



Columbia University

*Department of Economics
Discussion Paper Series*

**Macroeconomic Analysis Without the Rational Expectations
Hypothesis**

Michael Woodford

Discussion Paper No.: 1314-10

*Department of Economics
Columbia University
New York, NY 10027*

October 2013

Macroeconomic Analysis without the Rational Expectations Hypothesis*

Michael Woodford
Columbia University

August 16, 2013

Abstract

This paper reviews a variety of alternative approaches to the specification of the expectations of economic decisionmakers in dynamic models, and reconsiders familiar results in the theory of monetary and fiscal policy when one allows for departures from the hypothesis of rational expectations. The various approaches are all illustrated in the context of a common model, a log-linearized New Keynesian model in which both households and firms solve infinite-horizon decision problems; under the hypothesis of rational expectations, the model reduces to the standard “3-equation model” used in studies such as Clarida *et al.* (1999). The alternative approaches considered include rationalizable equilibrium dynamics (Guesnerie, 2008); restricted perceptions equilibria (Branch, 2004); decreasing-gain and constant-gain variants of least-squares learning dynamics (Evans and Honkapohja, 2001); rational belief equilibria (Kurz, 2012); and near-rational expectations equilibria (Woodford, 2010). Issues treated include Ricardian equivalence; the determinacy of equilibrium under alternative interest-rate rules; non-fundamental sources of aggregate instability; the trade-off between inflation stabilization and output-gap stabilization; and the possibility of a “deflation trap.”

*Prepared for *Annual Review of Economics*, volume 5. I would like to thank Klaus Adam, Ben Hebert, Mordecai Kurz, David Laibson, and Bruce Preston for helpful discussions, Savitar Sundaresan for research assistance, and the Institute for New Economic Thinking and the Taussig Visiting Professorship, Harvard University, for supporting this research.

A crucial methodological question in macroeconomic analysis is the way in which the expectations of decisionmakers about future conditions should be modeled. To the extent that behavior is modeled as goal-directed, it will depend (except in the most trivial cases) on expectations; and analyses of the effects of alternative governmental policies need to consider how expectations are endogenously influenced by one policy or another. Finding tractable ways to address this issue has been a key challenge for the extension of optimization-based economic analysis to the kinds of dynamic settings required for most questions of interest in macroeconomics.

The dominant approach for the past several decades, of course, has made use of the hypothesis of model-consistent or “rational expectations” (RE): the assumption that people have probability beliefs that coincide with the probabilities predicted by one’s model. The RE benchmark is a natural one to consider, and its use has allowed a tremendous increase in the sophistication of the analysis of dynamics in the theoretical literature in macroeconomics. Nonetheless, the assumption is a strong one, and one may wonder if it should be relaxed, especially when considering relatively short-run responses to disturbances, or the consequences of newly adopted policies that have not been followed in the past — both of which are precisely the types of situations which macroeconomic analysis frequently seeks to address.

While the assumption that an economy’s dynamics must necessarily correspond to an RE equilibrium may seem unjustifiably strong — and under some circumstances, is a heroic assumption indeed — it does not follow that we should then be equally willing to entertain all possible assumptions about the expectations of economic agents. It makes sense to assume that expectations should not be completely arbitrary, and have no relation to the kind of world in which the agents live; indeed, it is appealing to assume that people’s beliefs should be *rational*, in the ordinary-language sense, though there is a large step from this to the RE hypothesis.¹ We should like, therefore, to replace the RE hypothesis by some weaker restriction, that nonetheless implies a substantial degree of conformity between people’s beliefs and reality — that implies, at the least, that people do not make *obvious* mistakes.

The literature has explored two broad approaches to the formulation of a criterion for reasonableness of beliefs that is weaker than the RE hypothesis. One is to assume that people should correctly *understand the economic model*, and be able to form correct inferences from it about possible future outcomes. The other is to assume that that the probabilities that people assign to possible future outcomes should not be too

¹This is stressed, for example, by Kurz (2012, p. 1).

different from the probabilities with which different outcomes *actually occur*, given that experience should allow some acquaintance with these probabilities, whether or not people understand how these outcomes are generated. The former approach supposes that beliefs should be refined through a process of reflection, independent of experience and not necessarily occurring in real time, that Guesnerie (1992) calls *eduction*, while the latter supposes that beliefs should be refined over time through a process of *induction* from observed outcomes. Section 2 discusses the first approach, while examples of the second approach are taken up in sections 3 and 4.

Within the category of inductive approaches, one may distinguish two important sub-categories. A first class of approaches specifies the beliefs that should be regarded as reasonable by specifying the patterns that people should be able to recognize in the data on the basis of the rationality of the *procedure* used to look for such patterns. A different class of approaches specifies a degree of *correspondence* between subjective and model-implied probabilities that should be expected, without explicitly modeling the process of inference through which such beliefs are formed. The first class of approaches (models of econometric learning) is treated in section 3, while the second class (theories of partially or approximately correct beliefs) is taken up in section 4.

The different possible approaches to the specification of expectations are compared by illustrating the application of each in the context of the same general framework for macroeconomic analysis, introduced in section 1. In each case, it is shown that one can demand that the specification of beliefs satisfy quite stringent rationality requirements without, in general, being able to conclude that the predictions of the RE equilibrium analysis must obtain. I shall consider in particular the consequences of alternative specifications of expectations for several familiar issues: the conditions under which an interest-rate rule for monetary policy should be able to maintain a stable inflation rate; the nature of the trade-off between inflation stabilization and output stabilization; and the effects of the government budget on aggregate demand.

1 Temporary Equilibrium in a New Keynesian Model

I begin by introducing a framework for analysis of the determinants of aggregate output and inflation, in which subjective expectations can be specified arbitrarily, in order to clarify the role of alternative specifications of expectations. It allows the effects of both monetary and fiscal policies to be considered, along with a variety of

types of exogenous disturbances to economic “fundamentals” and possible shifts in expectations. The model is one in which households and firms solve infinite-horizon optimization problems, as in the DSGE models commonly used for quantitative policy analysis; in fact, under the assumption of rational expectations, the model presented here corresponds to a textbook New Keynesian model of the kind analyzed in Clarida *et al.* (1999), Woodford (2003, chap. 4), Gali (2008, chap. 3), or Walsh (2010, chap. 8). Essentially, the goal of this section is to show how “temporary equilibrium” analysis of the kind introduced by Hicks (1939) and Lindahl (1939) and further developed by Grandmont (1977, 1988) — in which a general competitive equilibrium is defined at each point in time, on the basis of the (independently specified) expectations that decisionmakers happen to entertain at that time — can be extended to a setting with monopolistic competition, sticky prices, and infinite-horizon planning, for closer comparison with the conclusions of conventional macroeconomic analysis. Some of the best-known conclusions from RE analysis of the model are also recalled, as a basis for comparison with the conclusions from alternative specifications of expectations in the later sections.

1.1 Expectations and Aggregate Demand

The economy is made up of a large number of identical households, and a large number of firms, each of which is the monopoly producer of a particular differentiated good, with each household owning an equal share of each firm.² At any point in time, a household has an infinite-horizon consumption (and wealth-accumulation) plan from that date forward, which maximizes expected discounted utility under certain *subjective* expectations regarding the future evolution of income and the rate of return on saving; and the household’s expenditure at that date is assumed to be the one called for by the plan believed to be optimal at that time. (The household may or may not continue to execute the same plan as time passes, depending on what is assumed about how expectations change.) While I wish for now to leave the subjective expectations unspecified, the expectations held at any date represent a well-behaved probability measure over possible future evolutions of the state variables. Because I

²These and other aspects of the model structure are explained in more detail in Woodford (2003). Here I focus only on the points at which an alternative model of expectations requires a generalization of the standard exposition.

shall not assume that subjective expectations are necessarily model-consistent, they are not necessarily the same across households; nor shall I necessarily assume that a household's later expectations are those that it previously expected to hold.

To simplify, I shall assume that the only traded asset is riskless nominal one-period government debt.³ I shall further assume that households have no choice but to supply the hours of work that are demanded by firms, at a wage that is fixed by a union that bargains on behalf of the households. A household then has a single decision each period, which is the amount to spend on consumption. Because its nonfinancial income (the sum of its wage income and its share of the profits of the firms) is outside its control, in order to analyze optimal expenditure, we need only specify households' expectations regarding the time path of their total nonfinancial income. If an equal amount of work is demanded from each household at the wage fixed by the union, and households similarly each receive an equal share of the profits of all of the firms,⁴ then each household's nonfinancial income will be the same each period, and equal to its share of the total value of output in that period; hence we can equivalently specify nonfinancial income expectations as expectations regarding the path of output.

A household's perceived intertemporal budget constraint then depends only on its expectations about the path of aggregate output, the path of aggregate tax collections (also assumed to be levied equally on each household), and the real return on the one-period debt. The consumption planning problem for an individual household at a given point in time is then the familiar one faced by a household with rational expectations and an exogenously given income process (discussed, for example, in Deaton, 1992), and the solution is correspondingly the same, except with subjective probabilities substituted for objective ones.⁵

³In many RE analyses, with a representative household and fiscal policy assumed to be Ricardian (in the sense defined in Woodford, 2001), the model dynamics are unaffected by allowing additional financial markets or even issuance of other types of government debt. The restriction to one-period debt is no longer innocuous, however, once one allows for more general hypotheses regarding expectations, as shown for example in Eusepi and Preston (2012b). Nonetheless, I here consider only the most simple case.

