendency for the remuneration of labor to vary within a narrower range than that of capital.

4.4 Solution of the Model

The model is solved using the deep learning-based methodology introduced in Chapter 3. Initially we have a given level of wealth W_0 , a fixed labor

²⁴But see the survey of investment manager opinion by Hammond and Leibowitz (2011) developed under the aegis of the CFA Institute, which produced a range of estimates for the equity risk premium from 0 to 7 percent and much disagreement about its stability over time.

supply \bar{H}_0 , and the allocations $\{\eta_{i,0}^K, \eta_{i,0}^L\}_{i=1}^N = \{a_{i,0}^K W_0, a_{i,0}^L H\}_{i=1}^N$.

We set $W_0 = 4000$ and $\bar{H}_0 = 8000$ on the intuition that the aggregates capitalize a 1:2 ratio of factor incomes. Because we evaluate rates of return and ratios of output and consumption relative to wealth the absolute levels of W and \bar{H} do not matter for the results, just their sizes relative to each other. Capital and labor are assumed uniformly distributed across all processes initially. All of the state variables are started at $Y_i = 1$. We consider the first year of a five-year time horizon. The number of post-solution simulations is reduced from 500 to 200 to save on computing time.

The model estimation procedure now proceeds as follows:

1. Draw $\Delta Z \approx dZ$ from a $\mathcal{N}(0, \sqrt{\Delta t})$ distribution, where Δt is the time step, and compute Y_{t+1} using the discretized SDE:

$$Y_{i,t+1} = Y_{i,t} + \mu_i(Y_t)\Delta t + s_i(Y_t)\Delta Z$$

2. Then for each time step $t = 1, \dots, T$:

- (a) Compute output for each process $d\eta_{i,t+1}$:

$$d\eta_{i,t+1} = dF^i(\eta_{i,t}^K, \eta_{i,t}^L, Y_{t+1})$$

- (b) Compute optimal policies \hat{a}_t^K , \hat{a}_t^L , \hat{C}_t , and \hat{S}_t :

$$\hat{a}_t^K = N_{a^K}(Y_{t+1}, d\eta_{t+1}, W_t | \Theta_{a^K})$$

$$\hat{a}_t^L = N_{a^L}(Y_{t+1}, d\eta_{t+1}, W_t | \Theta_{a^L})$$

$$\hat{C}_t = N_C(Y_{t+1}, d\eta_{t+1}, W_t | \Theta_C)$$

$$\hat{S}_t = N_S(Y_{t+1}, d\eta_{t+1}, W_t | \Theta_S)$$

where N_j denotes the function approximation achieved by the neural network and Θ_j the parameters of the neural network.

- (c) Capture the influence of optimal policies on future values of the

state variables. Compute W_{t+1} from the wealth evolution:

$$W_{t+1} = W_t + \sum_{i=1}^N d\eta_{i,t+1} - \hat{C}_t$$

set η_{t+1}^K based on the physical capital investment decision:

$$\eta_{i,t+1}^K = \hat{a}_i^K(t)W_t$$

and set η_{t+1}^L based on the human capital investment and labor supply decisions:

$$\eta_{i,t+1}^L = \hat{a}_i^L(t)\hat{S}_t\bar{H}$$

3. Compute the loss function \mathcal{L} and the gradient \mathcal{L}_Θ .

The output calculation in step 2(a) refers to the calculation in (4.15) above using a selected production function (4.16), the technological dependencies $\alpha(Y)$, and process-specific sensitivities to factor yield and technology shocks. Step 2(b) introduces the two new controls. The allocation decision now splits into separate decisions \hat{a}_t^K and \hat{a}_t^L regarding capital and labor, and an aggregate labor supply decision \hat{S}_t also appears. Notice that the controls are once again conditioned on a set of variables that is somewhat larger than necessary because they include the current level of output. The new controls produce an additional allocation calculation in step 2(c). Thus step 2(c) determines the level of aggregate wealth and the physical supplies of capital and labor available to each production process at the next time step.²⁵

In the following sections we study the equilibrium of the model, impulse responses to factor yield shocks, and impulse responses to technology shocks.

4.4.1 Dynamic equilibrium

In our model, a dynamic equilibrium consists of time-paths for capital allocations a_i^K , labor allocations a_i^L , labor supply S_L and consumption C . The

²⁵The efficiency equivalent capital and labor available to each production process are determined by the state evolution, while their yields are subject to technology and factor yield shocks.

time-paths are forward-looking plans made at $t = 0$ for the five-year horizon of our model, on the basis of contemporaneous information. In this section we examine the optimal control decisions that constitute equilibrium behavior, as well as the output-wealth ratios, risk-free rates and equity risk premiums that prevail in equilibrium.

Changes in equilibrium behavior and outcomes are driven by the evolution of technical progress. As time evolves, TFP, capital productivity and labor productivity grow from their starting values according to (4.28). Uncertainty concerning the rate of technical progress generates a distribution of potential states of technology. The state of technology shifts the drift terms of the production process yields in (4.15). Uncertainty about the yields of labor and capital in each process interact with uncertainty about the state of technology to generate distributions of output for each process. Our state-dependent equilibrium is defined over this distribution of states, outputs, and the corresponding levels of wealth.

We plot 25th-, 50th- and 75th-percentile outcomes for each feature of equilibrium over the first 12 months of our analysis horizon. The results from our model simulations may thus be read in two dimensions: the horizontal dimension captures time dependence in the equilibrium solution, while the vertical dimension captures dependence on the state of technology and the volatility of factor yields. Wide inter-quartile ranges in the figures indicate that the solution is highly state-dependent, while the location of the median within the inter-quartile range signals symmetry or skewness in the distribution of outcomes. In this way we seek to summarize the stochastic aspect of the dynamic equilibrium elegantly.

Figure 4.1 shows the evolution of capital and labor allocations to the four processes. The processes are ordered from top to bottom, with capital allocations on the left and labor allocations on the right.

Looking across rows, we first notice that capital and labor inputs tend to move together, though comparisons of the vertical axes show that the factors do not move in lock-step. Because the elasticities of substitution are well below unity for all four processes, weak complementarities between capital and labor exist and allocations to capital and labor follow similar time-paths

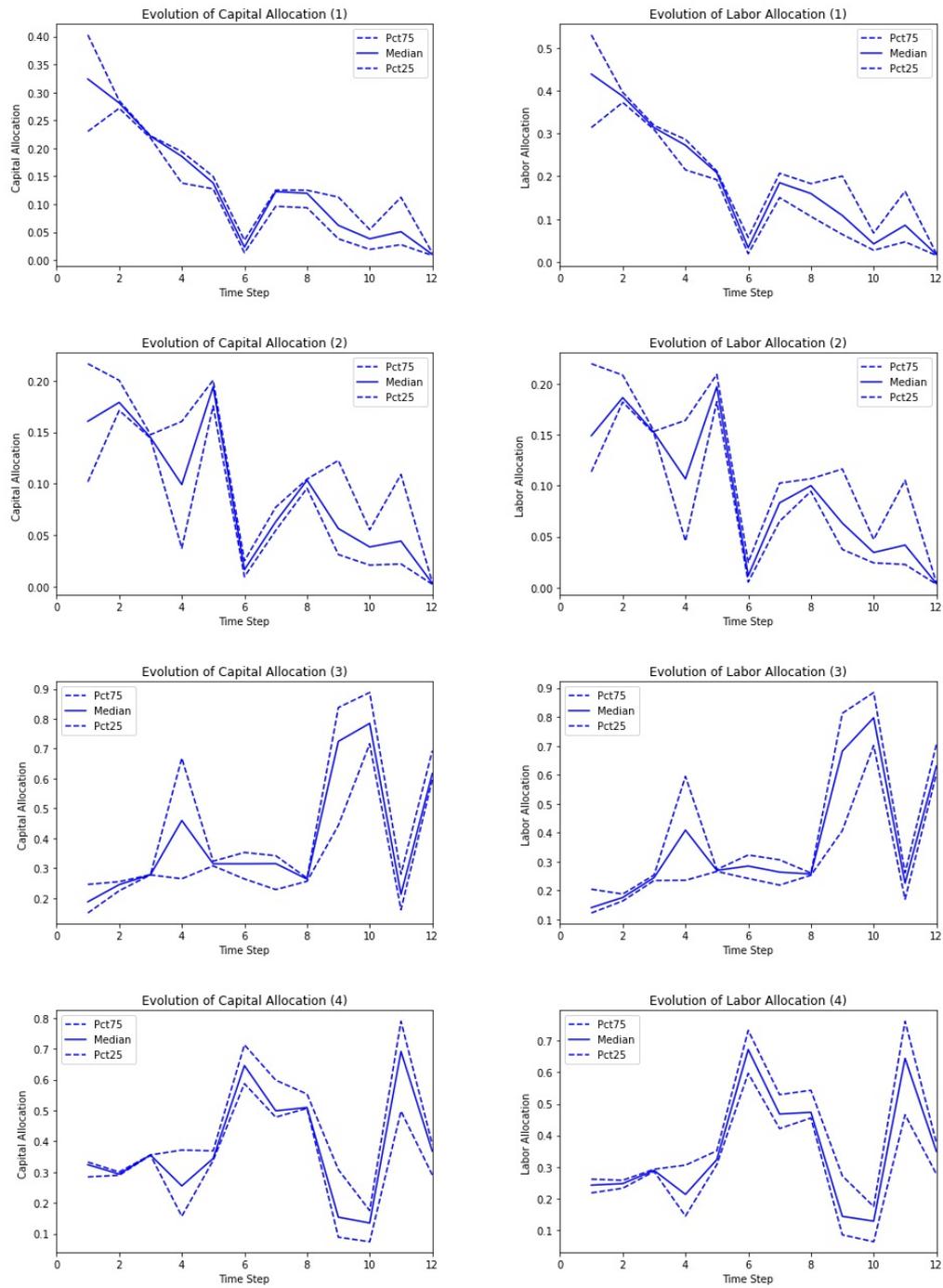


Figure 4.1: Evolution of equilibrium allocations a^K (left) and a^L (right)

in each process.

Comparing column-wise, we see that the relative proportions of capital and labor follow the ordering of the capital shares β_i . Agriculture (process 1) has the largest capital share, followed by manufacturing (2) and services (3 and 4). Accordingly the allocation of labor is large relative to capital in agriculture, comparable in size for manufacturing, and small relative to capital in services.

Overall, we see that capital and labor are re-allocated from agriculture and manufacturing to professional and non-professional services, accomplishing a transformation of the economy in 12 months that required generations in historical time. We suspect that this apparent tendency of time to run in fast-forward is a general property of our model, arising from the representative agent's knowledge of the state process and technology adoption rates δ^i . In reality, the evolution and uptake of technical progress are some of the most difficult aspects of the economy to discern. The speeds of reallocation observed also highlight the difference between our frictionless model and the practical difficulties of reallocating labor and capital across sectors in practice. For example, Davis, Haltiwanger and Schuh (1997) show the majority of labor flows occur within rather than across industries.

Nevertheless it stands that the shifting allocations of our model follow technical progress. Whereas RBC models generally constrain technical progress to be Hicks-neutral, our model lets the pattern of technical progress vary across production processes. Hicks-neutral technical progress dominates in Process 1 only; in the remaining processes technical change is overwhelmingly biased towards capital. Resources flow into those sectors where the efficiency of capital is most improved by technical change. Where capital goes, labor follows.

In addition we notice that uncertainty in the optimal allocation decisions expands and shrinks at the same time across processes. Very low degrees of uncertainty at $t = 2, 3, 5, 8$ alternate with periods of significant uncertainty at $t = 1, 4, 9, 10, 11$. Profiles of the allocation control decisions against the state variables (not shown) indicate that allocations do not depend in any obvious way on variations in the state of technology. The allocation decision

may vary based on shocks to factor yields, depend in a complex way on the state of technology, or some combination of these possibilities.

The strong reallocation of resources across sectors may be a sign that we have set the volatility of technical progress too low, or that we have ignored important correlations between the primary dimensions of technical progress. Time-series data on rates of technical progress are difficult to come by, and we are not aware of any time series of technical progress based on adequate quantity indices for labor and capital.²⁶ Latitude clearly exists for experimentation with different specifications of risk for the evolution of technical progress.

Figure 4.2 shows the equilibrium ratio of output to wealth, the ratio of consumption to wealth, and the share of available labor supplied by households. The top panel shows the ratio of output to wealth beginning at 53 percent, falling to 45 percent and then stabilizing at roughly 40 percent. The interquartile range of about 6-8 percent is fairly tight. Much of the interquartile range falls within our target range for the output-wealth ratio in Table 4.2. We could rationalize more of the distribution of output-wealth ratios implied by our model output by adopting a larger upper bound for the value of wealth.

The second panel of Figure 4.2 shows the consumption-wealth ratio has a central tendency of roughly 20 percent, stabilizing in a range of 18 to 24 percent. These results accord fairly well with the range of 17 to 21 percent targeted by our calibration in Table 4.2. The interquartile range of the state-dependent consumption-wealth ratio is much narrower than that of the output-wealth ratio, suggesting that additions to the capital stock are more volatile than consumption in our model.

The value of S_L , which we interpret as a combination of labor force participation and the number of hours worked, appears in the third panel of Figure 4.2. Equilibrium labor supply is centered around 70 percent with the interquartile range taking values between 64 and 74 percent for the twelve-

²⁶The United States Bureau of Labor Statistics furnishes annual estimates of multifactor productivity. Comparison with JHS reveals the estimates to be unreliable, overestimating the rate of TFP growth and underestimating growth in capital productivity.