⁴Equity ownership shares are assumed for simplicity to be non-tradeable.

⁵Note that the log-linear theory of aggregate demand derived here is the same as in a model where households are assumed simply to have an exogenous endowment of the consumption good, like that of Guesnerie (2008).

I shall log-linearize this and other structural relations of the model around a deterministic steady state, in which (i) all exogenous state variables are forever constant, (ii) monetary and fiscal policy are specified to maintain a constant zero rate of inflation and some constant positive level of public debt, and (iii) all subjective expectations are correct (*i.e.*, households and firms have perfect foresight). The log-linearized relations accordingly represent an approximation to the exact model, applying in the case that exogenous disturbances are sufficiently small, monetary and fiscal policies are sufficiently close to being consistent with this steady state, and expectational errors are sufficiently small.

A log-linear approximation to the consumption function takes the form

$$c_t^i = (1 - \beta)b_t^i + \sum_{T=t}^{\infty} \beta^{T-t} \hat{E}_t^i \{ (1 - \beta)(Y_T - \tau_T) - \beta\sigma(i_T - \pi_{T+1}) + (1 - \beta)s_b(\beta i_T - \pi_T) - \beta(\bar{c}_{T+1} - \bar{c}_T) \}. \quad (1.1)$$

Here c_t^i is total real expenditure by household i (on all of the differentiated goods) in period t , measured as a deviation from the steady-state level of consumption, and expressed as a fraction of steady-state output; b_t^i is the value of maturing bonds carried into period t by household i , deflated by the *period* $t - 1$ price level (rather than the period t price level) so that this wealth measure is predetermined at date $t - 1$;⁶ Y_t is the deviation of aggregate output from its steady-state value, as a fraction of that steady-state value; τ_t is net tax collections, also measured as a deviation from the steady-state level of tax revenues and expressed as a fraction of steady-state output; i_t is the riskless one-period nominal interest rate on debt issued in period t ; π_t is the inflation rate between periods $t - 1$ and t ; and \bar{c}_t is an exogenous shock to the utility from consumption in period t . In addition, $0 < \beta < 1$ is the utility discount factor, $\sigma > 0$ measures the intertemporal elasticity of substitution, and $s_b > 0$ is the steady-state level of government debt, expressed as a fraction of steady-state output.

The notation $\hat{E}_t^i\{\cdot\}$ indicates the expected value of the terms in the brackets, under the subjective expectations of household i in period t . The b_t^i and $Y_T - \tau_T$ terms then represent (a subjective version of) the usual permanent-income hypothesis;⁷ the

⁶The reduction in the real value of this wealth due to inflation between periods $t - 1$ and t is then reflected in the $-s_b\pi_t$ term inside the curly brackets.

⁷See, *e.g.*, Deaton (1992) for an exposition of the standard theory, under the assumptions of rational expectations, no preference shocks, and a constant real interest rate.

$\sigma(i_T - \pi_{T+1})$ terms indicate how expectations of a real interest rate different from the rate of time preference shift the desired time path of spending; the $s_b(\beta i_T - \pi_T)$ terms indicate the income effects of variations in nominal interest rates and inflation; and the \bar{c}_T terms indicate the effects of preference shocks on the desired time path of spending.

Equation (1.1) involves subjective expectations of a number of variables at many future horizons, but under our linear approximation we can write desired expenditure as a function of the household's forecast of a single variable,

$$c_t^i = (1 - \beta)b_t^i + (1 - \beta)(Y_t - \tau_t) - \beta[\sigma - (1 - \beta)s_b]i_t - (1 - \beta)s_b\pi_t + \beta\bar{c}_t + \beta\hat{E}_t^i v_{t+1}^i, \quad (1.2)$$

where the composite variable v_t^i is defined as

$$v_t^i \equiv \sum_{T=t}^{\infty} \beta^{T-t} \hat{E}_t^i \{ (1 - \beta)(Y_T - \tau_T) - [\sigma - (1 - \beta)s_b](\beta i_T - \pi_T) - (1 - \beta)\bar{c}_T \}. \quad (1.3)$$

The advantage of this notation is that we need only to specify how people forecast a single variable each period; note however that the variable that people must forecast is a *subjective* state variable, that will depend on their own future forecasts.

Aggregate demand is then given by $Y_t = \int c_t^i di + G_t$, where G_t is the departure of government purchases of the composite good from their steady-state level, also measured as a fraction of steady-state output. Substituting (1.2) for c_t^i in this expression, we obtain

$$Y_t = g_t + (1 - \beta)b_t + v_t - \sigma\pi_t, \quad (1.4)$$

where $g_t \equiv G_t + \bar{c}_t$ is a composite exogenous disturbance to “autonomous expenditure”,⁸ $b_t \equiv \int b_t^i di$ is the aggregate supply of public debt, and $v_t \equiv \int v_t^i di$ is the average value of the expectational variable defined in (1.3). We thus obtain an equation for aggregate demand that separates out the effects of the exogenous disturbances g_t , the wealth effect of public debt, and the average state of expectation captured by v_t .

The government's flow budget constraint implies a law of motion

$$b_{t+1} = \beta^{-1} [b_t - s_b\pi_t - s_t] + s_b i_t \quad (1.5)$$

⁸To simplify, I treat government purchases as an exogenous disturbance, rather than as a possibly endogenous policy choice.

for the public debt, where $s_t \equiv \tau_t - G_t$ is the real primary budget surplus in period t , measured as a deviation from its steady-state level and expressed as a fraction of steady-state output. The *aggregate demand block* of our model then consists of equations (1.4) and (1.5), together with monetary and fiscal policy equations that specify the evolution of i_t and s_t respectively (generally as a function of other endogenous variables), and a specification of the evolution of the forecasts $\{\hat{E}_t^i v_{t+1}^i\}$ (which determine the expectational variable v_t). We then have a system of four equations per period (plus the specification of expectations) to determine the paths $\{Y_t, b_t, i_t, s_t\}$ given a path for the price level, or the paths $\{\pi_t, b_t, i_t, s_t\}$ given a supply-determined path for output, along with the composite disturbance g_t and shocks to policy and expectations.

Definition (1.3) has a recursive form, so that the only subjective expectations involved in v_t^i are i 's forecast of the corresponding variable one period in the future. Specifically, (1.3) implies that

$$v_t^i = (1 - \beta)v_t + \beta(1 - \beta)(b_{t+1} - b_t) - \beta\sigma(i_t - \pi_t) + \beta\hat{E}_t^i v_{t+1}^i, \quad (1.6)$$

where in addition to substituting the forecast of v_{t+1}^i for the expectational terms, I have used (1.4) to substitute out Y_t and (1.5) to substitute out τ_t . Hence to specify the aspects of subjective expectations that are relevant for aggregate demand determination it suffices that we specify an evolution for the $\{v_t^i\}$, and subjective one-period-ahead forecasts of those variables, that are consistent with (1.6). This result is used below to characterize the possible temporary equilibrium (TE) dynamics, under various more specific assumptions about expectations.

1.2 Ricardian Expectations

Equation (1.5) implies that non-explosive dynamics for the real public debt require that

$$b_t = \sum_{T=t}^{\infty} \beta^{T-t} [s_T - s_b(\beta i_T - \pi_T)].$$

I shall say that households have *Ricardian expectations* if they expect that the path of primary surpluses will necessarily satisfy this condition, so that

$$\hat{E}_t^i \sum_{T=t}^{\infty} \beta^{T-t} [s_T - s_b(\beta i_T - \pi_T)] = b_t \quad (1.7)$$

at all times. It is not obvious that expectations must have this property, even under the RE hypothesis;⁹ it is even less obvious that reasonable expectations must have this property once one dispenses with the strict RE assumption. Nonetheless, this property is frequently assumed (at least tacitly) in non-RE analyses, and I shall mainly assume it in the discussion below, to simplify the analysis.¹⁰

Under this assumption, there is no longer a wealth effect of variation in the public debt; the b_t term in (1.4) is exactly canceled by the effects of the change in the expected path of primary surpluses on the v_t term. In fact, one can write (1.4) more simply as

$$Y_t = g_t + \bar{v}_t - \sigma\pi_t, \quad (1.8)$$

where $\bar{v}_t \equiv v_t + (1 - \beta)b_t$ is the average value of a subjective variable \bar{v}_t^i which (under the hypothesis of Ricardian expectations) can be defined simply as

$$\bar{v}_t^i \equiv \sum_{T=t}^{\infty} \beta^{T-t} \hat{E}_t^i \{ (1 - \beta)(Y_T - g_T) - \sigma(\beta i_T - \pi_T) \}. \quad (1.9)$$

In this case, aggregate demand determination is completely independent of the paths of both public debt and tax collections — so that “*Ricardian equivalence*” obtains — as long as these fiscal variables have no direct effect on people’s expectations regarding the evolution of the variables $\{Y_t, \pi_t, i_t, g_t\}$; or more simply, as long they have no direct effect on forecasts of the path of the variables $\{v_t^i\}$.¹¹

Under the assumption of Ricardian expectations, it is convenient to specify expectations by describing the evolution of the variables $\{\bar{v}_t^i\}$; these must be consistent

⁹It is possible for people to believe in a fiscal rule without this property, and yet for a RE equilibrium to exist, as shown for example in Woodford (2001). (In equilibrium, debt does not explode, and expectations are correct; but people believe that debt *would* explode in the case of certain paths for endogenous variables that do not occur in equilibrium.) Some have disputed whether such a specification of out-of-equilibrium beliefs should be considered to be consistent with the RE hypothesis; see Bassetto (2002) for a careful discussion.

¹⁰In many NK models with adaptive learning, TE relations are derived by simply inserting subjective expectations into equations that hold (in terms of model-consistent conditional expectations) in the RE version of the model; the fact that government liabilities are not perceived to be net wealth in the RE analysis then leads to an assumption that they are not in the learning analysis, without any discussion of the assumption. Evans and Honkapohja (2010) instead make the assumption of Ricardian expectations explicit. Eusepi and Preston (2012a, 2012b) and Benhabib *et al.* (2012) provide examples of analyses in which Ricardian expectations are not assumed.

¹¹See Evans *et al.* (2012) for a more detailed discussion of the conditions under which Ricardian equivalence obtains even without the RE hypothesis.

with a relation of the form

$$\bar{v}_t^i = (1 - \beta)\bar{v}_t - \beta\sigma(i_t - \pi_t) + \beta\hat{E}_t^i \bar{v}_{t+1}^i, \quad (1.10)$$

analogous to (1.6). The complete *aggregate demand block* of the model then consists of equation (1.8), a monetary policy equation, and a specification of the evolution of the expectational variable $\{\bar{v}_t^i\}$ consistent with (1.10). This system provides two equations per period to determine the paths of $\{Y_t, i_t\}$ given the evolution of the price level, the exogenous disturbances $\{g_t\}$, and shocks to policy and expectations.

1.3 Expectations and Aggregate Supply

Each of the monopolistically competitive firms sets the price for the particular good that it alone produces. Prices are assumed to remain fixed for a random interval of time: each period, fraction $0 < \alpha < 1$ of all goods prices remain the same (in monetary terms) as in the previous period, while the other prices are reconsidered; and the probability that any given price is reconsidered in any period is assumed to be independent of both the current price and the length of time that the price has remained fixed. It then follows that (again in a log-linear approximation) the rate of inflation π_t between periods $t - 1$ and t will be given by

$$\pi_t = (1 - \alpha)p_t^*, \quad (1.11)$$

where for each firm j , p_t^{*j} is the amount by which the firm would choose to set the log price of its good higher than p_{t-1} , the log of the general price index in period $t - 1$, were it to be one of the firms that reconsiders its price in period t ; and $p_t^* \equiv \int p_t^{*j} dj$ is the average value of this quantity across all firms.

Each firm that reconsiders its price in a given period chooses the new price that it believes will maximize the present value of its profits from that period onward. In a Dixit-Stiglitz model of monopolistic competition, with a single economy-wide labor market, profits in any period are the same function of a firm's own price and of aggregate market conditions for each firm. The single-period profit-maximizing log price p_t^{opt} is then the same for each firm, and a log-linear approximation to the first-order condition for optimal price-setting takes the form

$$p_t^{*j} = (1 - \alpha\beta) \sum_{T=t}^{\infty} (\alpha\beta)^{T-t} \hat{E}_t^j p_T^{opt} - p_{t-1}, \quad (1.12)$$

where β is again the utility discount factor (also the rate at which real profits are discounted in steady state), and $\hat{E}_t^j[\cdot]$ indicates a conditional expectation with respect to the subjective beliefs of firm j at date t . The recursive form of (1.12) implies that internally consistent expectations on the part of any firm must satisfy

$$p_t^{*j} = (1 - \alpha\beta)(p_t^{opt} - p_{t-1}) + \alpha\beta(\hat{E}_t^j p_{t+1}^{*j} + \pi_t). \quad (1.13)$$

Averaging this expression over firms j and using the resulting equation to substitute for p_t^* in (1.11), we see that inflation determination depends only on the average of firms' subjective forecasts of a single variable, namely each firm's own value for the expectational variable p_T^{*j} one period in the future.

Suppose furthermore that the union sets a wage each period with the property that at that wage, a marginal increase in labor demand would neither increase nor decrease average perceived utility across households, if for each household the marginal utility of the additional wage income is weighed against the marginal disutility of the additional work. This implies that (in a log-linear approximation)

$$w_t = \nu_t - \lambda_t,$$

where w_t is the log real wage, ν_t is the log of the (common) marginal disutility of labor, and λ_t is the average across households of λ_t^i , a household's subjective assessment of its marginal utility of additional real income. Since optimizing consumption demand implies that

$$\lambda_t^i = -\sigma^{-1}(c_t^i - \bar{c}_t),$$

we obtain

$$w_t = \nu_t + \sigma^{-1}(c_t - \bar{c}_t) = \nu_t + \sigma^{-1}(Y_t - g_t),$$

just as in a representative-household model with RE and a competitive economy-wide labor market. In a model in which labor is the only variable factor of production, both ν_t and the marginal product of labor can be expressed as functions of hours worked and hence as functions of output Y_t and the exogenous level of productivity in period t . We then obtain a relation of the form

$$p_t^{opt} = p_t + \xi(Y_t - Y_t^n) + \mu_t, \quad (1.14)$$

where y_t indicates the “output gap” (defined as $Y_t - Y_t^n$, where the “natural rate” of output Y_t^n is a composite exogenous disturbance, involving variations in g_t , produc-

tivity, and shocks to the disutility of labor), and μ_t measures exogenous variations in the desired markup of prices over marginal cost.¹²

Substituting (1.11) for π_t and (1.14) for p_t^{opt} in (1.13), we obtain

$$p_t^{*j} = (1 - \alpha)p_t^* + (1 - \alpha\beta) [\xi y_t + \mu_t] + \alpha\beta \hat{E}_t^j p_{t+1}^{*j}. \quad (1.15)$$

In order for beliefs to be internally consistent, the evolution of the expectational variables $\{p_t^{*j}\}$ and subjective one-period-ahead forecasts of those variables must satisfy (1.15) at all times. Given such beliefs, the inflation rate will be determined by (1.11). Thus equations (1.11) and (1.15) constitute the *aggregate-supply block* of the model, which determines the evolution of the general price level, given the evolution of output, exogenous disturbances, and expectations.

1.4 The Complete Model

In the complete TE system, then, the expectations that must be specified are paths for $\{\bar{v}_t^i\}$ for all households and $\{p_t^{*j}\}$ for all firms, consistent with relations (1.10) and (1.15) respectively. Given these expectations, paths for the exogenous disturbances $\{g_t, Y_t^n, \mu_t\}$, and a monetary policy rule for the evolution of $\{i_t\}$, the evolution of aggregate output and inflation are then given by equations (1.8) and (1.11).

Substituting for the variables v_t^i and p_t^{*j} in terms of current observables and forecasts of future conditions, (1.8) can be written as

$$Y_t = g_t - \sigma i_t + \int \hat{E}_t^i \bar{v}_{t+1}^i di, \quad (1.16)$$

and (1.11) can correspondingly be written as

$$\pi_t = \kappa y_t + u_t + (1 - \alpha)\beta \int \hat{E}_t^j p_{t+1}^{*j} dj, \quad (1.17)$$

where

$$\kappa \equiv \frac{(1 - \alpha)(1 - \alpha\beta)}{\alpha} \xi > 0, \quad u_t \equiv \frac{(1 - \alpha)(1 - \alpha\beta)}{\alpha} \mu_t.$$

We thus obtain an “IS equation” and “AS equation” to describe short-run output and inflation determination, given monetary policy, exogenous disturbances, and expectations.

¹²The separation of the effects of exogenous disturbances into Y_t^n and μ_t components in (1.14) is useful only in the case that stabilization of the output gap is a goal of policy.

If expectations are assumed not to change in response to policy changes or other shocks, the model makes predictions like those of a standard undergraduate textbook exposition. For example, an increase in the central bank's interest-rate target i_t should reduce output and inflation. An increase in government purchases (that increases g_t , and also increases Y_t^n by a smaller amount) should increase both output and inflation. And a "cost-push shock" $u_t > 0$ should increase inflation, but have no effect on output if the central bank's interest-rate target is unchanged; if i_t is instead raised in response to the increase in inflation, output should fall and the inflation increase will be more modest.

But the model also indicates the effects on output and inflation of changes in average expectations. For example, an increase in the average forecast by firms of the log price that they would wish to choose if reviewing their price one period in the future, relative to the current average price, should raise current inflation for any current level of output, just as in the case of an exogenous cost-push shock. And there is no general reason to suppose that expectations should be unaffected by shocks of the kind considered in the previous paragraph; hence a complete analysis even of short-run equilibrium requires a specification of how expectations are determined.

If expectations are Ricardian and monetary policy is unaffected by fiscal variables, the model implies that neither the size of the public debt nor the government budget matter for output and inflation determination. If instead expectations are not assumed to be Ricardian, it is more convenient to write the model in terms of expectations of v_{t+1}^i rather than \bar{v}_{t+1}^i . The "IS equation" can then alternatively be written in the form

$$Y_t = g_t - \sigma i_t + (1 - \beta)b_{t+1} + \int \hat{E}_t^i v_{t+1}^i di, \quad (1.18)$$

where one should recall from (1.5) that b_{t+1} is known at date t (though the debt matures at $t + 1$). In this case, since the endogenous public debt matters for output and inflation determination, we must adjoin (1.5) to the system that describes TE dynamics under a given specification of inflation.