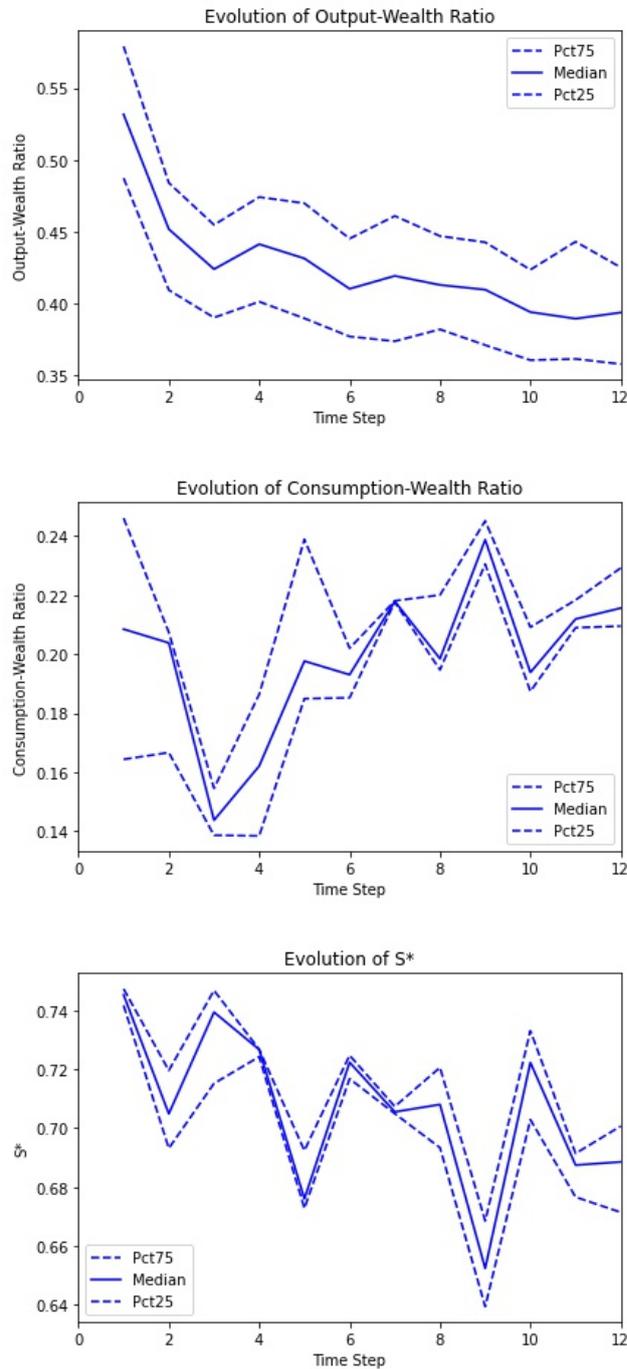


Figure 4.2: Evolution of equilibrium output/wealth ratio (top), consumption/wealth ratio (middle), and labor supply (bottom).

month horizon presented, versus our target of 62 to 68 percent in Table 4.2. Though the values we obtain are somewhat high, they are not unrealistic. Our results compare favorably with those of Cantore et al (2015), who obtain a posterior mean estimate of 53.4 percent. It is also clear from our results that labor supply is quite variable in dynamic equilibrium, suggesting that steady-state restrictions on labor supply may mask meaningful variations in equilibrium behavior.

Figure 4.2 gives us confidence that the calibrated parameter values in Table 4.3 yield results that are consistent with aggregate economic outcomes. In equilibrium the model does a reasonably good job of reproducing the level of gross output relative to wealth, the division between consumption and reinvestment, and the labor force participation rate during the first year of our five-year modeling horizon.

We gain further confidence from the predictions for the risk-free rate and the equity risk premium shown in Figure 4.3. Our equilibrium risk-free rate fluctuates in the central scenario from 2 to 6 percent per annum, while the equity risk premium ranges between 3 to 6 percentage points. Both are mean-reverting, volatile, and uncertain.²⁷

Most studies in the literature obtain numerical solutions by perturbation methods. Because perturbation solutions are fundamentally deterministic, only the mean of the distribution is available for analysis. We need to consider the mean of the distributions in Figure 4.3 to obtain an apples-to-apples comparison to other results.²⁸ In our genuinely dynamic, stochastic solution, we notice considerable variance and skewness in the risk-free rate and the equity risk premium. Downward skew in the distribution of risk-free rates suggests that the mean scenario in our model is below the median, likely in the 1 to 3 percent range. Conversely the upward skew in the distribution of the equity risk premium implies a mean value above the median, taking values between 4 and 8 percent. Upon comparing means to means, our results

²⁷We varied the parameter values in Table 4.3 to examine the sensitivity of our results to our chosen values. While certain configurations could impact our model's predictions for the targeted variables in Table 4.2, the qualitative predictions of the model were consistent with those discussed below.

²⁸We consider the winsorized mean of the values in the interquartile range.

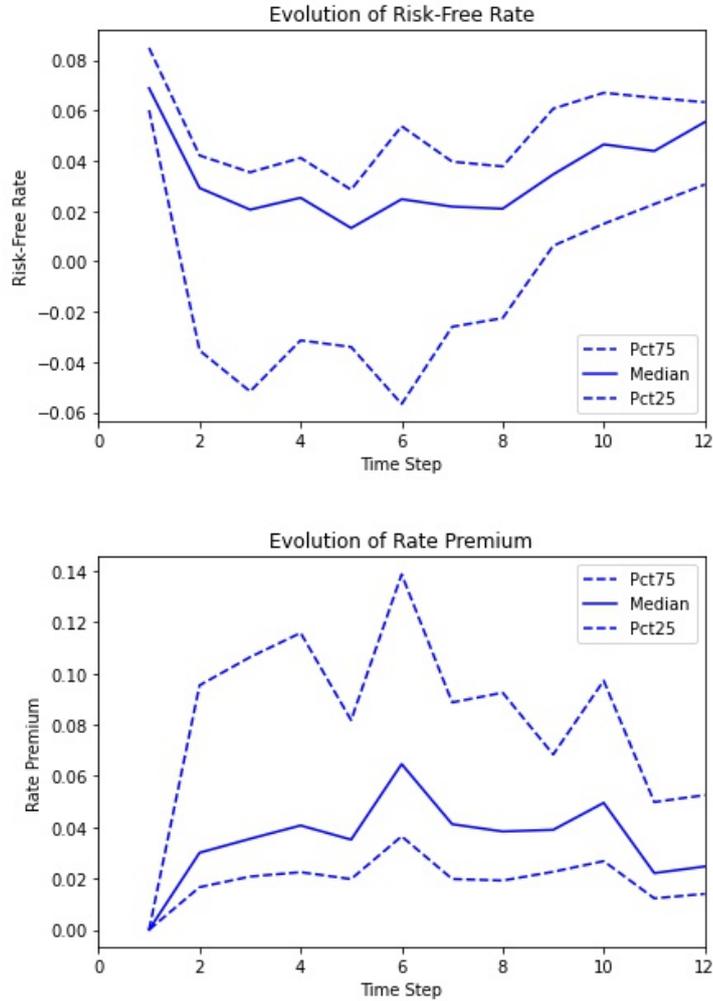


Figure 4.3: Evolution of equilibrium risk-free rate and equity risk premium.

are very much in line with other scholars' judgment about reasonable levels for the equilibrium risk-free rate and the equity risk premium.

Our results are all the more remarkable considering we have obtained them using a simple time-additive utility function and a coefficient of relative risk aversion below one, and without imposing a steady state.

Let us return to the uncertainty surrounding the risk-free rate and the equity risk premium. We emphasize that the distributions in Figure 4.3 are *cross-sectional* distributions across states, rather the moments of a *time*

series distribution. Researchers tend to compare the time-series variances of the risk-free rate and the equity risk premium to empirical time series, in part because they lack access to the cross-sectional distributions. The cross-sectional distributions are comparable to the volatility surfaces obtained from panels of options prices, however, suggesting that comparisons to empirical time series may be inapposite.

The cross-sectional variance of the risk-free rate and equity risk premium distributions is large, while the skewness of the distributions we obtain is consistent with the option volatility ‘smile.’ The lower panel of Figure 4.3 shows that more probability is attached to larger equity risk premiums than smaller ones, which implies that the downside risk of equity prices is greater than the upside risk.²⁹ As a result, we would expect the Black-Scholes implied volatilities of out-of-the-money puts to be larger than those for in-the-money puts, a state of affairs that is commonly observed in the options markets. Because we know that risk aversion plays a relatively minor role in our model, our results suggest the nature of production risk and optimizing behavior play an under-appreciated role in generating the volatility smile.

The risk-free rate, on the other hand, exhibits a negative skew. Negative skewness is a feature of the non-central chi-square distribution for the risk-free rate derived in CIR85b for a simple linear technology. Our model preserves this feature of the CIR85b risk-free rate process, while also allowing the risk-free rate to take negative values with substantial probability. Since our model is developed entirely in terms of real quantities, the risk-free rate is the *real* risk-free rate.³⁰

Thus with our chosen parameter values we achieve a low risk-free rate and a reasonably-sized equity risk premium, both in line with received wisdom concerning their long-run levels. The cross-sectional variance and skewness of the risk-free rate and the equity risk premium, each of which reflect fundamental uncertainty about asset prices, appear to be consistent with the cross-section of equity options prices.

²⁹Recall that prices move inversely to risk premia.

³⁰This observation has far-reaching consequences for the conduct of monetary policy, which we reserve for future research.

Below, we will undertake an impulse response analysis to better understand the drivers of the equity risk premium in our model. The analysis shows that the equity risk premium reflects capital yield risk *within* individual processes more so than productivity risk affecting *all* processes. Shocks to labor and capital yields have far more material effects on equilibrium and asset prices than shocks to productivity.

4.4.2 Responses to factor yield shocks

Because the disturbance vector in our model now has $2N + K = 11$ terms the number of impulse responses that can be studied in the model can quickly become unwieldy, forcing us to be selective in presenting results. In this section we study impulse responses for shocks to factor yields comprising the first $2N$ disturbances, starting with a careful analysis of a single factor yield and progressing to an analysis of the equity risk premium in terms of all factor yields. The next section studies impulse responses for shocks to the state of technology in the last K disturbances.

Changes in equilibrium behavior

We begin by studying a negative one-standard deviation shock to the yield of capital in process 1. Figure 4.4 shows the response of output, consumption, and labor supply. while Figure 4.5 shows the response of capital and labor allocations.

In Figure 4.4 consumption falls by more than output, and then quickly recovers and stabilizes at its baseline level. The drop in consumption permits greater investment in capital, to be carried into $t = 2$. The supply of labor then surges in $t = 3$ after falling slightly. Thus in three periods the optimal response to a shortfall in the yield of capital is to call forth greater supplies of productive factors. By increasing investment and labor effort, the level of output, wealth and consumption is quickly stabilized.

Figure 4.5 shows that changes in labor and capital allocation decisions are small relative to the changes in investment and labor supply. Apart from a delayed response in $t = 10, 11$, allocation decisions are little disturbed. The

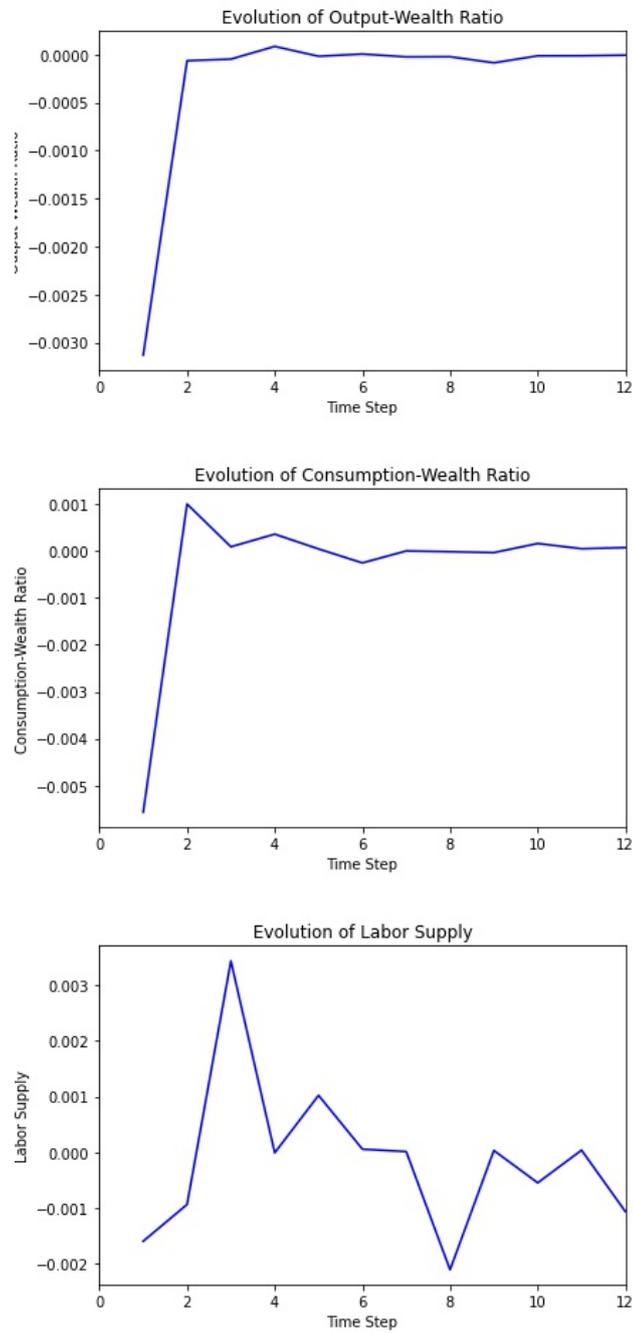


Figure 4.4: Response of output/wealth ratio (top), consumption/wealth ratio (middle), and labor supply (bottom) to capital yield shock in process 1.