The complete TE dynamics for the endogenous variables $\{\pi_t, Y_t, b_{t+1}\}$ are then given in the non-Ricardian case by the system consisting of (1.5), (1.17) and (1.18), given expectations $\{v_t^i, p_t^{*j}\}$ consistent with (1.6) and (1.15). One observes that a larger budget deficit (or smaller surplus) should increase output and inflation, to the extent that it does not cause a reduction in the final expectational term in (1.18), of

a magnitude as large as the increase in $(1 - \beta)b_{t+1}$.

1.5 RE Equilibrium

Thus far, I have been completely agnostic about the nature of subjective forecasts. Under the *RE hypothesis*, all agents' probability beliefs are identical, and coincide with the probabilities predicted by one's model, given the choices that people make on the basis of those beliefs. Under the hypothesis that all beliefs are identical, we can replace the operators $\{\hat{E}^i[\cdot]\}$ and $\{\hat{E}^j[\cdot]\}$ by the single expectation operator $\hat{E}[\cdot]$. And since in this case $v_t^i = v_t$ for all i and $p_t^{*j} = p_t^*$ for all j , equations (1.10) and (1.15) reduce to

$$\bar{v}_t = -\sigma(i_t - \pi_t) + \hat{E}_t \bar{v}_{t+1}, \quad (1.19)$$

$$p_t^* = (1 - \alpha)^{-1} [\kappa y_t + u_t] + \beta \hat{E}_t p_{t+1}^* \quad (1.20)$$

respectively. If we assume a monetary policy rule (or central-bank reaction function) of the form¹³

$$i_t = \phi_\pi \pi_t + \phi_y y_t + \epsilon_t^i, \quad (1.21)$$

where ϵ_t^i is an exogenous disturbance to monetary policy, and use equations (1.8), (1.11), and (1.21) to substitute for i_t , π_t and y_t in equations (1.19)–(1.20), we obtain a system that can be written in the form

$$z_t = B \hat{E}_t z_{t+1} + b \xi_t, \quad (1.22)$$

where z_t is the vector of endogenous variables (\bar{v}_t, p_t^*) , ξ_t is the vector of (composite) exogenous disturbances $(g_t - Y_t^n, u_t, \epsilon_t^i)$, and B and b are matrices of coefficients.

If subjective probabilities must coincide with objective probabilities, we can replace (1.22) by

$$z_t = B E_t z_{t+1} + b \xi_t, \quad (1.23)$$

where $E_t[\cdot]$ denotes an expectational conditional on the state of the world at date t , under the probability distribution over future paths that represents the equilibrium outcome. An RE equilibrium (REE) is then a stochastic process $\{z_t\}$ consistent with

¹³This version of the “Taylor rule” (Taylor, 1993) reflects an implicit inflation target of zero. In particular, if all exogenous disturbances are zero forever, this policy is consistent with the steady state around which the model equations have been log-linearized.

(1.23). (Note that any solution for the process $\{z_t\}$ completely determines the stochastic evolution of the variables $\{\pi_t, Y_t, i_t\}$, using equations (1.8), (1.11), and (1.21).) We shall here restrict attention only to the possibility of bounded solutions $\{z_t\}$, on the assumption that the disturbances $\{\xi_t\}$ are bounded, since our log-linearized equations are derived under this assumption.

The RE hypothesis does not necessarily determine a unique set of model-consistent probability beliefs; because it is a consistency criterion, rather than a hypothesis about how beliefs are formed, it essentially defines a fixed-point problem. RE beliefs are a fixed point of the mapping from possible subjective probabilities into the implied objective probabilities; see Evans and Honkapohja (2001) for further discussion of this “*T-mapping*.” Such a fixed-point problem may or may not have a solution, and the solution may or may not be unique. In the case of a linear system such as (1.23), Blanchard and Kahn (1980) establish that the existence and uniqueness of bounded solutions depend on the eigenvalues of the matrix B . Because no elements of z_t are predetermined, RE equilibrium is *determinate* (a unique bounded RE solution exists) if and only if both eigenvalues of B lie inside the unit circle. Under the assumption that the response coefficients in (1.21) satisfy $\phi_\pi, \phi_y \geq 0$, one can show (Woodford, 2003, chap. 4) that this condition is satisfied if and only if the response coefficients also satisfy

$$\phi_\pi + \frac{1 - \beta}{\kappa} \phi_y > 1, \quad (1.24)$$

sometimes called the “Taylor Principle.”¹⁴ If monetary policy satisfies (1.24), (1.23) can be “solved forward” for z_t as a linear function of current and expected future values of the exogenous disturbances; in the case that the exogenous disturbances follow linear-Markovian dynamics, so that $E_t \xi_{t+1} = \Lambda \xi_t$ for some stable matrix Λ , this solution is given by

$$z_t = Z \xi_t, \quad (1.25)$$

¹⁴See Proposition 4.3 in Woodford (2003, chap. 4). Woodford analyzes a system of the form (1.23) but in which the vector z_t has inflation and the output gap as elements. But since this vector is a non-singular linear transformation of the vector z_t used here (plus an exogenous term, which does not affect the determinacy calculation), the eigenvalues of the matrix B in Woodford (2003) are the same as those of the matrix B here.

where¹⁵

$$Z \equiv \sum_{j=0}^{\infty} B^j b \Lambda^j.$$

The implied responses of inflation, output and interest rates to exogenous disturbances of various kinds are discussed further in Woodford (2003, chap. 4) and Gali (2008, chap. 3).

If instead (1.24) is not satisfied, there are an infinite number of REE (even restricting our attention to bounded solutions), including solutions in which inflation and output fluctuate in response to “sunspots” (random events with no consequences for the economic fundamentals ξ_t) or in which the fluctuations in inflation and output are arbitrarily large relative to the magnitude of the exogenous disturbances. In the case of such a policy, the economy may be vulnerable to instability due purely to the volatility of expectations, even under the assumption that the economy must evolve in accordance with an REE. Because instability of this kind is undesirable, it is often argued that a policy commitment should be chosen that ensures the existence of a determinate RE equilibrium (see, e.g., Woodford, 2003, chaps. 2, 4); in the context of New Keynesian models of the kind sketched here, this provides an argument for the desirability of an interest-rate rule that conforms to the Taylor Principle (see, e.g., Clarida, *et al.*, 2000).¹⁶ Apart from the use of the determinacy result as a criterion for the choice of a monetary policy rule, the predicted character of the fluctuations due to self-fulfilling expectations when the Taylor Principle is not satisfied has been proposed by some as a positive theory of the aggregate fluctuations observed during periods when monetary policy has arguably been relatively “passive” (e.g., Clarida *et al.*, 2000; Lubik and Schorfheide, 2004). As discussed below, however, relaxation of the RE hypothesis opens up additional possibilities for instability under regimes that fail to pin down expectations sufficiently precisely.

Some argue that avoidance of indeterminacy of REE need not be a concern when

¹⁵Note that this infinite sum must converge, because we have assumed that both B and Λ have all eigenvalues inside the unit circle.

¹⁶Some object to this argument for the Taylor Principle on the ground that it ensures only a *locally* unique REE — there is only one equilibrium in which the endogenous variables remain forever near the target steady state — but does not exclude the possibility of other REE, including sunspot equilibria, that do not remain near the steady state (e.g., Benhabib *et al.*, 2001). This is issue is beyond the scope of the current review, owing to our reliance on a local log-linear characterization of equilibrium dynamics; but see Woodford (2003, chap. 2, sec. 4) for further discussion.

choosing a monetary policy rule, on the ground that even in the indeterminate case, there is no reason to expect people’s expectations to coordinate on a sunspot equilibrium, or even on one with excessive fluctuations in response to fundamental disturbances (*e.g.*, McCallum, 1983). Such authors argue that a model’s positive prediction should be based on some further refinement of the set of equilibria, such as a restriction to *Markovian* equilibria, in which endogenous variables depend only on those aspects of the state of the world that affect either the equilibrium relations that determine those variables, or the conditional probabilities of states that will be relevant for equilibrium determination in the future.¹⁷

If $\{\xi_t\}$ is Markovian in the above example, this would mean restricting attention to REE of the form (1.25), for some matrix Z . Since (1.25) represents an REE if and only if the matrix Z satisfies

$$Z = BZ\Lambda + b,$$

and these are a system of 6 linear equations for the 6 elements of Z , there is a unique solution of this form for generic parameter values, even when the monetary policy rule fails to satisfy (1.24). (McCallum, 1983, calls this the “minimum-state-variable solution.”) But the question whether, or under what circumstances, we should expect people to coordinate on the particular expectations specified by the Markovian solution is a question that cannot be answered by the RE hypothesis itself; and the consideration of plausible restrictions on expectations that do not simply assume RE can be of help in justifying a particular selection from among the set of REE, as proposed by McCallum.