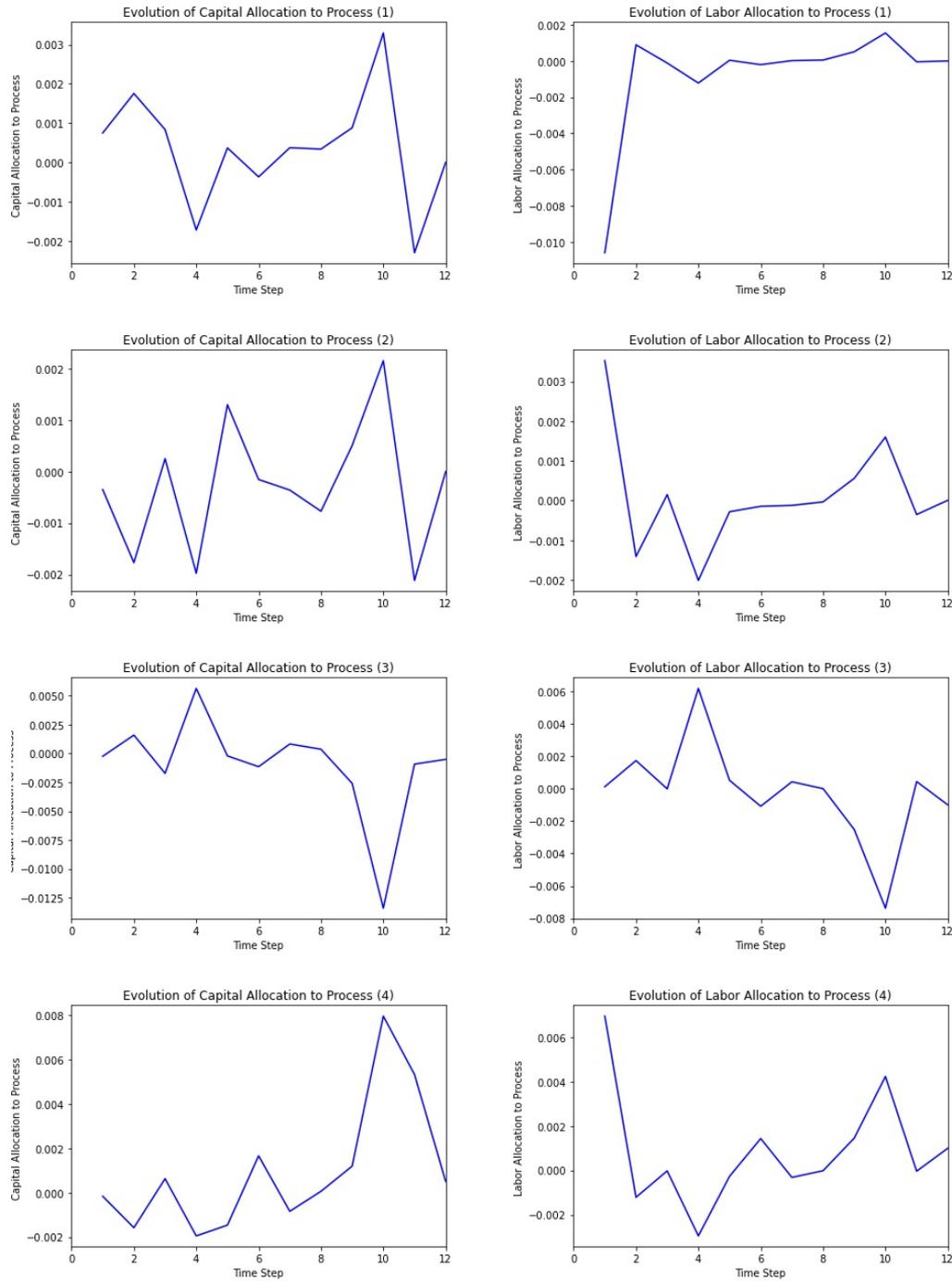


Figure 4.5: Response of a^K and a^L to capital yield shock in process 1.

optimal policy response to the shortfall in capital yield is not to rearrange investment or labor effort, but to supply more of it.

The pattern in responses is similar across processes, but the size of the response varies widely. Shocks to the yield of labor and capital in process 3 cause labor supply to increase by 2 percentage points and lead to large but temporary reallocations of capital and labor to processes 1, 2, and 4. Shocks to labor and capital yields in process 2 produce responses similar to the scale of process 3, while shocks to yields in process 4 produce small responses comparable to those seen in Figures 4.4 and 4.5.

Why do factor yield shocks in processes 2 and 3 produce such large responses, versus the relatively small responses to factor yield shocks in processes 1 and 4? Recall from (4.15) that the contribution of capital yield to the conditional variance of output in process i is $F_K^i K_i \sigma_i^K$, while labor contributes $F_L^i L_i \sigma_i^L$. In our calibration we have set σ_i^K and σ_i^L equal across all processes. Accordingly differences in the magnitude of a reaction to a shock must be driven by differences in the marginal productivity of capital and labor (F_K^i and F_L^i) at the chosen levels of resource allocation (K_i and L_i).

In equilibrium the chosen levels of resource allocation set marginal factor products equal to their compensation. In our model, compensation is given by the expected returns given in (4.19) and (4.20). We conjecture that the magnitude of response to the shock is driven by the curvature terms F_{KK} and F_{LL} that appear in (4.19) and (4.20). which are 2-3 times larger (more negative) for processes 2 and 3 at the equilibrium allocations than for process 4, and an order of magnitude larger than for process 1.³¹ The rate at which the marginal products of labor and capital are changing is greater at the equilibrium allocation for processes 2 and 3, so a shortfall in output requires a larger adjustment to the equilibrium allocation. This occurs because allocation decisions are the primary mechanism our model has to conform factor incomes to productivity.

Equilibrium allocations return to their baseline levels because they are determined by technical progress, not factor yields. The complexion of tech-

³¹The calculation uses the parameter values in Table 4.1 and the initial allocations in Figure 4.1.

nological change appears to furnish a dynamic stability path in our model. Resources will inevitably flow to those uses that are most favored by technological change.

Impact on the equity risk premium

Having traced the dynamics of the equilibrium response to factor yield shocks, we now want to study the effect of factor yield shocks on the equity risk premium. Figure 4.6 considers the impulse response of the equity risk premium to shocks in each of the factor yield inputs. The rows of the figure correspond to processes 1 to 4, while the columns represent the capital yield and labor yield, respectively.

In the dynamic equilibrium of our multifactor model, the equity risk premium (4.27) evolves in response to the accumulation of wealth, labor supply decisions, and factor allocation decisions. The levels of relative risk aversion that determine the market prices of risk are determined by the properties of the utility function at given levels of consumption and labor supply.

Figure 4.6 shows that shocks to capital yields in processes 2 and 3 produce an immediate increase ($t = 2, 3$) in the equity risk premium of close to 200 basis points and 50 basis points, respectively. This immediate increase in the equity risk premium is consistent with the increase in labor supply and curtailment of consumption needed to stabilize output, as shown in Figure 4.4. Accordingly we believe the initial spike is determined by changes in the market price of risk.

A second, delayed increase ($t = 6 - 10$) is evident in Figure 4.6 for shocks to the capital yield in all processes, ranging from 50 to 150 basis points. By this time consumption and labor supply have returned to their baseline levels, so we believe market prices of risk have also returned to normal. Instead, changes at this time scale are more likely to be driven by changes in equilibrium factor allocations prompted by the evolution of technology, and differences in the path of wealth relative to its baseline.³²

³²Increases in the productivity of capital relative to labor will increase the effective capital-labor ratio, leading to increases in the marginal rate of substitution of capital for labor along a convex frontier. The same pattern of factor-augmenting technical change in

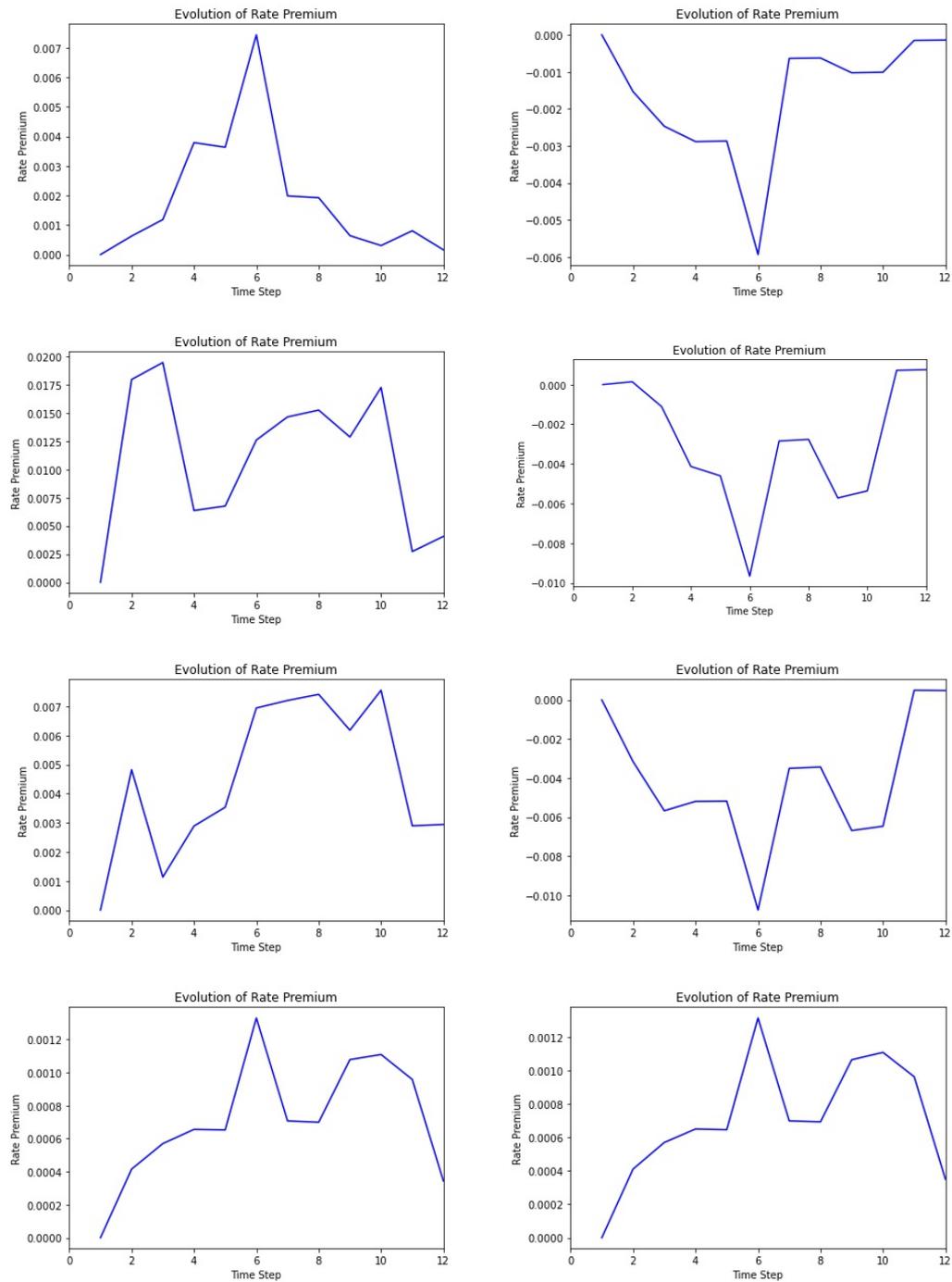


Figure 4.6: Response of equity risk premium to capital (left) and labor yield (right) shocks in processes 1 through 4 (top to bottom).

a Cobb-Douglas production function would overstate the rate of increase in the marginal rate of substitution.

Though allocations of capital and labor have also returned to their baseline levels as of $t = 6$, a second change in equilibrium allocations is evident from $t = 6$ to $t = 10$ in Figure 4.5. Process 1 experiences the shock, but resources are moved away from process 3 and into processes 1, 2, and 4, unwinding movements of resources into process 3 that occurred at $t = 4$. The second spike in risk premiums occurs at $t = 6$. We believe this is an indication that there is too much capital and labor employed in process 3. Asset prices fall until $t = 10$, when the movement of resources out of process 3 is complete.³³

Therefore we discern two effects in the response of the equity risk premium to factor yield impulses, similar to income and substitution effects. The first response is driven by utility losses arising from increased labor effort and reduced consumption. The second response arises from distortions in the allocation of productive resources.

Our explanation goes deeper than the decomposition into ‘payout uncertainty’ and ‘valuation’ effects by Jermann (1998: 267-9), as we tie valuation effects to elastic labor supply and payout uncertainty to production risks arising during the reallocation of resources.

4.4.3 Response to technology shocks

In contrast to the large and meaningful equilibrium responses to factor yield shocks, equilibrium responses to technology shocks are fairly tame. After sifting a few dozen plots of impulse responses, we present the most significant in Figure 4.7. We see that the equity risk premium actually falls by about 40 basis points after an adverse shock to capital productivity.

The relative size of the responses of the equity risk premium to factor yield shocks and productivity shocks is unexpected. Aren’t factor yield shocks ‘idiosyncratic risks’ that do little to determine asset prices, while productivity shocks are ‘systematic risks’ because they impact all production?

Our model challenges us to rethink the facile systematic-idiosyncratic

³³Shocks to labor yield in processes 1 through 3 tend to reduce the equity risk premium after a delay, an effect for which we lack a straightforward explanation.

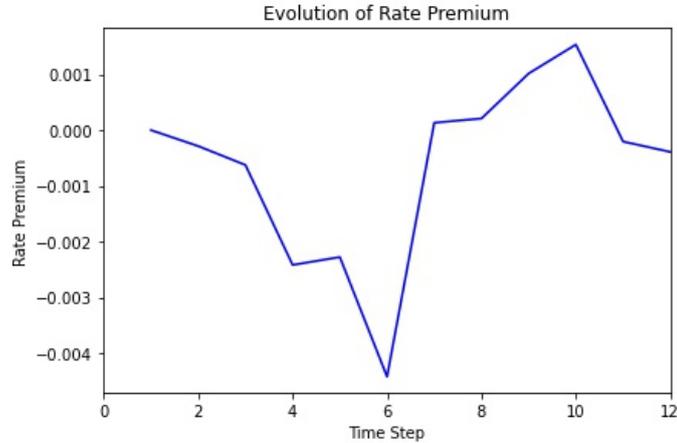


Figure 4.7: Response of equity risk premium to capital productivity shock.

risk classification of the CAPM, which is predicated on the importance of a monolithic ‘market risk factor.’ The expression we have derived for the equity risk premium in (4.27) makes it clear that the equity risk premium depends not only on the variables that comprise the state of technology, but also on individual risks to production (captured in the matrices G_K and G_L), labor supply decisions, and resource allocation decisions. The ‘market risk factor’ is the aggregate of all individual risks, so process-specific risks affect the equity risk premium. Granted, each production process accounts for about a quarter of all output in our model. In the real world, where production is far more finely divided, we would expect individual factor yield shocks to matter less. Nevertheless, some production processes do come to encompass a great deal of economic activity—technology startups at the beginning of the century, housing construction and mortgage finance in the 2000s—and adverse outcomes in these sectors can have a systemic effect on asset prices.

Indeed, we could go further and say that our model offers an alternative explanation for aggregate fluctuations. In contrast to the RBC literature, which emphasizes shocks to technology that affect aggregate output, our model shows that aggregate fluctuations can be generated by shocks to the yields of *inputs* to individual production processes. The sudden collapse of a factor’s yield in a significant production process generates large movements in

output and asset prices consistent with a recession. Unless shocks to technical progress are much larger, more correlated or more persistent than what we have modeled here, they are not sufficient to generate meaningful aggregate fluctuations because the evolution of technology influences allocations, not supplies of productive factors.