2 Rationalizable TE Dynamics

As discussed in the introduction, one broad approach to the formulation of a criterion for reasonableness of beliefs — without simply postulating an exact correspondence between people’s forecasts and those that are correct (according to one’s model of the economy) — is to assume that people should correctly *understand the economic model*, and be able to form correct inferences from it about possible future outcomes. This approach supposes that beliefs should be refined through a process of reflection,

¹⁷This refinement is closely related to the idea of restricting attention to the Markov perfect equilibria of dynamic games (Maskin and Tirole, 2001).

independent of experience and not necessarily occurring in real time, that Guesnerie (1992) calls *eduction*. While RE beliefs would certainly *withstand* a process of scrutiny of this kind, such beliefs need not be the *only* ones that could be rationalized in this way.

2.1 “Eductive Stability” Analysis

For the sake of concreteness, let us consider the case of Ricardian expectations and monetary policy specified by (1.21). The assumption that people in the economy “understand the model” means, in the present context, that people understand that the TE dynamics of inflation, output and interest rates will be determined by equations (1.16), (1.17) and (1.21) each period, given the average one-period-ahead forecasts of others. In order for someone’s expectations regarding the paths of the endogenous variables to be consistent with this knowledge, the expected paths must be able to be generated by these equations, under *some* supposition about the average expectations of others.

Let a possible conjecture about the evolution of average one-period-ahead forecasts be specified by a vector stochastic process $\mathbf{e} \equiv \{e_t\}$, where at any date the two elements of e_t specify $\int \hat{E}_t^i \bar{v}_{t+1}^i di$ and $\int \hat{E}_t^j p_{t+1}^{*j} dj$. For any such evolution of average forecasts, the structural equations determine unique TE processes $\{\pi_t, Y_t, i_t\}$. Hence any agent (household or firm) that expects average expectations to evolve in the future in accordance with the process \mathbf{e} , and that understands the model, must (in order to have internally consistent beliefs) forecast precisely this particular evolution of the variables $\{\pi_t, Y_t, i_t\}$. There is then a unique internally consistent anticipated evolution of this agent’s own one-period-ahead forecasts $\{\hat{E}_t^i \bar{v}_{t+1}^i\}$ (in the case of a household), implied by (1.10), and similarly a unique internally consistent anticipated evolution $\{\hat{E}_t^j p_{t+1}^{*j}\}$ (in the case of a firm), implied by (1.15). This allows us to determine a new vector stochastic process $\mathbf{e}' \equiv \{e'_t\}$ that describes the one-period-ahead forecasts that must be made by agents who understand the model and believe that the average forecasts of others will evolve according to \mathbf{e} .

Let Ψ denote the mapping that determines $\mathbf{e}' = \Psi(\mathbf{e})$ in the way just described. Then individual beliefs \mathbf{e}^i are consistent with knowledge of the model only if there exists some specification of average beliefs \mathbf{e} such that $\mathbf{e}^i = \Psi(\mathbf{e})$. In this case we can say that the beliefs \mathbf{e}^i can be *rationalized* by the conjecture \mathbf{e} about average beliefs.

Because of the linearity of the mapping Ψ , it is evident that if all agents understand the model, a specification of average beliefs \mathbf{e} is rationalizable if and only if there exists some conjecture about average beliefs \mathbf{e}^1 such that $\mathbf{e} = \Psi(\mathbf{e}^1)$. But if in addition all agents understand that all agents understand the model, their conjectures about average beliefs are consistent with this knowledge only if \mathbf{e}^1 can itself be similarly rationalized, which is to say, if there exists a conjecture \mathbf{e}^2 such that $\mathbf{e}^1 = \Psi(\mathbf{e}^2)$. Any number of levels of rationalization might be demanded in a similar way.

Even if we assume that agents' forecasts should be grounded in reasoning of this kind, it may be reasonable only to demand some finite number of levels of rationalization, either because it is only assumed that people understand that others understand ... that others understand the model, to some finite order of recursion; or because it is not considered practical for agents to check their beliefs for this degree of internal consistency beyond some finite number of levels.¹⁸ In this case, k th-order beliefs (for some finite k) are allowed to be specified arbitrarily (required to be internally consistent, but not necessarily consistent with understanding the model). In this case, obtaining definite conclusions requires a specific theory of k th-order beliefs (perhaps some fairly simple specification), or at least some bounds on the class of possible specifications of k th-order beliefs that may be entertained.

Alternatively, one may, as in the literature on “rationalizable equilibria” in game theory (Bernheim, 1984; Pearce, 1984), require that beliefs be consistent with an infinite hierarchy of beliefs, each level of which is rationalized by the next higher level of beliefs. RE beliefs represent one possible type of rationalizable beliefs in this sense; but not all rationalizable beliefs need be RE beliefs, and even when REE is determinate, there may be a large multiplicity of rationalizable beliefs, even under the requirement that beliefs satisfy some uniform bound at all levels. Guesnerie (1992, 2005) calls an investigation of whether the REE beliefs are the unique rationalizable beliefs “eductive stability analysis.” If the REE *is* eductively stable, he considers the REE outcome to be a reasonable prediction of one’s model; but if not, the other rationalizable paths are taken to be equally plausible predictions.

The existence of a large set of possible equilibrium outcomes, including the possibility of fluctuations in response to sunspots, or large fluctuations in response to small

¹⁸See, e.g., Phelps (1983) or Evans and Ramey (1992) for proposals of this kind. Evans and Ramey propose to endogenize the number of levels of rationalization in terms of “calculation costs” involved in iterating the mapping Ψ another time.

changes in fundamentals, is regarded as an undesirable form of instability. Hence Guesnerie proposes as a criterion for policy choice the desirability of finding a policy under which the REE is eductively stable.¹⁹ This is in the spirit of the proposal, discussed above, that policy be designed to ensure determinacy of RE equilibrium, but is an even stronger requirement, since uniqueness of rationalizable equilibrium necessarily implies determinacy of REE, while the converse is not true.

2.2 The Taylor Principle and Determinacy Reconsidered

As an illustration, let us consider whether a monetary policy rule of the form (1.21) ensures unique (uniformly bounded) rationalizable dynamics. To simplify, as in Guesnerie (2008), let us consider the limiting case of perfectly flexible prices. In this limit, output Y_t is determined by exogenous fundamentals, and (1.16) and (1.21) jointly determine i_t and π_t given expectations and exogenous fundamentals. Calculations are also simplified if there are assumed to be no exogenous disturbances (including no variation in Y_t), as the issue of the multiplicity of solutions is unaffected by the amplitude of disturbances. Hence the monetary policy rule (1.21) reduces simply to $i_t = \phi_\pi \pi_t$. REE is determinate if and only if $\phi_\pi > 1$.²⁰ When this condition is satisfied, the unique bounded REE has $\bar{v}_t = 0$ for all t , implying that $i_t = \pi_t = 0$ for all t .

While $\phi_\pi > 1$ is therefore also a necessary condition for uniqueness of the rationalizable dynamics, it is not sufficient. Guesnerie (2008) shows that if $1/2 < \beta < 1$ and

$$\phi_\pi > (2\beta - 1)^{-1} > 1, \quad (2.1)$$

then even though REE is determinate, there is a large multiplicity of uniformly bounded rationalizable equilibria. For example, consider the hierarchy of beliefs that may support a rationalizable TE at some date t . The requirement of rationalizability does not establish any necessary linkages between what happens at different dates: only between what happens at date t and what people expect at date t that people will expect at later dates about what people will expect ... will happen at still later dates. Thus for each date t we may separately specify what happens at that date, along with the hierarchy of forecasts of forecasts that rationalize it.

¹⁹Since the REE is necessarily a rationalizable equilibrium, uniqueness of rationalizable equilibrium requires that the REE be the only such equilibrium.

²⁰This is the implication of (1.24) in the limit in which $\kappa \rightarrow \infty$.

Substituting the policy rule for i_t in (1.10), and then using (1.8) to eliminate π_t , we obtain the requirement that

$$\bar{v}_t^i = (1 - \beta\phi_\pi)\bar{v}_t + \beta\hat{E}_t^i\bar{v}_{t+1}^i \quad (2.2)$$

for all i . In a rationalizable TE, not only must this hold at date t , but everyone must expect anyone else to expect anyone else ... to expect it to hold at any future date.

One such specification of the hierarchy of beliefs is given by $\bar{v}_t^i = \epsilon$, and

$$\begin{aligned} \hat{E}_t^{i_1}\hat{E}_{t_1}^{i_2}\dots\hat{E}_{t+j_n-1}^{i_n}\bar{v}_{t+j_n}^{i_n} &= (-\mu)^{1-n}\phi_\pi\epsilon, \\ \hat{E}_t^{i_1}\hat{E}_{t_1}^{i_2}\dots\hat{E}_{t+j_n-1}^{i_n}\bar{v}_{t+j_n}^{i_{n+1}} &= (-\mu)^{-n}\epsilon \end{aligned}$$

for any sequences of households and dates of the kind assumed above, where ϵ is an arbitrary real number and

$$\mu \equiv \frac{\beta\phi_\pi - 1}{(1 - \beta)\phi_\pi}.$$

These beliefs satisfy all of the requirements for rationalizability for any real number ϵ , as discussed further in the Appendix (available online). If (2.1) is satisfied, $\mu > 1$, and forecasts of all orders also satisfy a uniform bound. There is thus (at least) a continuum of uniformly bounded rationalizable TE. Moreover, ϵ may represent the realization of a “sunspot” event unrelated to fundamentals, so that there are seen to exist bounded sunspot equilibria, despite the fact that monetary policy satisfies the Taylor Principle.