We have also seen that the equity risk premium responds to increased investment and labor supply prompted by shortfalls of resources, as well as to inefficient allocations of resources. In our model, adverse shocks to the state of technology do not create such conditions. A lower rate of technological improvement at $t = 0$ reduces output relative to baseline, but does so far less than a shortfall in the yields of capital and labor. And because rates of technology adoption are perfectly correlated across sectors in our model, an adverse shock to the state of technology has minimal consequences for resource allocation. The expected output of all processes falls simultaneously and proportionately.

Shocks to technology may be thought of as perturbations of the dynamically-stable path associated with the most efficient allocation of resources. The response to technological shocks is therefore of second order relative to factor yield shocks. Where factor yield shocks change optimal plans for resource supply, technological shocks change allocation decisions. In a different model with allocation frictions, these second-order considerations may be promoted to the first order.

Our model may understate the risks to technological progress. The variances and covariances of our state transition matrix may be too small, and it is unlikely that the impacts of adverse shocks to technology are perfectly correlated across production processes. A more refined understanding of technological risks and the relative sizes of state and production yield risks may shift explanations of the equity risk premium away from what developments in individual production processes to developments in the state of technology.³⁴

³⁴Certainly there is room to refine and operationalize our understanding of what is meant by technological shocks. Romer (2016) calls TFP “phlogiston” to underscore the mysteriousness of what is actually meant by a negative TFP shock.

4.5 Conclusions

In this chapter we have derived, calibrated and solved a model of risky multifactor production. The model has labor and capital entering a constant elasticity of substitution production process, with multifactor technical change and elastic labor supply. The representative agent has a rich set of choices available to formulate an optimal response to dynamic and uncertain economic conditions. The presence of labor in our model and the match of our modeling outcomes to fundamental stylized facts about both asset prices and the more general economy represent substantial improvements over the CIR85a baseline.

We solve for dynamic equilibrium in our model and study its impulse responses. With our calibrated parameter values, the match between the equilibrium of the model and macro-financial reality is quite good. Our model produces reasonably-sized risk-free rates and equity risk premiums, along with realistic levels for gross output, consumption, savings and labor supply.

Our model suggests an alternative explanation for aggregate fluctuations. We saw that the sudden collapse of capital yields in an intensively-utilized production process can generate significant responses in output, consumption, labor supply and resource allocation when technical progress evolves with relatively little volatility. Conversely, aggregate shocks in TFP and factor-augmenting technical progress produced small responses with signs opposite to those expected. We believe this explanation of aggregate fluctuations is more tangible and provides a more intelligible basis for policy response than fluctuations in disembodied technical know-how.

Instead of driving aggregate fluctuations, the evolution of technical progress in our model serves to steer the optimal allocation of resources. Our model displays the consequences of capital-biased technical progress proceeding at different rates in multiple production processes. Non-neutral technical change dominates factor-neutral changes in TFP. Capital allocations follow technical progress, and labor follows capital due to weak complementarities between the factors. The equity risk premium remains elevated during these

resource movements, reflecting elevated exposure to changes in the investment opportunity set. Thus we conjecture that ‘uneven’ growth leads to higher expected excess returns on equity as the factor bias creates increased uncertainty about the optimal factor allocation. In a representative scenario of an open economy undergoing rapid technical change with a bias toward capital, expected returns on capital will be high, which is to say that the cost of capital for new investment will be elevated.

We have achieved our results by augmenting the treatment of production relative to models cast in the RBC and consumption-based asset pricing molds. In place of an aggregate Cobb-Douglas production function we disaggregate production into multiple CES processes with uncertain factor yields and non-neutral technological change. We allow labor supply to be determined in equilibrium, providing a new means by which behavior can respond to adverse shocks. Solving for a dynamic equilibrium without imposing a steady state, we find that reasonable risk-free rates and equity risk premia may be generated with a coefficient of relative risk aversion below one and other parameter values that are well-supported by empirical research.

We believe these results present three substantial challenges to RBC models and the consumption-based asset pricing paradigm. First, we have shown that the equity risk premium can be generated by a time-separable utility function with a reasonable degree of risk aversion. Neither habit-formation preferences nor recursive Epstein-Zin-Weil preferences are necessary to ‘solve’ the equity risk premium puzzle. Their lack of necessity weighs strongly in favor of eliminating time-dependent preferences from models. As Stigler and Becker (1977) argued decades ago, one should be skeptical about elaborate preference specifications offered as resolutions to economic problems. Complex utility functions try to ‘naturalize’ a problem that has obvious origins, while threatening useful properties like aggregation, time-consistency and invariance to scale. We find it more convincing to explain the origins of the equity risk premium in terms of production decisions than to construct an argument about why people dislike the time series properties of aggregate consumption so much. We replace a litany of preoccupations with the time path of consumption with a theory of dynamic allocation based on the time-

paths for technology and the marginal products of capital and labor. RBC models can't see this because all time paths lead back to the same steady state.

Second, the blindness of the consumption-based paradigm to the disutility of labor effort is simply unsupportable. Instead of supplementing a simple utility function with a preference for leisure, the consumption-based paradigm has spun a theory based on subjective elasticities of intertemporal substitution. Further, ignoring labor supply has meant that much time has been spent on pseudo-problems concerning the hedging of outside labor income with assets in incomplete markets. Including labor supply in a general equilibrium theory resolves the equity risk premium puzzle and dissolves problems created by 'exogenous' labor incomes in a single stroke.

Third, we show that TFP is nearly useless as a basis for the explanation of aggregate fluctuations, which implies that production cannot be reduced to an endowment process in general equilibrium models. Once one breaks the habit of measuring 'Solow residuals' with poorly-constructed capital and labor aggregates, one sees that TFP growth rates are too small to matter very much. And as a theoretical matter, it is difficult to see why output should suddenly fall because of a 'TFP shock'. We treated technical change sympathetically and showed it has meaningful consequences for the dynamic allocation of resources. However technical change was not sufficient to generate aggregate fluctuations. In our model, shocks to the yield of inputs drive aggregate fluctuations, especially when the shock occurs in an intensively-utilized process. Such shocks are more intelligible and better aligned with recent experience. When a stock of capital is suddenly poorly-aligned with market demand – like the redundant assets of internet firms in the dot-com boom or the automated underwriting and mortgage securitization operations of the housing boom – its usefulness in producing valuable output collapses suddenly and unceremoniously.

On a methodological level, augmenting the CIR85a construct to include labor completes the "proof of concept" for our proposed macro-finance paradigm shift. The expanded theory furnishes a laboratory that displays the advantages of our new numerical solution method, and in which the explicit model-

ing of the production side pays clear dividends. We see how important state- and time-dependent solutions are for characterizing dynamic equilibrium and studying the response of the economy to shocks. More importantly, we have shown that our model reproduces a range of macroeconomic and financial market stylized facts without using an elaborate specification of utility. And where we have made changes to the utility function, we found that elastic labor supply lets reasonable levels of constant relative risk aversion be consistent with low risk free rates and an empirically plausible equity risk premium. In our view, modifications of the utility function along the lines proposed here are far more tangible, testable and compelling than the tired and overstretched Epstein-Zin-Weil orthodoxy.

Chapter 5

General Conclusions: A Paradigm for Macro-Finance

We have delivered a proof-of-concept for a new paradigm in macro-finance. We have demonstrated a flexible numerical solution procedure for macroeconomic models that solves models formulated as systems of nonlinear stochastic differential equations. We developed a benchmark general equilibrium model with risky multifactor production, solved it using our numerical methods, and showed the calibrated model can reproduce the risk-free rate and the equity risk premium when the vast majority of the model parameters are pinned down to values estimated in the empirical literature.

The methodological changes we proposed bore fruit. We have expanded the empirical content of macro-finance beyond the theory of the consumer by freeing it from the treatment of production in real business cycle (RBC) models. A model has empirical content not because it has been proven true, but because it leads to propositions that are falsifiable on the basis of experience. We have put forward several new testable statements about general equilibrium and the relationship between production and asset prices. We make new conjectures about the micro-foundations of asset prices in a theory of multifactor production and proffer a competing conception of aggregate fluctuations. Challenging the RBC theory inevitably leads to a differentiated understanding of real business cycles.

At the same time, a proof-of-concept is only that. Our understanding of the proposed methodology and its empirical content remain to be developed by working through extensions, examining the robustness of our conclusions, and discovering ways to confront the predictions of the model with empirical evidence. We have signposted many opportunities to extend the model in previous chapters, and return to catalog them in this section. Our methodology also has limitations and faces many open questions. Some of the most pressing methodological, theoretical and empirical questions are sketched below.

A thesis concerned with methodology is inevitably abstract. Thus I wish to conclude by connecting the insights afforded by our model to some thoughts on economic policy. Our view of the business cycle is sufficiently different from the standard RBC theory that some challenges to the orthodox ‘supply-side’ vision of economic policy are in order.

5.1 Contributions to the Literature

We make three primary contributions to the literature. The first is methodological. We provide a numerical solution procedure for macroeconomic models formulated as systems of stochastic differential equations in continuous time. Our solution procedure characterizes a dynamic state-dependent equilibrium that is inaccessible with standard numerical methods.

The second is theoretical. We provide a benchmark continuous-time general equilibrium model that connects production, consumption and asset prices. Our model highlights decisions about the supply and allocation of productive resources and lets asset prices be determined as known functions of equilibrium behavior.

The third is empirical. The results of our calibrated model suggest a new explanation for aggregate fluctuations, more commonly known as the business cycle. Our calibrated model implies that aggregate shocks to technology are insufficient to explain the business cycle. Instead, we offer shocks to factor yields within particular production processes as a more tangible and plausible explanation of aggregate fluctuations.

5.1.1 Dynamic numerical equilibrium

In the Introduction we pointed out two ways in which the numerical methods that constitute the workhorse DSGE methodology may be constraining macroeconomic modeling. First, we noted that DSGE modeling puts a premium on models for which the steady state may be determined analytically. Second, we suggested that linearization and perturbations calculated by Taylor expansion around the steady state might disfigure the very risks that stochastic models are intended to capture. More to the point, one wonders how DSGE methods can be advertised as ‘dynamic’ when they assume a steady-state equilibrium and ‘stochastic’ when deterministic perturbations around a steady state are expected to capture all of the state-dependent features of the model.

In place of DSGE methods we have put forward a powerful numerical solution procedure that actually characterizes a dynamic, state-dependent equilibrium. The procedure works directly with the model in the form of a system of nonlinear stochastic differential equations. No pre-solution analysis or linearization is necessary. Optimal policies are learned by deep neural networks based on stochastic simulations of the economy over a finite horizon. Thus we trade large, distortionary errors from linearization and local Taylor expansions for relatively benign Monte Carlo error. Our solution procedure is supplemented with an impulse response analysis that enables a study of true comparative dynamics.

Dropping the steady-state restriction from the solution of a macroeconomic model is a big deal. DSGE methods permit a growing economy, but by anchoring the economy to steady-state growth they prevent the economy from *evolving* in any real sense. That the economy an agent encounters after a shock would not look exactly like it did before the shock is inconceivable within the DSGE toolkit. Our numerical method allows the production possibilities of the economy to change over time. Without changing production possibilities one cannot adequately capture risk premia in general equilibrium, as Merton (1971, 1973) and Breeden (1979) and CIR85a have shown. Steady-state equilibrium entails static allocations of resources, while

in dynamic equilibrium, dynamic allocations and risks of changing production possibilities go hand in hand. Now that a choice of methods exists, it is hard to see why one would continue to work with DSGE methods in macro-finance.

When behaviors are not constrained to a steady state, variations in behavior through time can improve outcomes relative to a time-independent strategy – indeed, this is the observation that first motivated study of the intertemporal portfolio problem. For our specification of the economy in Chapter 4, our solution method allows changes in labor effort, resource allocation, and consumption to become part of the character of equilibrium. Our solution method also shows that responses to shocks are far more complex than the smooth returns to the steady state fabricated by perturbation methods. Responses to a shock may occur in multiple waves, and the dynamic equilibrium of the economy that prevails at the end of those responses need not correspond to the economy that existed when the shock occurred.

It is also a big deal that our solution is actually stochastic. The simulations underlying the solution procedure explore the entire state space of the model. There is no need to impose stationarity on the dynamics of the model or to discretize the state space. As a result, the equilibrium behaviors we discover in our solution method characterize dated, state-dependent asset prices, à la Arrow and Debreu. We obtain a full distribution of forward-looking asset prices under the risk-neutral measure – that is, the model generates all of the data one would need to generate theoretical prices of futures, options and other derivative securities. This seems a vast improvement over simulating the moments of the risk-free rate and the equity risk premium from a steady-state model solution.

Finally, though we have been concerned with the deficiencies of DSGE methods, we should not forget that our numerical solution method provides a general strategy for solving the dynamic portfolio problem when asset returns depend on state variables in a predictable way. The problem has long stumped the financial economics profession, as pointed out by Cochrane (2014). It may fairly be said that the current state of the art is still defined by the single state-variable model of Campbell and Viceira (1999). Our

method is far more general, scalable and powerful, and we look forward to exploring the potential advantages of a dynamic asset allocation strategy with our new tool.

5.1.2 Equilibrium with risky production

While our challenge to DSGE methods eliminates many of the numerical obstacles to a macro-financial analysis of risk, better numerical methods alone cannot remedy a lack of risk in the assumptions of the standard RBC model. In the RBC framework production is subject to only one risk: the risk that output suddenly declines due to a total factor productivity (TFP) shock, leaving households with less investable output at their desired level of consumption.

We identified CIR85a as a promising alternative to the RBC framework because it connects asset prices to the risks of production, as we showed in Chapter 2. There was already much to like about CIR85a, including its ability to characterize the risk-free rate, the equity risk premium, and contingent claims risk premia within a consistent equilibrium framework, as well as its emphasis on resource allocation decisions as a primary determinant of equilibrium. In Chapter 3 we addressed the absence of a procedure in CIR85a to obtain an explicit solution of the model. But it remained for us to remedy two theoretical deficiencies in Chapter 4: the absence of labor, and the relatively abstract treatment given to the production process.

Thus in chapter 4 we introduce labor as the stochastic yield of human capital, paralleling the treatment of physical capital in CIR85a. The representative agent supplies his endowment of human capital elastically, choosing how much to work before deciding how to allocate his labor effort across multiple processes of production. By giving the representative agent the ability to control his labor supply, we were able to lower the representative agent's coefficient of relative risk aversion to levels that are consistent with empirical evidence, but unheard-of in the post-Mehra and Prescott (1985) literature on macro-finance. The elastic supply of labor provides another link between asset prices and production in general equilibrium.