If instead

$$1 < \phi_\pi < (2\beta - 1)^{-1}, \quad (2.3)$$

one can show that the determinate REE is also the only uniformly bounded rationalizable TE. Note that (2.2) implies that

$$\bar{v}_t = \phi_\pi^{-1} \int \hat{E}_t^i \bar{v}_{t+1}^i di. \quad (2.4)$$

Hence if it is common knowledge that there exists some finite bound κ such that $|\bar{v}_t^i| \leq \kappa$ for all i and t , it follows from (2.4) that it must also be common knowledge that $|\bar{v}_t| \leq \phi_\pi^{-1}\kappa$ for all t . Using this bound, (2.2) then implies that

$$|\bar{v}_t^i| \leq |1 - \beta\phi_\pi|\phi_\pi^{-1}\kappa + \beta\kappa = \lambda\kappa,$$

where

$$0 < \lambda \equiv \max\{\phi_\pi^{-1}, 2\beta - \phi_\pi^{-1}\} < 1.$$

Hence common knowledge that $|\bar{v}_t^i| \leq \kappa$ implies that it must be common knowledge that $|\bar{v}_t^i| \leq \lambda\kappa$. By the same reasoning, it must then be common knowledge that $|\bar{v}_t^i| \leq \lambda^2\kappa$, and so on, until a bound smaller than any positive quantity is established. Thus it must be common knowledge that $v_t^i = 0$ for all i and all t , and hence that $i_t = \pi_t = 0$ for all t .

In such a case, Guesnerie says that the REE is “eductively stable,” and argues that there is reason (in this case, and this case only) to expect this equilibrium to obtain. While this is possible, the restrictions on the coefficient ϕ_π are much more stringent than those required for determinacy under the RE hypothesis. If, for example, $\beta = 0.99$ (a common calibration for quarterly New Keynesian models), (2.3) requires that $1 < \phi_\pi < 1.02$. This very tight bound is violated by the rule recommended by Taylor (1993), as well as by most estimated central-bank reaction functions.

If the avoidance of instability due to self-fulfilling expectations of this particular type is a design criterion for monetary policy, it follows that one must be careful about seeking to stabilize inflation in the face of real disturbances simply by using a rule of the form (1.21) with a very strong inflation response coefficient. Instead, the simultaneous achievement of eductive stability and stable inflation despite exogenous disturbances is possible only in the case of a policy rule that directly responds to the underlying *determinants* of inflation, namely, the exogenous disturbances and observed subjective forecasts.

3 Learning Dynamics

An alternative way of disciplining the specification of expectations does not demand that they be consistent with a correct structural model of the variables that are forecasted, but instead requires that the probabilities assigned to possible future outcomes are not too different from the probabilities with which outcomes actually occur. The idea is that it is not reasonable to suppose that people should fail to notice predictable regularities in economic data, whether or not they understand why those regularities exist.

But which regularities are the ones that one can reasonably expect people to take into account? A common answer assumes that forecasts should be derived through extrapolation from prior observations. Such approaches, based on explicit models of learning, have the advantage of explaining how the postulated similarity between subjective beliefs and actual patterns in the data comes about. They also reduce the problem of indeterminacy of the model's predictions, pervasive in the case of approaches like those discussed thus far, which demand only that beliefs be a fixed point of a certain mapping. While there may be many possible asymptotic states of belief under an explicit model of learning, with the one that is reached depending on initial conditions and/or random events along the way, a model of learning often makes a unique prediction conditional upon initial conditions and the subsequent history of shocks.

The most common approach of this general type assumes that agents' forecasts at any time t are derived from an econometric model, estimated using the data observed up until that date.²¹ Let the model be specified by a vector θ of parameters, and suppose that any model θ implies that forecasts e_t should be some function of the current state ζ_t . Then in any period t , new estimates $\hat{\theta}_t$ and $\hat{\zeta}_t$ are formed of the parameters and of the state, based on the data available at that point. Under a recursive estimation scheme, the new estimates are functions of the prior estimates and of the new data observed since the formation of the prior estimates,

$$\hat{\theta}_t = \Theta_t(\hat{\theta}_{t-1}, \zeta_{t-1}, x_t), \quad (3.1)$$

$$\hat{\zeta}_t = Z_t(\hat{\theta}_{t-1}, \hat{\zeta}_{t-1}, x_t) \quad (3.2)$$

where x_t is some vector of new data.²² Given these new estimates, period t forecasts are given by a function²³

$$e_t = \Psi(\hat{\zeta}_t; \hat{\theta}_t). \quad (3.3)$$

²¹See Evans and Honkapohja (2001, 2009) and Sargent (2008) for reviews of work of this kind.

²²The time subscripts on the functions $\Theta_t(\cdot)$ and $Z_t(\cdot)$ allow for the possibility that the updating rules may depend on the size of the existing dataset. Equation (3.15) below for the evolution of the mean estimates provides a simple example.

²³Here, for simplicity, I assume that every household i and every firm j forecasts in the same way, using the same observed data, so that average forecasts are simply the common forecasts, given by a function of the common estimates. One may, however, allow each household and to have its own estimated model $\hat{\theta}_t^i$, evolving according to a separate equation of the form (3.1), and then define $\int \hat{E}_t^i \bar{v}_{t+1}^i di$ as a function of the entire probability distribution of estimates $\{\hat{\theta}_t^i\}$, rather than simply as a function of a single estimate $\hat{\theta}_t$; and similarly with the firms.

Given these beliefs, the TE values of the variables z_t are determined by (1.22), given the exogenous disturbances ξ_t . This system, possibly along with additional structural equations, determines the new data x_t .²⁴ Thus the system of equations (1.22) and (3.1)–(3.3) jointly determines $\hat{\theta}_t, \hat{\zeta}_t, e_t$, and x_t , given the lagged estimates and the disturbances ξ_t . Solution of these equations in each of a succession of periods yields the predicted dynamics of both beliefs and endogenous variables, as a function of the history of exogenous disturbances.

3.1 Restricted Perceptions Equilibrium

A focus of much of the literature on TE dynamics with learning has been to ask whether learning dynamics should converge asymptotically to REE; indeed, much of the early literature (beginning with Bray, 1982) was concerned more with the foundations of the REE concept — seeking to provide a causal explanation for how the postulated coincidence between subjective and objective probabilities could come about — than with the provision of an alternative model of economic dynamics. Obviously this is only possible if the class of forecasting models that are contemplated includes a model θ^{REE} that produces the forecasts associated with the REE. If one does not assume that economic agents are endowed with knowledge of the structural model, and hence with the information required to compute the REE, it is not obvious that their forecasting approach should even entertain as a possibility the precise forecasting rule implied by the REE; but if no value of θ results in forecasts of this kind, convergence to REE beliefs (and hence to the REE dynamics) is obviously impossible.

There might, however, still be convergence of beliefs to some fixed point $\bar{\theta}$ with the property that under the TE dynamics generated by the beliefs $\bar{\theta}$, $\bar{\theta}$ is the model (among the class of models considered) that yields the best forecasts (under some

²⁴It is common in the literature on learning dynamics to specify a recursive causal structure by assuming that the data x_t are actually determined in period $t - 1$ (*i.e.*, they include π_{t-1} rather than π_t , and so on). In this case, all of the arguments of the functions in (3.1)–(3.2) are predetermined, so that these equations determine the new estimates $(\hat{\theta}_t, \hat{\zeta}_t)$ independently of the period- t shocks; equation (3.3) then determines the forecasts e_t ; and finally equations (1.22) determine the endogenous variables z_t given the shocks. But it is not obvious why, in the logic of the NK model presented here, one should suppose that period t forecasts must be made prior to the observation of period t endogenous variables.

criterion that is used for the estimation). For example, suppose that the class of models considered consists of those in which both \bar{v}_{t+1} and p_{t+1}^* are linear functions of ζ_t and unforecastable disturbances at $t + 1$; that the elements of ζ_t are all part of the history of the observables $\{x_t, x_{t-1}, \dots\}$, so that ζ_t is observable; and that the coefficients of the linear model are estimated so as to minimize the mean squared error of the forecasts of \bar{v}_{t+1} and p_{t+1}^* . Then forecasts will be of the form $e_t = \hat{\theta}'\zeta_t$, where the vector ζ_t is assumed to include an element equal to 1 each period; and with an infinite sequence of data generated by the TE dynamics under beliefs $\bar{\theta}$, the estimated coefficients will satisfy

$$\hat{\theta}' = E[z_{t+1}\zeta_t']E[\zeta_t\zeta_t']^{-1}, \quad (3.4)$$

where $E[\cdot]$ indicates an unconditional expectation under the ergodic TE dynamics resulting from some beliefs θ . Alternatively, we can write $e_t = P_t[z_{t+1}]$, where $P_t[\cdot]$ denotes the linear projection of the random variable inside the brackets on the space spanned by the elements of ζ_t .

The beliefs $\bar{\theta}$ constitute a *restricted perceptions equilibrium* (RPE) if the optimal estimate $\hat{\theta}$ given by (3.4) is equal to $\bar{\theta}$ when the unconditional expectations are the ones implied by the TE dynamics generated by beliefs $\theta = \bar{\theta}$ (Evans and Honkapohja, 2001, chap. 13; Branch, 2004). This is a weaker requirement than that of an REE, as forecasts are assumed to be optimal only within a particular class of linear models, rather than within the class of all forecasts that might be made on the basis of information available in period t . Note that in the special case that the optimal forecast of z_{t+1} is indeed a linear function of ζ_t , so that

$$E_t[z_{t+1}] = P_t[z_{t+1}] = T(\theta)'\zeta_t \quad (3.5)$$

when the TE dynamics are generated by beliefs θ , then (3.4) implies that $\hat{\theta} = T(\theta)$. Hence in this case, $\bar{\theta}$ describes RPE beliefs if and only if $T(\bar{\theta}) = \bar{\theta}$, which is also the condition for REE beliefs. More generally, however, when the forecasting variables ζ_t do not span a large enough space, or at any rate not the correct one, RPE beliefs will differ from REE beliefs.

Conditions can be established under which the learning dynamics resulting from repeated re-estimation of an OLS forecasting equation (*least-squares learning* dynamics) converge with probability 1 to an RPE as the length of the observed data set grows large enough. But even in such a case, the dynamics need not coincide, even

asymptotically, with the model's REE dynamics. Fuster *et al.* (2010, 2011, 2012) provide examples in which more complex dynamics of asset prices, consumption and investment are implied by RPE dynamics than would be associated with REE dynamics of the same models; here the suboptimality of forecasts results from estimation of lower-order autoregressive models of the data than the correct model.

3.1.1 Application to the NK Model

Suppose that we equate the subjective expectations in (1.22) with linear projections, to obtain

$$z_t = BP_t z_{t+1} + b\xi_t. \quad (3.6)$$

This is the set of conditions that must be satisfied in order for the evolution of the expectational variables $\{z_t\}$ to represent an RPE, under either of two possible interpretations of how forecasts for horizons more than period in the future are formed.

On the one hand, we might assume, as Preston (2005) does, that forecasts for arbitrary future horizons are formed by estimating a vector-autoregressive system

$$P_t x_{t+1} = \Lambda_x \zeta_t, \quad P_t \zeta_{t+1} = \Lambda_\zeta \zeta_t,$$

where x_t is the vector of variables that must be forecasted, other than the future values of the forecasting variables ζ_t themselves. Forecasts for arbitrary future horizons can then be computed as

$$\hat{E}_t^i x_{t+j} = \hat{E}_t^i \hat{E}_{t+1}^i \cdots \hat{E}_{t+j-1}^i x_{t+j} = \Lambda_x \Lambda_\zeta^{j-1} \zeta_t$$

for any $j \geq 1$. That is, forecasts for horizons more than one period in the future are formed by forecasting one's own future one-period-ahead forecasts, while one-period-ahead forecasts are given by linear projections on the forecasting variables ζ_t .²⁵ Given forecasts of this kind, the definitions (1.9) and (1.12) of the expectational variables then imply that their evolution must satisfy (3.6).

Alternatively, we might assume, as Evans and McGough (2009) propose, that people estimate the values of the expectational variables z_t , not using the definitions (1.9) and (1.12) of these variables in terms of long-horizon forecasts of variables with

²⁵An advantage of this method is that forecasts $\hat{E}_t^i x_{t+j}$ can be formed for arbitrarily large j , using coefficients that can be estimated using finite data sets, as there is no need to actually regress observed values of x_{t+j} on the prior forecasting variables ζ_t .

objective definitions, but instead using the recursive relations (1.10) and (1.15) to estimate values on the basis of one's current forecasts of one's own estimates in the next period. These forecasts of one's own future estimates can be obtained by collecting data on what one's estimates have been, and regressing them on the previous period's values of the forecasting variables ζ_t . Also under this assumption, if there is convergence to an RPE, the expectational variables will necessarily satisfy (3.6). Note that these two approaches to "least-squares learning" are not mathematically equivalent, but if in each case there is convergence to an RPE, then the RPE is the same in both cases.²⁶

Given a solution for the dynamics of the expectational variables $\{z_t\}$, the dynamics of inflation, output, and interest rates are then determined by equations (1.8), (1.11) and (1.21), as under other specifications of expectations. The difference between REE and RPE predictions then stems entirely from the difference between (3.6) and (1.23). When (3.5) holds, the predictions are necessarily the same.

3.1.2 Failure of Ricardian Equivalence

As an illustration of how macroeconomic dynamics in an RPE may differ from the REE dynamics predicted in the case of a given policy rule, let us reconsider the argument for Ricardian equivalence. Suppose that expectations are *not* Ricardian, i.e., that people do not assume that the future path of primary surpluses must satisfy (1.7), and instead estimate an econometric model to forecast future surpluses that does not impose this condition as an *a priori* restriction.²⁷ The TE dynamics are then determined by the system consisting of (1.5), (1.17) and (1.18), given expectations of the future evolution of the variables $\{p_t^{*j}, v_t^i\}$ defined by (1.6) and (1.15), and specified paths for the policy variables $\{i_t, s_t\}$.

Suppose further that monetary policy is described by a Taylor rule of the form (1.21), the coefficients of which satisfy the Taylor Principle (1.24), while fiscal policy

²⁶Since the learning dynamics outside the RPE are in general not identical, the conditions for convergence to an RPE are not always the same in the two cases.

²⁷Note that each household needs only to forecast *its own* tax obligations in excess of the value of government purchases; its use of a forecasting model that violates (1.7) does not necessarily imply that it believes that *aggregate* tax collections do not satisfy the present-value relation. While I assume that the tax obligations of all households are the same, this is not necessarily known to the households.

is described by a feedback rule of the form

$$s_t = \phi_b b_t + \epsilon_t^s, \quad (3.7)$$

where

$$1 - \beta < \phi_b < 1, \quad (3.8)$$

and ϵ_t^s is an exogenous disturbance. This specification, together with (1.5), implies debt dynamics that remain bounded in the case of *any* bounded processes for inflation and the nominal interest rate, and hence that model-consistent expectations will be Ricardian, in the case of any REE involving bounded fluctuations. Such a specification of monetary and fiscal policy implies the existence of a determinate REE in the case of any bounded disturbance processes, and in this REE, the fiscal shocks $\{\epsilon_t^s\}$ have no effects on the evolution of output, inflation or interest rates.²⁸ Hence in such a model, rational expectations imply Ricardian equivalence.

Let us consider instead the possible character of RPE dynamics. Suppose, for example, that the vector of forecasting variables ζ_t consists only of the single state variable s_t .²⁹ Then RPE forecasts are of the form $e_t = \psi s_t$,³⁰ for some vector of coefficients ψ . Substituting these forecasts into (1.17) and (1.18), one can solve for the TE values of π_t, y_t, i_t , and b_{t+1} as linear functions of b_t, s_t and ξ_t .

The calculations are especially simple if we consider a case in which $\phi_y = s_b = \kappa = 0$, and assume that there are no exogenous disturbances other than the fiscal shock $\{\epsilon_t^s\}$, which is assumed to be unforecastable (white noise). In this limiting case, there are no equilibrium fluctuations in π_t or i_t , and the solutions for y_t and b_{t+1} are given by

$$y_t = (\beta^{-1} - 1)(b_t - s_t) + \psi_v s_t, \quad (3.9)$$

²⁸See Woodford (2001) for further discussion of the consequences of a “locally Ricardian” fiscal policy of this kind, under an REE analysis.

²⁹Note that under the REE dynamics, if the disturbances are all unforecastable (white noise), none of the state variables that must be forecasted by households or firms are forecastable, *except* the primary surplus (that must be forecasted in order to estimate v_t^i). It is perhaps not implausible to suppose that households forecast future primary surpluses using only the current level of the surplus. This would not, however, constitute a model-consistent forecast, as (1.5) and (3.7) imply that an optimal forecast of future primary surpluses depends on the value of b_{t+1} , or alternatively upon both s_t and b_t .

³⁰In this non-Ricardian case, the first element of e_t is assumed to be $\int \hat{E}_t^i v_{t+1}^i di$ rather than $\int \hat{E}_t^i \bar{v}_{t+1}^i di$.

$$b_{t+1} = \beta^{-1}(b_t - s_t), \quad (3.10)$$

where ψ_v (to be determined) is the first element of ψ . It then follows from (1.6) that

$$v_t^i = y_t - (1 - \beta)b_t \quad (3.11)$$

for all i .