Following the introduction of labor, we develop a model of multifactor production in which the stochastic yields of labor and capital produce stochastic outputs according to a production function with a constant elasticity of substitution between labor and capital. Both inputs and outputs are characterized as Itô processes, emphasizing that production is subject to input and output risk. Then we introduced technology in a way that amplified the risks on both sides of production. The risk of shortfalls in input yield is amplified by variations in capital- and labor-augmenting technical progress, while risks to output yield are amplified by variations in TFP. Finally, we disaggregated production into multiple processes to make decisions about the allocation of labor and capital a central feature of general equilibrium.

Like CIR85a, our general equilibrium model characterizes the risk-free rate, the equity risk premium, and the expected excess returns on contingent claims as functions of equilibrium behavior. Our extended model goes beyond CIR85a by grounding asset prices in a theory of multifactor production. We show that risk premia can be represented in equilibrium as a function of marginal factor productivities and their rates of change, process-specific factor yield risks, equilibrium allocations of capital and labor, and the supply of labor. The corresponding market prices of risk incorporate preferences for leisure as well as consumption. The model thus provides micro-foundations for asset prices in terms of equilibrium production decisions. We believe our micro-foundations for risk premia and rates are far richer and more compelling than the time-inseparable preference specifications that have been emphasized in the macro-finance literature.

By emphasizing allocation decisions, our general equilibrium construct can bring macroeconomics into closer contact with research on economic dynamism (Davis, Haltiwanger and Schuh 1996, Davis, Faberman and Haltiwanger 2006, Hsieh and Klenow 2009, Decker et al., 2014) and macroeconomic restructuring (Hopenhayn 1992, Caballero 2007). In addition, our disaggregated specification of production addresses concerns about the existence of aggregate production functions (Fisher 2005) and provides an avenue for research on industrial organization and firm-level productivity growth (Syverson 2011) to inform macroeconomic analysis.

5.1.3 Aggregate fluctuations

The empirical content of our modeling effort consists of the predictions our model makes for responses to adverse shocks in general equilibrium. In short, our challenge to the RBC/DSGE paradigm leads to a different narrative about the origins of business cycles and the dynamics of economic recovery following a recession.

We diverge from the standard RBC theory based on our views on the relative importance of technology shocks and what we call factor yield shocks. We rely on the work of Jorgenson, Ho and Stiroh (2005), which shows that increases in gross output at the industry level are mostly attributable to increases in inputs of labor, capital and intermediate goods. Accounting for the quantity and heterogeneity of inputs (an exercise in index number construction) shows that the residual growth attributable to technological improvement is relatively small. In particular the rate of TFP growth is an order of magnitude smaller than what is estimated by reduced-form growth accounting methods. On this basis, we are confident that rates of technological change are small and exert a bias in favor of capital use. We also believe that the volatilities of underlying technological variables are small because it would be odd for years of technological progress to alternate frequently with years of technological regress.

The calibration of our model to risk-free rates, the equity risk premium, the output-wealth ratio, the consumption-wealth ratio, and the labor force participation rate suggests that expected yields from capital and labor are far larger and more volatile than rates of technological progress. When we compare factor yield shocks to technological shocks in the impulse responses of Chapter 4, we find that only the former are sufficient to generate aggregate fluctuations. Shocks to technology produce comparatively benign responses, changing dynamic equilibrium paths but doing little harm to consumption.

We trace reactions to a negative capital yield shock within a single production process; in all other production processes, capital continues to provide services as expected. The process-specific shock to the yield of capital precipitated a complex response. Initially, agents find themselves unexpectedly

short of capital resources. Agents increase their labor supply and reduce consumption, thereby increasing investment. These efforts to bring a greater supply of productive resources online generate an ‘income’ effect that pushes the equity risk premium higher. The income effect arises because the temporary combination of lower consumption and increased labor effort raises market prices of risk.

The capital yield shock does not disturb the evolution of technology, which continues to change the production possibilities available to the economy as it recovers from the shock. The evolution of technology thus implies a set of dynamic equilibrium paths for the economy, which in turn imply equilibrium allocations of resources. Once agents have increased their endowment of capital and returned to desired levels of labor effort, they find that a different dynamic equilibrium path is optimal. Thus a reallocation of labor effort and capital follows the initial response. The equity risk premium rises again, producing a ‘substitution’ effect. The substitution effect recognizes the prevailing uncertainty about optimal allocations and the risk that further changes in production possibilities may occur.

We hasten to acknowledge the difference between empirical *content* and empirical *verification*; our model’s prediction of such effects is not sufficient to conclude that this is the way the world works. However we believe our analysis is eminently plausible, and far more satisfying – indeed, more *economic* – than frankly psychological explanations that rely on consumption falling below a habituated level, or uncertainty about consumption being resolved later than agents would prefer. We return to the predictions of our model in the final section of this chapter, where we consider how policy interventions may help or hinder the economy’s dynamic response to a recession.

5.2 Roads Not Taken: Potential Extensions

In developing a proof of concept for our proposed paradigm, we have consistently made modeling decisions that favor simplicity and fidelity to the CIR85a framework. These decisions have left us with a long list of possibilities for further development of our paradigm.

5.2.1 Production risk

Though we have emphasized the importance of production risk for equilibrium asset prices, the representation of production risk we chose in Chapter 4 was fairly limited. We assumed that disturbances to factor yields were not correlated with different factor yields in the same process, with the same factor yield in different processes, or with the state of technology.¹ A more general model of production risk would accommodate wider possibilities of dependence between factor yields and the state of the economy.

In addition, while we allowed factor productivities and TFP to evolve over time, we treated the production technology itself as essentially fixed. One could allow various elements of the production function to vary over time, such as the factor weights, the elasticity of substitution and returns to scale. Letting these aspects of the production function vary over time would introduce new possibilities for modeling the evolution of technology and market structure. If these properties of the production function were defined by the state of the economy, production would become subject to additional risks.

We held the endowment of human capital fixed in Chapter 4 to avoid dealing with the complication of a second endogenous state variable. As our most valuable endowment, human capital deserves a deeper treatment comparable to the treatment of physical capital. Human capital could accumulate as a function of experience via past employment or as a function of education via a specialized production process, and depreciate with time or unemployment. Agents would be obliged to manage their human capital endowment much as they manage their wealth. They would undertake activities that add to their human capital without producing current output, and possibly take on sub-optimal employment to limit losses from depreciation. Risks to the accumulation of human capital would also contribute to equity and contingent claims risk premia.

In any reasonable treatment human capital would not be fungible with

¹We are speaking only of the disturbances. *Expected* factor yields depend on the state of technology.

physical capital. Accordingly the possibility of accumulating human capital would require dealing with heterogeneity more deeply than we have thus far.

5.2.2 Heterogeneity

Like CIR85a, our model features a single good that may be consumed or reinvested as capital or an intermediate good. Corresponding to this single good is a single endowment of wealth that changes endogenously with the state of the economy. Endogenous changes are the result of the optimizing behavior of a single representative agent.

Working with one good is a convenience. Distinguishing capital goods and intermediate inputs could have far-reaching consequences. With two or three goods that are neither fungible in production nor fungible in consumption, separate production and allocation decisions would attend to each good. Demand for final output would generate derived demands for capital and intermediate goods.² Capital and intermediate goods would accumulate as separate stores of wealth.

We could pursue these distinctions further, allowing capital and intermediate goods to accumulate separately for each production process. In Chapter 2 we considered an extension to irreversible investment in which capital would accumulate independently in each production process. In this case, the capital allocation decision made by the representative agent would concern *additions* to the stock of capital in each process from current output, while previously-invested capital would remain in place. We could also allow past investment to be reversible at a cost. The representative agent would be able to choose disinvestment amounts for existing processes, but recover only a fraction of the amount disinvested. In a similar way, one could let endowments of human capital accumulate within each production process and reverse at a cost.

²Some modifications to the production function are necessary to floor the levels of capital and intermediate good input. Experimentation with a multi-good model showed that agents will choose to produce consumption goods exclusively with labor if given the chance. With the exception of the Leontief specification, standard production functions do not prohibit such corner solutions.

Yet another helpful disaggregation would extend the menu of productive inputs beyond labor and capital to include renewable and depletable resources. Including renewable and depletable resources would introduce new constraints into the economy from the population dynamics of renewable resources and the sharply rising costs of accelerated resource depletion.

Our model also relies on a representative agent to control and deploy all of the resources of the economy based on his knowledge of the structure of production and the law of motion followed by the state vector. It is hardly a surprise that economists sometimes slip and call him a social planner. We could disaggregate our representative agent into heterogeneous agents along two dimensions. In either case we would confront multiple optimizations, questions about sequencing, and a need for some kind of social choice function that aggregates heterogeneous preferences.

First, we could have different agents control different resources – a representative rentier and a representative laborer, for example, who earn incomes from capital and labor, respectively. Different income sources would motivate investigations into the distribution of income. We defined the income distribution in Chapter 4 according to intuition about competitive equilibrium, but we pointed out that this is by no means the only solution. Entrepreneurs could further de-risk workers from income fluctuations, or shifts in bargaining power could be allowed to influence the division of income.

Second, we could allow agents to have different opinions about the structure of production and the law of motion followed by the state of the economy. Heterogeneous expectations are especially challenging to handle in our model because forward-looking expectations essentially generate the dynamic equilibrium. Agents with different views about the structure of the economy would formulate different equilibrium plans. We would expect some core of feasible plans to emerge, among which there may be no determinate choice. However the existence of heterogeneous expectations would motivate contingent claims trading in equilibrium. The representative agent formulation of CIR85a ensures that holdings of contingent claims are collectively *and* individually zero, which is a somewhat disappointing outcome in a model which otherwise motivates the economic function of contingent claims admirably.

The presence of multiple goods, irreversible investment, and heterogeneous agents may counteract the rapid changes in equilibrium allocations that make time seem to run in fast-forward in our model. We would prefer to see ‘frictions’ emerge from structural choices like those sketched above, rather than being bolted on as artificial costs in the model of production or arbitrary penalties built into the optimization process.

5.2.3 Structural estimation

We base the empirical predictions of our model on a calibration exercise. The calibration could be made more rigorous – or eliminated all together – by turning it into a structural estimation exercise.

We can imagine an estimation method for our model that follows the method of simulated moments. Start with a set of initial values for the parameters to be estimated. In Chapter 4, these are the parameters describing the yields of the capital and labor inputs. Solve the dynamic model – learn the parameters of the neural networks describing optimal policy – given these initial parameter values. Then simulate the model and calculate moments of the model output over some horizon. In Chapter 4 these are the predicted mean values for the output-wealth ratio, the consumption-wealth ratio, the labor force participation rate, the risk-free rate, and the equity risk premium. Perturb all of the input parameters to calculate a gradient, and then take a step in the direction of the empirical moments chosen as targets – i.e., the actual means of the series being predicted. Repeat until convergence. The process would be computationally expensive, but not obscenely so. It is not unusual for structural estimations to require days of computer time.

But even with a feasible estimation method, several problems remain. One must choose data that match the concepts of the model, and if inferences are desired, one must understand the properties of the chosen estimator and formulate appropriate specification tests. Leaving difficult questions about inference aside for now, let us consider the choice of data and the related problem of identification.

The optimal policies in our model are forward-looking plans as of an

initial date, and the macroeconomic outcomes predicted by the model are forward-looking expectations based on those plans. As time moves forward, initial conditions change and agents adjust their plans and expectations. If this is the process generating the empirical data, there is no reason why the simulated time-series moments of the expectations generated by the model at any point in time should match the moments of realized outcomes recorded in empirical time series, especially in a dynamic equilibrium. Instead, it seems that we should be comparing the cross-sectional moments of the expectations generated by the model to the moments of risk-neutral distributions implied by contingent claims prices in financial markets. Thus we should not expect a well-estimated model to explain historical aggregate time series.³ Rather, it will explain the macroeconomic expectations embedded in current asset prices. Thus a consequence of our proposed shift in paradigm may shift the empirical basis of macro-finance from quarterly time series to quoted derivatives prices.

Identification becomes a serious problem as the number of parameters to be estimated grows. In addition to substantial data and pre-processing requirements, one must pin down maximal representations for the state vector and the production processes. In Chapter 3 we referred to the work of Dai and Singleton (2000), who first got financial economists talking about maximal representations. Dai and Singleton worked out the maximum number of free parameters one could have in a three-factor CIR85b model for the discount curve and still identify the model. Similar analysis would be in order for any specification of the production and state evolution processes.

5.2.4 Extensions to include money and banking

The general absence of institutions from our model forces it to be a model of real business cycles, albeit one that is at odds with the RBC theory and methodology. We view the absence of money and credit from our model as lacunae of the same order as the absence of labor in CIR85a. The economy

³Still, it may behoove us to use historical time series for some applications as a means of obtaining results that are comparable to the extant literature.

we live in is a monetary economy with credit. For all the insight our model affords, it is not nearly enough to establish conclusively that business cycles are, always and everywhere, non-monetary phenomena.

Our model determines a real risk-free interest rate in equilibrium. Upon introducing a monetary authority, we would have a nominal risk-free interest rate which ought to equal the real interest rate. Gaps produce inflation. Such a setup disturbs the usual causal ordering running through the Fisher equation from nominal rates to inflation to real rates, though it is consistent with the monetary theory of Wicksell (1935).⁴

The monetary authority is constrained by its connections to the government and the banking system, each of which needs to be modeled. The government determines a budget deficit exogenously, which entails time paths for net transfers and outstanding government debt. The banking system maximizes its output by controlling the size of its lending, subject to constraints on liquidity and capital adequacy. Lending allows the supply of capital to be elastic; investment in capital need not be preceded by saving of output in kind. At the same time lending creates an asset that allows agents to carry output through time without loss. It may be the case that precautionary demands for such an asset exist. The behavior of the banking system adds yet another facet of heterogeneity to the model.

We can imagine the monetary authority learning a ‘reaction function’ to the government’s finances, the state of the banking system, and the state of the economy, subject to a given mandate. Neural networks may discover more complicated policy rules suitable for a complex economy, which could hold obvious interest for the design of monetary policy.

The primary goal of elaborating a monetary version of our model would be to study the effect of monetary shocks. Monetary shocks may originate in the banking sector, as banks engage in a lending boom or reverse course in a crisis of confidence. They may originate in the government, with an unexpected change in the dynamics of the deficit. And they may originate with the behavior of the monetary authority, due to uncertainty about the state of the economy or the behavior of actors in other sectors. How such

⁴Much more so than the ‘neo-Wicksellian’ theory of Woodford (2003).

shocks propagate through the economy and affect the real economy and asset prices are questions of fundamental interest.⁵

With the advent of quantitative easing as a tool of monetary policy, the analysis of actions by the monetary authority extends beyond the correspondence between the policy rate and the equilibrium interest rate to the conditions of balance in the markets for contingent claims. Our CIR-based setting provides a framework in which such policies may be analyzed. Quantitative easing allows for the mass transfer of contingent claims from private agents to the monetary authority in exchange for an infusion of wealth. Given an intervention of this kind, agents are free to choose different paths of consumption, labor supply and factor allocation than those which would prevail otherwise. The changes in plans afforded by quantitative easing lead, in turn, to the repricing of equities and contingent claims via (2.20) and (2.27) and their extensions in Chapter 4. At the moment, however, we have not characterized equilibrium in the economy with $b(W, Y, t) \neq 0$. To do so, we would have to define specific contingent claims in the basis of our model, and set bounds on the values that may be taken by elements of b , which may be endogenous.

5.2.5 Jumps

Our modeling has assumed that the state of the economy follows a continuous multidimensional diffusion process. How might our results change in the presence of jumps?

It would be fairly straightforward to simulate jumps according to a compound Poisson process, for instance, while the simulation of wider classes of Lévy processes is considerably more challenging. Some care would be needed due to the possibility of observing multiple jumps within a discrete time interval. Thus we believe we could continue using our deep learning-based solution method when the state of the economy follows a jump-diffusion process, following suitable modifications to the SDE simulation strategy.

⁵The approach we imagine is quite different than the intermediary or liquidity-based theory of asset pricing elaborated by He and Krishnamurthy (2013, 2019) and Brunnermeier and Sannikov (2014, 2016, 2017).

But what sort of solution is likely to be found for our stochastic control problem in the presence of jumps in the state? In Chapter 4 we saw that the evolution of the state vector traced out a dynamic equilibrium path for factor allocations, labor effort and consumption. Gaussian disturbances to the state vector generated a zone of uncertainty around the optimal policies characterizing the dynamic equilibrium path, reflecting the representative agent's ability to adjust behavior for different realizations of uncertainty in the economy. In the presence of jumps this zone of uncertainty would expand. The memoryless property of the compound Poisson process implies that the agent must be ready at all times for additional jumps, regardless of which jumps have been realized.

Particular interest attaches to the case in which jumps are asymmetric. The representative agent must be particularly on guard against jumps that may suddenly reduce wealth. Given the choice between production processes with equal Sharpe ratios but different downward jump risks conditional on the state vector, the agent will allocate more productive factors to the process with lower jump risk. Thus we expect exposures to jump risks to 'tilt' factor allocations relative to the diffusion benchmark such that the agent is able to carry out his plans with a low probability of exhausting his wealth. Accordingly we anticipate that incorporating jumps into our modeling framework would only reinforce our conviction that RBC/DSGE models arrive at unreliable predictions for asset prices, while suppressing important dynamic equilibrium behaviors.

However in the presence of jumps we face the fundamental difficulty of defining an appropriate contingent claims basis to maintain complete markets. It is a challenging task to define the aggregate risk of the economy in this setting, and to work out the consequences for the equity risk premium, contingent claims risk premia, and contingent claims pricing. Such an investigation would provide a useful structural interpretation for recent work in derivatives pricing, and may shed light on the fraught question of choosing an appropriate change of measure (Cont and Tankov 2004, Eberlein and Kallsen 2019).

5.3 Limitations and Open Questions

As we speculate about the many open possibilities for extending our paradigm, it is important not to lose sight of its limitations. Indeed, many of the extensions suggested above could be construed as limitations arising from the specification of production, collapsed distinctions between goods and agents, and the institution-less set up of the economy. Our reflections on structural estimation have also pointed to questions about the relevant empirical basis for testing the predictions of our model.

5.3.1 Methodological limitations

The introduction of a novel numerical method is likely to be met with suspicion and distrust. We have attempted to assuage concerns about the reliability of our method by following sound mathematical practice, showing that the numerical solutions we obtain are consistent with economic intuition, and profiling the policies learned by the neural network. Nevertheless, much more investigation is needed to provide a completely rigorous defense of our solution method.

Our solution method allows more complicated models to be solved because one need not obtain the first-order conditions of the model analytically. Though we see this as a virtue, it may raise questions about whether the economy is operating in an equilibrium state, as well as whether the equilibrium behaviors learned by the model are indeed optimal. We also define optimality by maximizing a measure of cardinal utility, whereas economists are more comfortable with the existence of ordinal utilities, and the properties of workhorse utility functions are more often justified on ordinal grounds.

The reliability of the solutions we obtain depends in the first instance on the soundness of our calibration method. In our discussion of structural estimation above, we mentioned an open question concerning the choice of data to be targeted. Calibrating to the volatility surface for options on the S&P 500, for instance, is considerably more challenging than matching the mean and variance of quarterly aggregate data. Our calibration is also conditioned, implicitly, on an array of estimates gathered from the empirical

literature. While we believe these estimates are sensible and reliable, that does not necessarily mean that they can be treated as ‘structural’ in a model of this form, and changes in these parameter values may require large changes in calibrated parameter values. We expect that further experimentation with the modeling framework will help to clarify best practice principles for calibration so they reach a level of shared acceptance comparable to calibrations of RBC and neo-Keynesian models.

Solving the model also involves choosing several ‘hyperparameters’ that have purely mathematical importance. In our solution method, choices about the length of the analysis horizon and its division into periods (the time-step) are likely to be the most consequential decisions for analysis, much as they are in the numerical solution of time-dependent partial differential equations. Choices about the number of layers and neurons to include in the neural network architecture may also be important. Research in deep learning shows that optimal neural network architectures can be very problem-specific, and the manner in which neural network approximations converge to a correct functional solution is an open research problem.⁶

Questions about hyperparameters are best resolved by the study of convergence in controlled settings. In a problem with a known solution, one can calculate the error of a numerical solution for a chosen hyperparameter configuration, and examine whether errors behave predictably when varying hyperparameters. For example, errors may usually be expected to decline as the time-step is shortened, and the order of convergence for the errors may be estimated from a plot of the errors against the number of steps. Analyses like this give the modeler confidence that the choice of hyperparameters will lead to a solution in which numerical errors have been controlled suitably.

As a means of getting a tighter grip on the numerical properties of our solution framework, it may be useful to recast some workhorse macroeconomic models in our continuous-time framework. For models with known analytical solutions, we can quantify numerical errors, examine convergence

⁶The universal approximation properties of neural networks are much like the central limit theorem. Just as the central limit theorem can’t tell you how an estimator will perform in finite samples, the mere fact that neural networks are universal approximators does not tell you whether a chosen architecture is adequate to model a particular function.

rates and obtain hyperparameter settings that are compatible with known results for the optimal policies and the character of equilibrium. We can also verify that the solutions our numerical methods obtain for such models are qualitatively consistent with analytical solutions.

DSGE solutions of workhorse models would provide another useful basis of comparison. A battery of numerical solutions using DSGE and deep learning techniques would help to quantify the benefits of trading off log-linearization and Taylor approximation errors for Monte Carlo error. More substantively, we would achieve a deeper appreciation of the differences between steady state and dynamic equilibrium solutions. Discovering how dynamic equilibrium behavior departs from behavior in the steady state, and how the response to shocks in dynamic equilibrium contrasts with the responses predicted by perturbation methods will provide a more detailed accounting of the distortions introduced by DSGE analysis.

5.3.2 Theoretical limitations

On the theoretical side, our model employs objects that are not typically encountered in macroeconomic analysis. Labor and capital inputs are not frequently described in terms of stochastic yields. The apparent strangeness of the factor yields in our model may leave one with the feeling that we have traded one mysterious construct – shocks to total factor productivity – for an equally mysterious construct. It is also not clear what data one might use to obtain reasonable estimates of the expectation and variance of factor yields in empirical production processes.⁷ In Chapter 4 we chose parameter values that were consistent with reasonable scales of output and consumption, but these points of reference are not sufficient to determine parameter values precisely.

Our model has also become fairly complex. For all its problems, there is something appealing about the plunging-your-head-into-cold-water austerity of the basic RBC model. The model is waiting for you around every corner, ready to refer each of your questions back to the same set of assumptions.

⁷Variations in utilization rates and variable compensation come to mind.

Everything that seems to be missing is already there. The RBC model is a zen koan to be repeated to quiet the mind of the working economist. It invites us to meditate on the fundamental truth: there are sudden, persistent drops in output that define business cycles. Let's not worry about why drops in output happen; let's see how rational people react to them. Our formulation lacks this austerity. A generation raised on social media may not have the patience for it.

Though we see the indeterminacy of the state vector as a virtue, others may view it as a liability. Many years of research on dynamic asset pricing have produced little consensus about the number or identity of latent state vectors in even a market as simple as that for Treasury debt.⁸ Pinning down state vectors is important to work out observable consequences of the model for derivative prices, as well as for developing a basic consensus about model specifications.

5.3.3 Empirical limitations

Turning to the data, we face the problem of confronting an empirical basis of historical outcomes aggregated over individuals and time with objects in our model that represent forward-looking expectations. If we were willing to collapse the distinction between history and expectations beloved of 'rational expectations' theories, our only task would be to match the moments of historical time series with those obtained from simulations of our model, perhaps with time aggregation or averaging over states. But a 'solution' of this kind would only serve to destroy the dynamic, forward-looking, state-dependent character of equilibrium in our model, which we view as one of its chief merits.

The cost of maintaining the character of our model is losing contact with the established terms of debate in the asset pricing literature. As in any conflict of paradigms, it is not clear that decisions about the truth of one paradigm can be settled using the testing procedures of another. But this observation does not relieve us of the responsibility to put our results in

⁸See, for example, Rudebusch (2010) and Duffee (2013).

closer dialogue with the consumption-based literature, or the responsibility to identify new protocols for testing the predictions of our model. We have suggested at several points in this thesis that panels of asset prices – interest rate term structures, the risk-neutral measures implicit in options prices, etc. – may furnish the best testing ground for the implications of our model. Condensing these phenomena to a short list of crucial stylized facts would potentially go a long way in clarifying expectations for the empirical performance of our model.

One way in which we might reconcile the predictions of our model to the literature is to replicate some econometric studies using simulations from our model as ‘data’. The empirical predictions of our model would be far more convincing if they could be shown to account for well-known phenomena documented in the empirical literature. To show that the relationship between consumption and risk premia in our model looks a lot like the relationship in a model with time-dependent preferences, for example, would reinforce our claims about the model’s ability to explain asset price movements.

5.4 Implications for Economic Policy

We conclude by considering the implications of our analysis for economic policy, contrasting them with the ‘supply-side’ orthodoxy.

Neither our model nor the standard RBC model contains anything that looks like an economic institution. In both settings, the economy recovers from a shock without the intervention of policy. Alternative economic policies are only analyzable in terms of their consistency with the economy’s return to equilibrium. Policy either helps or hinders recovery, but is not constitutive of the recovery. This is the view of economic policy our model affords, whether or not we endorse it as a philosophical matter.⁹

We are concerned with shocks that generate aggregate fluctuations, i.e., the business cycle. In the standard RBC model, aggregate fluctuations originate from negative, persistent shocks to total factor productivity. Recovery

⁹Spoiler: We believe institutions matter.

occurs as the supply of labor responds to the shock, declining first and then climbing above its steady-state level. Through this lens, policy fails when it gets in the way of the ‘supply-side’ response to the shock in the labor market. Unemployment insurance is an object of particular ire, as it allows for increased consumption without increased labor effort. For adherents of the RBC view, supporting workers with unemployment insurance only deepens recessions and prolongs recoveries.

Our model reproduces the labor market implications of the RBC model, while locating the origins of aggregate fluctuations in particular production processes and showing additional responses relating to the allocation of labor and capital. When a negative shock to capital’s yield in production occurs, the dynamic equilibrium path of labor supply is to fall slightly and then temporarily increase above its baseline level.

Our model decouples the origin of the business cycle from technology. Following the initial capital yield shock, technological change continues apace and its factor biases change the relative desirability of alternative production processes. Accordingly optimal resource allocations continue to change while the economy is recovering from the capital yield shock. Economic policy should support the migration of resources out of adversely affected sectors. As a result, a desire to keep workers in jobs does not necessarily mean keeping them in their current positions.

The optimal response to a capital yield shock also includes an increase in investment (reduced consumption). Policy makers often fret that transfers made in recessions will be saved rather than spent. Our model suggests that we should want people to do exactly that, with two important caveats. First, saving in our model entails new capital formation rather than the accumulation of bank balances or the purchase of shares, which are claims to returns on *existing* capital. Political discussions of ‘investment’ often confuse the former with the latter. Second, newly-formed capital should not be deployed according to existing allocations. Just as optimal labor allocations continue to evolve according to technical change, optimal capital allocations will change, too.

Thus, in a recession, the economic policy response must resist the temp-

tation to ‘bail out’ or support badly-impacted industries. While no politician wants to openly advocate bankruptcy, a more positive way to support the reallocation of capital and labor in a recession would be to support the formation of new ventures. Policymakers should not pick winners and losers but try, for a time, to relax obstacles in the way of business and capital formation. Our analysis of asset prices reproduces the stylized fact that expected excess returns to capital are elevated during recessions. If economic policy were to support entrepreneurs and owners of capital during the recovery from a recession, that support would arguably produce windfall profits. A levy on those profits once recovery is underway would not be unseemly.

Monetary policy is particularly prone to reinforce the existing configuration of capital in a recession. Central banks lend against security, as do commercial banks. Accordingly supportive monetary policy is likely to support commercial bank loans secured against assets in place at existing firms at the expense of unsecured loans and new loans secured against newly-acquired assets. Some support for unsecured lending that supports hiring and overtime payments may be a better idea than the classic Lombard Street discounting of good collateral, though we freely admit that making such an arrangement operational is fraught with unsolved problems.

Though these tentative conclusions are open to criticism, the policy lessons suggested here highlight the need for greater attention to the composition of the capital stock in ‘supply-side’ thinking, and underscore the differences between prescriptions based on a dynamic equilibrium theory and those predicated on a return to a steady state. In a dynamic economy, economic policy must take pains to avoid distorting the allocation of resources by resisting the urge to favor incumbents and particular sectors of activity.

We expect that continued development of our paradigm to refine the core theory and model the activity of economic institutions will reveal further insights about the origins of aggregate fluctuations and the conduct of economic policy.

References

Adda, Jerome and Russell Cooper. 2003. *Dynamic economics: quantitative methods and applications*. MIT Press.

Adjemian, Stéphane, Houtan Bastani, Michel Juillard, Frédéric Karamé, Junior Maih, Ferhat Mihoubi, Willi Mutschler, George Perendia, Johannes Pfeifer, Marco Ratto, and Sébastien Villemot. 2020. *Dynare reference manual version 4*. Dynare Working Papers 1, CEPREMAP.

Al-Aradi, Ali, Adolfo Correia, Danilo Naiff, Gabriel Jardim, and Yuri Saporito. 2018. Solving nonlinear and high-dimensional partial differential equations via deep learning. *arXiv preprint*.

Arrow, Kenneth J. 1971. *Essays in the theory of risk-bearing*. Amsterdam: North-Holland.

Arrow, Kenneth J., Hollis B. Chenery, B. S. Minhas, and Robert M. Solow. 1961. Capital-labor substitution and economic efficiency. *Review of Economic Studies* 43(3): 225-250.

Bansal, Ravi and Amir Yaron. 2004. Risks for the long run: a potential resolution of asset pricing puzzles. *Journal of Finance* 59(4): 1481-509.

Bansal, Ravi, Dana Kiku, Ivan Shaliastovich, and Amir Yaron. 2013. Volatility, the macroeconomy and asset prices. *Journal of Finance* 69(6): 2471-511.

Barberis, Nicholas. 2000. Investing for the long run when returns are predictable. *Journal of Finance* 55(1): 225-264.

Barro, Robert. 2007. Rare disasters, asset prices, and welfare costs. NBER Working Paper 13690.

Bartelsman, Eric J. and J. Joseph Beaulieu. 2004. A consistent accounting of U.S. productivity growth. Federal Reserve Board Finance and

Economics Discussion Series 2004-55.

Beck, Christian, Weinan E, and Arnulf Jentzen. 2019. Machine learning approximation algorithms for high-dimensional fully nonlinear partial differential equations and second-order backward stochastic differential equations. *Journal of Nonlinear Science* 29: 1563-1619.

Bensoussan, Alain, Jens Frehse, and Philip Yam. 2013. Mean field games and mean field type control theory. Springer-Verlag.

Berk, Jonathan B. 2018. What I learned from Steve Ross. *Journal of Portfolio Management* 44(6): 70-73.

Bertsimas, Dimitris and Andrew W. Lo. 1998. Optimal control of execution costs. *Journal of Financial Markets* 1(1): 1-50.

Black, Fischer and Myron Scholes. 1973. The pricing of options and corporate liabilities. *Journal of Political Economy* 81: 637-359.

Boguth, Oliver and Lars-Alexander Kuehn. 2013. Consumption volatility risk. *Journal of Finance* 68(6): 2589-615.

Bosworth, Barry P. and Jack E. Triplett. 2003. Services productivity in the United States: Griliches' services volume revisited. Brookings Institution.

Brav, Alon, George M. Constantinides, and Christopher C. Gezcy. 2002. Asset pricing with heterogeneous consumers and limited participation. *Journal of Political Economy* 110(4): 793-824.

Breeden, Douglas T. 1979. An intertemporal asset pricing model with stochastic consumption and investment opportunities. *Journal of Financial Economics* 7: 265-296.

Breeden, Douglas T. 1986. Consumption, production, inflation and interest rates: a synthesis. *Journal of Financial Economics* 16(1): 3-39.

Breeden, Douglas T. and Robert H. Litzenberger. 1978. Prices of state-contingent claims implicit in options prices. *Journal of Business* 51(4): 621-51.

Breeden, Douglas T., Robert H. Litzenberger, and Tingyan Jia. 2015a. Consumption-based asset pricing, part 1: Classic theory and tests, measurement issues, and limited participation. *Annual Review of Financial Economics* 7: 35-83.

Breeden, Douglas T., Robert H. Litzenberger, and Tingyan Jia. 2015b. Consumption-based asset pricing, part 2: Habit formation, conditional risks, long-run risks, and rare disasters. *Annual Review of Financial Economics* 7: 85-131.

Brennan, Michael J., Eduardo S. Schwartz and Ronald Lagnado. 1997. Strategic asset allocation. *Journal of Economic Dynamics and Control* 21: 1377-1403.

Browning, Martin, Lars Peter Hansen, and James J. Heckman. 1999. Micro data and general equilibrium models. Chapter 8 in John B. Taylor and Michael Woodford, eds. *Handbook of Macroeconomics, Volume 1*. North-Holland, 543-633.

Brunnermeier, Markus K. and Yuliy Sannikov. 2014. A macroeconomic model with a financial sector. *American Economic Review* 104(2): 379-421.

Brunnermeier, Markus K. and Yuliy Sannikov. 2016. The I theory of money. NBER Working Paper 22533.

Brunnermeier, Markus K. and Yuliy Sannikov. 2017. Macro, money and finance: a continuous-time approach, in John B. Taylor and Harald Uhlig, eds., *Handbook of Macroeconomics, Volume 2B* (North-Holland, Amsterdam), 1497-1546.

Caballero, Ricardo J. 2007. *Specificity and the macroeconomics of restructuring*. MIT Press.

Campbell, John Y. 1986. Bond and stock returns in a simple exchange model. *Quarterly Journal of Economics* 101(4): 785-804.

Campbell, John Y. 2003. Consumption-based asset pricing, in G. M. Constantinides, M. Harris, and R. Stulz, *Handbook of the Economics of Finance, Volume 1B* (North-Holland, Amsterdam), 803-887.

Campbell, John Y. 2018. *Financial decisions and markets: a course in asset pricing*. Princeton University Press.

Campbell, John Y. and John H. Cochrane. 1999. By force of habit: a consumption-based explanation of aggregate stock market behavior. *Journal of Political Economy* 107(2): 205-51.

Campbell, John Y. and John H. Cochrane. 2000. Explaining the poor performance of consumption-based asset pricing models. *Journal of Finance*

55(6): 2863-78.

Campbell, John Y. and Luis Viciera. 1999. Consumption and portfolio decisions when expected returns are time varying. *Quarterly Journal of Economics* 114(2): 433-495.

Cassasus, Jaime and Pierre Collin-Dufresne. 2005. Stochastic convenience yield implied from commodity futures and interest rates. *Journal of Finance* 60(5): 2283-331.

Chen, Zhanhui. 2016. Time-to-produce, inventory, and asset prices. *Journal of Financial Economics* 120: 330-345.

Chernozhukov, Victor, Mert Demirer, Esther Duflo and Ivan Fernandez-Val. 2018. Generic machine learning inference on heterogeneous treatment effects in randomized experiments. NBER working paper 24678.

Chetty, Raj. 2006. A new method of estimating risk aversion. *American Economic Review* 96(5): 1821-34.

Chetty, Raj. 2012. Bounds on elasticities with optimization frictions: a synthesis of micro and macro evidence on labor supply. *Econometrica* 80(3): 969-1018.

Cochrane, John H. 1988. Production-based asset pricing. NBER working paper 2776.

Cochrane, John H. 1991. Production-based asset pricing and the link between stock returns and economic fluctuations, *Journal of Finance* 46(1): 209-237.

Cochrane, John H. 1993. Rethinking production under uncertainty. Working paper, University of Chicago Department of Economics, February.

Cochrane, John H. 1996. A cross-sectional test of an investment-based asset pricing model, *Journal of Political Economy* 104(3): 572-621.

Cochrane, John H. 2014. A mean-variance benchmark for intertemporal portfolio theory. *Journal of Finance* 69(1): 1-49.

Cochrane, John H. 2017. Macro-finance. *Review of Finance* 2017: 945-985.

Cohen, Avi J. and G. C. Harcourt. 2003. Whatever happened to the Cambridge capital theory controversies? *Journal of Economic Perspectives* 17(1): 199-214.

Constantinides, George M. 1990. Habit formation: a resolution of the equity premium puzzle. *Journal of Political Economy* 98(3): 519-43.

Constantinides, George M. and Darrell Duffie. 1996. Asset pricing with heterogeneous consumers. *Journal of Political Economy* 104(2): 219-40.

Cont, Rama and Peter Tankov. 2004. *Financial modelling with jump processes*. Chapman and Hall.

Costa, Celso Jose. 2016. *Understanding DSGE models: theory and applications*. Vernon Press.

Cox, John C. and Chi-fu Huang. 1989. Optimal consumption and portfolio policies when asset prices follow a diffusion process. *Journal of Economic Theory* 49: 33-83.

Cox, John C. and Stephen A. Ross. 1976. The valuation of options for alternative stochastic processes. *Journal of Financial Economics* 3: 145-166.

Cox, John C. Jonathan E. Ingersoll, and Stephen A. Ross. 1981. The relation between forward prices and futures prices. *Journal of Financial Economics* 9: 321-346.

Cox, John C. Jonathan E. Ingersoll, and Stephen A. Ross. 1985a. An intertemporal general equilibrium model of asset prices. *Econometrica* 53(2): 363-384.

Cox, John C. Jonathan E. Ingersoll, and Stephen A. Ross. 1985b. A theory of the term structure of interest rates. *Econometrica* 53(2): 385-408.

Croce, Mariano M. 2014. Long-run productivity risk: A new hope for production-based asset pricing? *Journal of Monetary Economics* 66: 13-31.

Cuoco, Domenico. 1997. Optimal consumption and equilibrium prices with portfolio constraints and stochastic income. *Journal of Economic Theory* 72: 33-73.

Cybenko, G. 1989. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems* 2: 303-314.

Dai, Qiang and Kenneth J. Singleton. 2000. Specification analysis of affine term structure models. *Journal of Finance* 55(5): 1943-78.

Davis, Steven J., R. Jason Faberman, and John C. Haltiwanger. 2006. The flow approach to labor markets: new data sources and micro-macro links. *Journal of Economic Perspectives* 20(3): 3-26.

- Davis, Steven J., John C. Haltiwanger, and Scott Schuh. 1996. Job creation and destruction. MIT Press.
- Dawkins, Christina, T. N. Srinivasan, and John Whalley. 2001. Calibration. Chapter 58 in James J. Heckman and Edward Leamer, eds. Handbook of Econometrics, Volume 5. North-Holland, 3653-3703.
- Decker, Ryan, John Haltiwanger, Ron Jarmin, and Javier Miranda. 2014. The role of entrepreneurship in US job creation and economic dynamism. *Journal of Economic Perspectives* 28(3): 3-24.
- Duffee, Gregory R. 2013. Bond pricing and the macroeconomy. Chapter 13 in George Constantinides, Milton Harris, and Rene Stulz, eds. Handbook of the Economics of Finance, Volume 2B. North-Holland, 907-67.
- Duffee, Gregory R. and Richard H. Stanton. 2012. Estimation of dynamic term structure models. *Quarterly Journal of Finance* 2(2): 1250008.
- Duffie, Darrell. 2001. Dynamic asset pricing theory. Princeton University Press.
- Duffie, Darrell and Rui Kan. 1996. A yield-factor model of interest rates. *Mathematical Finance* 6(4): 379-406.
- E, Weinan, Jiequn Han, and Arnulf Jentzen (2017) Deep learning-based numerical methods for high-dimensional parabolic partial differential equations and backward stochastic differential equations. arXiv:1706.04702v1.
- E, Weinan, Chao Ma and Lei Wu. 2019. Barron spaces and the compositional function spaces for neural network models. arXiv:1906.08039v1.
- Eberlein, Ernst and Jan Kallsen. 2019. Mathematical finance. Springer.
- Engle, Robert. 2008. Anticipating correlations: A new paradigm for risk management. Princeton University Press.
- Epstein, Larry G. and Stanley E. Zin. 1989. Substitution, risk aversion, and the temporal behavior of consumption and asset returns: a theoretical framework. *Econometrica* 57(4): 937-69.
- Eraker, Bjorn. 2008. Affine general equilibrium models, *Management Science* 54(12): 2068-2080.
- Eraker, Bjorn and Ivan Shaliastovich. 2008. An equilibrium guide to designing affine pricing models, *Mathematical Finance* 18(4): 519-543.
- Fama, Eugene F. 1970. Multiperiod consumption-investment decisions.

American Economic Review 60(1): 163-74.

Fernandez-Villaverde, Jesus, Juan F. Rubio-Ramirez, and Frank Schorfheide. 2016. Solution and estimation methods for DSGE models. Chapter 9 in John B. Taylor and Harald Uhlig, eds. *Handbook of Macroeconomics*, Volume 2A. North-Holland, 527-724.

Ferson, Wayne E. and George M. Constantinides. 1991. Habit persistence and durability in aggregate consumption: empirical tests. *Journal of Financial Economics* 29(2): 199-240.

Ferson, Wayne E. and Campbell R. Harvey. 1991. The variation of economic risk premiums. *Journal of Political Economy* 99(2): 385-415.

Fisher, Franklin M. 2005. Aggregate production functions – a pervasive, but unpersuasive fairytale. *Eastern Economic Journal* 31(3): 489-491.

Forsyth, Peter and George Labahn. 2007. Numerical methods for controlled Hamilton-Jacobi-Bellman PDEs in finance. *Journal of Computational Finance* 11(2): 1-43.

Friedman, Milton. 1957. *A theory of the consumption function*. NBER and Princeton University Press.

Gali, Jordi. 2015. *Monetary policy, inflation, and the business cycle. An introduction to the New Keynesian framework and its implications*, 2nd ed. Princeton University Press.

Gilboa, I. and D. Schmeidler. 1989. Maxmin expected utility with non unique prior. *Journal of Mathematical Economics* 18(2): 141-153.

Gomes, J. F., L. Kogan and L. Zhang. 2003. Equilibrium cross section of returns. *Journal of Political Economy* 111: 693-732.

Goodfellow, Ian, Yoshua Bengio, and Aaron Courville. 2016. *Deep learning*. MIT Press.

Gullickson, William and Michael J. Harper. 2002. Bias in aggregate productivity trends revisited. *Monthly Labor Review* 125(3): 32-40.

Gurkaynak, Refet S. and Jonathan H. Wright. 2012. Macroeconomics and the term structure. *Journal of Economic Literature* 50(2): 331-367.

Hakansson, Nils H. 1970. Optimal investment and consumption strategies under risk for a class of utility functions. *Econometrica* 38(5): 587-607.

Hammond, P. Brett and Martin L. Leibowitz. 2011. Rethinking the

equity risk premium: an overview and some new ideas, in P. Brett Hammond, Martin L. Leibowitz and Laurence B. Siegel, eds. Rethinking the equity risk premium (Research Foundation of CFA Institute), 1-17.

Han, Jiequn and Weinan E. 2016. Deep learning approximation for stochastic control problems. Neural Information Processing Systems Workshop on Deep Reinforcement Learning. arxiv.org/abs/1611.07422

Han, Jiequn, Arnulf Jentzen, and Weinan E. 2018. Solving high-dimensional partial differential equations using deep learning. *Proceedings of the National Academy of Sciences* 115(34): 8505-8510.

Hansen, Lars Peter. 1987. Calculating asset prices in three example economies. Chapter 6 in Truman Bewley, ed. *Advances in Econometrics Fifth World Congress*, vol. I. Cambridge University Press, 207-43.

Hansen, Lars Peter and Kenneth J. Singleton. 1982. Generalized instrumental variables estimation of nonlinear rational expectations models. *Econometrica* 50(5): 1269-86.

Hansen, Lars Peter and Kenneth J. Singleton. 1983. Stochastic consumption, risk aversion, and the temporal behavior of asset returns. *Journal of Political Economy* 91(2): 249-65.

Hansen, Lars Peter and Kenneth J. Singleton. 1984. Generalized instrumental variables estimation of nonlinear rational expectations models. *Econometrica* 52(1): 267-8.

He, Zhiguo and Arvind Krishnamurthy. 2013. Intermediary asset pricing. *American Economic Review* 103(2): 732-770.

He, Zhiguo and Arvind Krishnamurthy. 2019. A macroeconomic framework for quantifying systemic risk. *American Economic Journal: Macroeconomics* 11(4): 1-37.

Heaton, John and Deborah Lucas. 1992. The effects of incomplete insurance markets and trading costs in a consumption-based asset pricing model. *Journal of Economic Dynamics and Control* 16(3): 601-20.

Higham, Desmond J. 2001. An algorithmic introduction to numerical simulation of stochastic differential equations. *SIAM Review* 43(3): 525-546.

Hopenhayn, Hugo A. 1992. Entry, exit and firm dynamics in long-run equilibrium. *Econometrica* 60(5): 1127-50.

Hornik, K., M. Stinchcombe, and H. White. 1989. Multilayer feedforward networks are universal approximators. *Neural Networks* 2: 359-366.

Hsieh, Chang-Tai and Peter J. Klenow. 2009. Misallocation and manufacturing TFP in China and India. *Quarterly Journal of Economics* 124(4): 1403-48.

Hure, Come. 2019. Numerical methods and deep learning for stochastic control problems and partial differential equations. Doctoral thesis, Université Sorbonne.

Ioffe, Sergey and Christian Szegedy. 2015. Batch normalization: accelerating deep network training by reducing internal covariate shift. In *Proceedings of The 32nd International Conference on Machine Learning (ICML)*, June 2015.

Jermann, Urban. 1998. Asset pricing in production economies. *Journal of Monetary Economics* 41: 257-275.

Jermann, Urban. 2010. The equity premium implied by production. *Journal of Financial Economics* 98: 279-296.

Jermann, Urban. 2013. A production-based model for the term structure. *Journal of Financial Economics* 109: 293-306.

Jiang, Daniel R. and Warren B. Powell. 2015. An approximate dynamic programming algorithm for monotone value functions. *Operations Research* 63(6): 1489-1511.

Jorgenson, Dale W. and Zvi Griliches. 1967. The explanation of productivity change. *Review of Economic Studies* 34(99): 249-280.

Jorgenson, Dale W., Frank M. Gollop, and Barbara M. Fraumeni. 1987. *Productivity and U.S. economic growth*. Harvard University Press.

Jorgenson, Dale W., Mun S. Ho and Kevin Stiroh. 2005. *Information technology and the American growth resurgence*. MIT Press.

Judd, Kenneth L. 1998. *Numerical methods in economics*. MIT Press.

Kaji, Tetsuya, Elena Manresa and Guillaume Pouliot. 2020. An adversarial approach to structural estimation. arXiv preprint, July.

Kaltenbrunner, G. and L. A. Lochstoer. 2010. Long-run risk through consumption smoothing. *Review of Financial Studies* 23: 3141-3189.

Knight, Frank H. 1921. *Risk, uncertainty and profit*. Houghton Mifflin.

- Kim, Tong Suk and Edward Omberg. 1996. Dynamic nonmyopic portfolio behavior. *Review of Financial Studies* 9: 141-61.
- King, Robert G., Charles I. Plosser, and Sergio T. Rebelo. 1988a. Production, growth and business cycles I: The basic neoclassical model. *Journal of Monetary Economics* 21: 195-232.
- King, Robert G., Charles I. Plosser, and Sergio T. Rebelo. 1988b. Production, growth and business cycles II: New directions. *Journal of Monetary Economics* 21: 309-341.
- King, Robert G., Charles I. Plosser, and Sergio T. Rebelo. 2002. Production, growth and business cycles: Technical appendix. *Computational Economics* 20: 87-116.
- Kloeden, Peter E. and Eckhard Platen. 1999. Numerical solution of stochastic differential equations. Springer-Verlag.
- Kocherlakota, Narayana R. 1996. The equity risk premium: it's still a puzzle. *Journal of Economic Literature* 34: 42-71.
- Koesler, Simon and Michael Schymura. 2015. Substitution elasticities in a constant elasticity of substitution framework—empirical estimates using nonlinear least squares. *Economic Systems Research* 27(1): 101-121.
- Kogan, Leonid and Dimitris Papanikolaou. 2014. Growth opportunities, technology shocks, and asset prices. *Journal of Finance* 69(2): 675-718.
- Kogan, Leonid and Dimitris Papanikolaou. 2018. Equilibrium analysis of asset prices: Lessons from CIR and APT. *Journal of Portfolio Management* 44(6): 59-69.
- Kreps, David M. and Evan L. Porteus. 1978. Temporal resolution of uncertainty and dynamic choice theory. *Econometrica* 46(1): 185-200.
- Kung, Howard. 2015. Macroeconomic linkages between monetary policy and the term structure of interest rates. *Journal of Financial Economics* 115: 42-57.
- Kushner, Harold and Paul Dupuis. 2002. Numerical methods for stochastic control in continuous time. Springer-Verlag.
- Kydland, Finn and Edward Prescott. 1982. Time to build and aggregate fluctuations. *Econometrica* 50: 1345-1370.
- Lettau, Martin and Sydney C. Ludvigson. 2001a. Consumption, aggre-

- gate wealth, and expected stock returns. *Journal of Finance* 56(3): 815-49.
- Lettau, Martin and Sydney C. Ludvigson. 2001b. Resurrecting the (C)CAPM: a cross-sectional test when risk premia are time-varying. *Journal of Political Economy* 109(6): 1238-87.
- LeVeque, Randall. 2007. Finite difference methods for ordinary and partial differential equations. SIAM.
- Liang, Shiyu and R. Srikant. 2017. Why deep neural networks for function approximation? *Proceedings of the 5th International Conference on Learning Representations*.
- Lintner, John. 1965. The valuation of risk assets and the selection of risky investments in stock portfolios and capital budgets. *Review of Economics and Statistics* 47(1): 13-37.
- Liu, Jun. 2007. Portfolio selection in stochastic environments. *Review of Financial Studies* 20: 1-39.
- Ljungqvist, Lars and Thomas Sargent. 2018. Recursive macroeconomic theory, 4th ed. MIT Press.
- Lucas, Robert E. 1978. Asset prices in an exchange economy. *Econometrica* 46(6): 1429-1445.
- Ludvigson, Sydney C. 2013. Advances in consumption-based asset pricing: empirical tests. In George M. Constantinides, M. Harris, and Rene M. Stulz, eds., *Handbook of the Economics of Finance*, vol. 2, 799-906. Amsterdam: North-Holland.
- Mankiw, N. Gregory and Stephen P. Zeldes. 1991. The consumption of stockholders and nonstockholders. *Journal of Financial Economics* 29(1): 97-112.
- Markowitz, Harry. 1952. Portfolio selection. *Journal of Finance* 7(1): 77-91.
- Mehra, Rajnish. 2012. Consumption-based asset pricing models. *Annual Review of Financial Economics* 4: 385-409.
- Mehra, Rajnish and Edward C. Prescott. 1985. The equity premium: a puzzle. *Journal of Monetary Economics* 15: 145-161.
- Merton, Robert C. 1969. Lifetime portfolio selection under uncertainty: the continuous-time case. *Review of Economics and Statistics* 51(3): 247-257.

- Merton, Robert C. 1971. Optimal consumption and portfolio rules in a continuous-time model. *Journal of Economic Theory* 3(4): 373-413.
- Merton, Robert C. 1973a. An intertemporal capital asset pricing model, *Econometrica* 41: 867-887.
- Merton, Robert C. 1973b. Theory of rational option pricing. *Bell Journal of Economics and Management Science* 4(1): 141-183.
- Merton, Robert C. 1974. On the pricing of corporate debt: the risk structure of interest rates. *Journal of Finance* 449-470.
- Miao, Jianjun. 2014. *Economic dynamics in discrete time*. MIT Press.
- Mossin, Jan. 1966. Equilibrium in a capital asset market. *Econometrica* 34(4): 768-783.
- Oksendal, Bernt and Agnes Sulem. 2019. *Applied stochastic control of jump diffusions*, 3rd ed. Springer-Verlag.
- Papanikolaou, Dimitris. 2011. Investment shocks and asset prices. *Journal of Political Economy* 119(4): 639-685.
- Parker, Jonathan A. and Christian Julliard. 2005. Consumption risk and the cross section of expected returns. *Journal of Political Economy* 113(1): 185-222.
- Petkova, Ralitsa and Lu Zhang. 2005. Is value riskier than growth? *Journal of Financial Economics* 78(1): 187-202.
- Piketty, Thomas. 2014. *Capital in the 21st century*. Harvard-Belknap.
- Pratt, John W. 1964. Risk aversion in the small and the large. *Econometrica* 32(1/2): 122-136.
- Raissi, M., P. Perdikaris, and G. E. Karniadakis. 2019. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics* 378: 686-707.
- Reichling, Felix and Charles Whalen. 2012. Review of estimates of the Frisch elasticity of labor supply. Congressional Budget Office Working Paper Series 2012-13.
- Rietz, Thomas A. 1988. The equity risk premium: a solution. *Journal of Monetary Economics* 22: 117-131.
- Rogers, L. C. G. 2013. *Optimal investment*. Springer-Verlag.

- Romer, David. 2016. The trouble with macroeconomics. Manuscript.
- Ross, Stephen A. 1976. The arbitrage theory of capital asset pricing. *Journal of Economic Theory* 13: 341-360.
- Rubinstein, Mark. 1976. The valuation of uncertain income streams and the pricing of options. *Bell Journal of Economics and Management Science* 7(2): 407-25.
- Rudebusch, Glenn D. 2010. Macro-finance models of interest rates and the economy. *The Manchester School (Supplement 2010)*: 25-52.
- Rudebusch, Glenn D. and Eric T. Swanson. 2008a. The bond premium in a DSGE model with long-run real and nominal risks. Federal Reserve Bank of San Francisco Working Paper 2008-31.
- Rudebusch, Glenn D. and Eric T. Swanson. 2008b. Examining the bond premium puzzle with a DSGE model. *Journal of Monetary Economics* 55: S111-26.
- Samuelson, Paul A. 1969. Lifetime portfolio selection by dynamic stochastic programming. *Review of Economics and Statistics* 51(3): 239-46.
- Sangvinatsos, Antonios and Jessica Wachter. 2005. Does the failure of the expectations hypothesis matter for long-term investors? *Journal of Finance* 60: 179-230.
- Santos, Tano and Pietro Veronesi. 2006. Labor income and predictable stock returns. *Review of Financial Studies* 19(1): 1-44.
- Sato, K. 1967. A two-level constant-elasticity-of-substitution production function. *Review of Economic Studies* 34(2): 201-218.
- Schwartz, Eduardo S. 1997. The stochastic behavior of commodity prices: implications for valuation and hedging. *Journal of Finance* 52(3): 923-973.
- Sethi, Suresh P. 1997. *Optimal consumption and investment with bankruptcy*. Kluwer Academic Publishers.
- Sharpe, William F. 1964. Capital asset prices: A theory of market equilibrium under conditions of risk. *Journal of Finance* 19(3): 425-442.
- Sirignano, Justin and Konstantinos Spiliopoulos. 2018. DGM: A deep learning algorithm for solving partial differential equations. arXiv preprint.
- Staum, Jeremy. 2008. Incomplete markets. Chapter 12 in John R. Birge and Vadim Linetsky, eds. *Handbooks in Operations Research and Manage-*

- ment Science, Volume 15: Financial Engineering. North-Holland: 511-63.
- Stigler, George J. and Gary S. Becker. 1977. De gustibus non est disputandum. *American Economic Review* 67(2): 76-90.
- Syverson, Chad. 2011. What determines productivity? *Journal of Economic Literature* 49(2): 326-365.
- Tobin, James. 1969. A general equilibrium approach to monetary theory. *Journal of Money, Credit and Banking* 1(1): 15-29.
- Torres, Jose L. 2015. Introduction to dynamic macroeconomic general equilibrium models. Vernon Press.
- Uhlig, Harald. 1999. A toolkit for analyzing nonlinear dynamic stochastic models easily. In Ramon Marimon and Andrew Scott, eds. *Computational methods for the study of dynamic economies*. Oxford University Press: 30-61.
- Valentinyi, Akos and Berthold Herrendorf. 2008. Measuring factor income shares at the sectoral level. *Review of Economic Dynamics* 11: 820-835.
- Villemot, Sebastien. 2011. Solving rational expectations models at first order: what Dynare does. Dynare Working Paper 2, CEPREMAP.
- Weil, Philippe. 1989. The equity premium puzzle and the risk-free rate puzzle. *Journal of Monetary Economics* 24(3): 401-21.
- Wicksell, Knut. 1935. *Lectures in political economy, volume II: Money*. Routledge.
- Woodford, Michael. 2003. *Interest and prices: foundations of a theory of monetary policy*. Princeton University Press.
- Wu, Kailiang and Dongbin Xiu. 2020. Data-driven deep learning of partial differential equations in modal space. *Journal of Computational Physics* 408: 109307.
- Yong, Jiongmin and Xun Yu Zhou. 1999. *Stochastic controls: Hamiltonian systems and HJB equations*. Springer-Verlag.
- Zariphopoulou, Thaleia. 1994. Consumption-investment models with constraints. *SIAM Journal of Control and Optimization* 32(1): 59-85.
- Zhang, L. 2005. The value premium. *Journal of Finance* 60: 67-103.