From this one can show that

$$\begin{aligned} e_{1t} &= P_t[y_{t+1} - (1 - \beta)b_{t+1}] \\ &= (\beta^{-1} - 1)(1 - \beta)P_t[b_{t+1}] + (\psi_v + 1 - \beta^{-1})P_t[s_{t+1}] \\ &= [(1 - \beta - \phi_b)(\beta^{-1} - 1) + \phi_b\psi_v]P_t[b_{t+1}], \end{aligned}$$

where for any variable x_{t+1} known at date t , $P_t[x_{t+1}]$ denotes the linear projection of x_{t+1} on s_t . (Here the first line uses the fact that (3.11) must also hold at date $t + 1$; the second line uses the fact that (3.9) must also hold at $t + 1$; and the third line uses the fact that s_{t+1} is determined by (3.7).) Writing $P_t[b_{t+1}] = \Lambda_b s_t$, where the coefficient Λ_b depends only on β and ϕ_b (not the assumed value of ψ_v), one then observes that ψ_v must satisfy the consistency condition

$$\psi_v = [(1 - \beta - \phi_b)(\beta^{-1} - 1) + \phi_b\psi_v]\Lambda_b. \quad (3.12)$$

Under assumption (3.8), this equation has a unique solution for ψ_v , and implies that³¹

$$\psi_v < \beta^{-1} - 1. \quad (3.13)$$

Equation (3.9) then implies that in the unique RPE, an exogenous positive innovation in the size of the primary surplus s_t lowers current output y_t . It will also reduce the debt b_{t+1} carried into the next period, with persistent effects on economic activity in later periods as well. Hence Ricardian equivalence does not hold in the RPE, even though the specification of fiscal policy implies that in the model's unique bounded REE, fiscal shocks have no effect on output, either immediately or subsequently. While (3.10) implies that even under the RPE dynamics, a *correct* forecast would satisfy (1.7) at all times, households' forecasts *do not* satisfy this condition, as a result of forecasting future surpluses purely on the basis of the current primary surplus; and because of this systematic forecasting error, Ricardian equivalence fails.

³¹See the Appendix for details.

3.2 “Learnability” of REE

Even when the class of contemplated forecasting models does include the REE forecasts, and even when the estimator used to determine $\hat{\theta}$ is one that should be asymptotically consistent, in the case of a sufficiently long series of data generated by the REE, it need not follow that $\hat{\theta}$ must converge asymptotically to θ^{RE} under the TE dynamics with learning. The reason is that, at each point in time, the observed data will actually be generated by the behavior that results from current beliefs $\hat{\theta}_t$, and not by REE behavior. If a departure of people’s estimates from θ^{RE} gives rise to patterns in the data that justify estimates even *farther* from θ^{RE} , the learning dynamics may diverge from RE beliefs almost surely, even if people start out with beliefs quite near to RE beliefs. The question whether the REE can in fact be reached as the asymptotic outcome of a learning process of the kind described above is therefore a non-trivial one. Authors such as Bullard and Mitra (2002) and Evans and Honkapohja (2003, 2006) propose as a design criterion for a monetary policy rule not only that the rule should be consistent with a desirable REE, but that the rule should imply that learning dynamics should converge to that REE, so that the desirable equilibrium is “learnable.”

3.2.1 Adaptive Estimation of Means

As a simple example, suppose that the disturbances ξ_t are all independently and identically distributed (i.i.d.) random variables, with mean zero. In this case, there is a unique Markovian REE, in which $z_t = b\xi_t$ each period, and the RE forecasts satisfy $E_t z_{t+1} = 0$ at all times. In this equilibrium, π_t, y_t and i_t will also each be a linear function of ξ_t , and the RE forecast of each variable will be zero (*i.e.*, the constant steady-state value) at all times. Suppose furthermore that the class of forecasting models considered by decisionmakers consists of all models under which the forecast of each variable is a constant (that is, people believe that each of these is an i.i.d. random variable, and seek only to estimate its mean). This simple class of models includes the forecasting rule used in the Markovian REE, so the assumption of such a restricted class does not in itself rule out the possibility of convergence to RE beliefs.

Finally, suppose that people estimate the means of each of the stationary variables

using the sample mean of the values observed to date, so that

$$\hat{x}_t = t^{-1} \sum_{s=1}^t x_s, \quad (3.14)$$

where x_t refers to any of the variables π_t, y_t, i_t or to any of the elements of ξ_t ; \hat{x}_t is the estimate of the variable's mean at date t (common to all agents); and 1 is the date at which the available data series begins.³² This can be written recursively in the form

$$\hat{x}_t = \hat{x}_{t-1} + \gamma_t(x_t - \hat{x}_{t-1}), \quad (3.15)$$

where the “gain” $\gamma_t = 1/t$ indicates the degree to which estimates are adjusted in response to an observation that differs from what has been forecasted. Note that if the case that the economy were to reach the REE, each of the variables x_t would indeed be i.i.d., and the estimators (3.14) would almost surely converge asymptotically to the true means (and hence to the REE beliefs), by the law of large numbers. Hence the estimation strategy is not inherently incompatible with learning the REE beliefs.

Any set of estimates of the means implies forecasts given by³³

$$\hat{E}_t^i \bar{v}_{t+1}^i = \hat{y}_t - \hat{g}_t - \frac{\sigma}{1 - \beta}(\beta \hat{\lambda}_t - \hat{\pi}_t), \quad (3.16)$$

$$\hat{E}_t^j p_{t+1}^{*j} = \frac{1}{1 - \alpha\beta} \hat{\pi}_t + \xi \hat{y}_t + \hat{\mu}_t, \quad (3.17)$$

for each household and each firm. Hence in vector form we can write

$$e_t = C \hat{\mathbf{x}}_t, \quad (3.18)$$

for a certain matrix of coefficients C , where $\hat{\mathbf{x}}_t$ is the vector of estimates (3.14). The system consisting of (1.16), (1.17) and (1.21) allows us to solve for TE values

$$\mathbf{x}_t = D e_t + d \xi_t, \quad (3.19)$$

for certain matrices D and d , where \mathbf{x}_t is the vector of actual values of the variables x_t . Equations (3.15), (3.18) and (3.19) then completely describe the TE dynamics of

³²This is an example of least-squares learning, in which the vector s_t has a single element, 1, each period.

³³Here I use the notation \hat{g}_t for the current estimate of the mean of the composite disturbance $g_\tau - Y_\tau^n$, for simplicity.

actual values and forecasts with adaptive learning, given the exogenous disturbances $\{\xi_t\}$ and initial prior estimates $\hat{\mathbf{x}}_{t-1}$.

Combining these equations, we obtain a law of motion

$$[I - \gamma_t DC] \hat{\mathbf{x}}_t = (1 - \gamma_t)\hat{\mathbf{x}}_{t-1} + \gamma_t d\xi_t \quad (3.20)$$

for the estimates; thus the TE dynamics are uniquely defined as long as the matrix in square brackets is non-singular, as I shall assume.³⁴ In fact, all that matters about these estimates is the implied forecasts e_t ; so we can reduce the dimension of the system (3.20) by pre-multiplying by the matrix C , yielding

$$[I - \gamma_t A] \Delta e_t = -\gamma_t [I - A] e_{t-1} + \gamma_t a \xi_t, \quad (3.21)$$

where $A \equiv CD$, $a \equiv Cd$. This equation determines the dynamics of the forecasts $\{e_t\}$ given initial forecasts and the evolution of the exogenous disturbances; the paths of the other relevant variables are then given by (3.19). The TE dynamics converge asymptotically to the REE dynamics (and subjective expectations coincide asymptotically with the REE forecasts) if and only if $e_t \rightarrow 0$ for large t .

Using the stochastic approximation methods introduced by Marcet and Sargent (1989) and expounded in detail in Evans and Honkapohja (2001), one can show that in the case of a decreasing gain sequence $\{\gamma_t\}$ like the one implied by (3.14), the path for $\{e_t\}$ implied by the stochastic law of motion (3.21) eventually converges to one of the trajectories of the ordinary differential equation (ODE) system

$$\dot{e} = -[I - A] e(\tau), \quad (3.22)$$

where τ is a re-scaled time variable defined by $\tau_t \equiv \sum_{s=1}^t \gamma_s$, and the dot indicates a derivative with respect to τ .

The ODE system (3.22) has a unique rest point $\bar{e} = 0$ (corresponding to REE forecasts) if and only if the 2×2 matrix A has no eigenvalue exactly equal to 1 (so that $I - A$ is non-singular); and the trajectories of (3.22) converge asymptotically to this rest point if and only if both eigenvalues of A have real parts less than 1 (so that both eigenvalues of $I - A$ have positive real parts). If the latter condition holds, (3.21) implies that $e_t \rightarrow 0$ with probability 1 as t grows large; beliefs asymptotically

³⁴For any matrices C and D , this will be true for all small enough values of the gain γ_t . We are here only concerned with TE dynamics in the low-gain case.

approach the REE forecasts, and one may say that the REE is “learnable” following the procedure postulated above. If, instead, A has an eigenvalue with real part greater than 1, the trajectories of (3.22) diverge from the rest point for almost all initial conditions, and correspondingly, one can show that there is zero probability of the beliefs implied by (3.21) remaining forever within a neighborhood of the REE beliefs, even if people begin with initial beliefs near (or exactly equal to) the REE beliefs.

Hence the learnability of the REE depends on the eigenvalues of A . In the case of a policy (1.21) with $\phi_\pi, \phi_y \geq 0$, one can show that both eigenvalues of A have real part less than 1 (implying learnability) if and only if the response coefficients satisfy (1.24) — that is, policy conforms to the Taylor Principle.³⁵ This is identical to the condition for determinacy of the REE dynamics, and the connection between the two results is not accidental. The matrix A has an eigenvalue equal to 1 if and only if the model’s steady-state inflation rate and output gap are *indeterminate*: the associated right eigenvector \bar{e} indicates the direction in which a constant forecast ($e_t = \bar{e}$ for all t) may differ from zero and still constitute a perfect foresight equilibrium if the disturbances equal zero.³⁶ (Any multiple of \bar{e} is also a possible perfect-foresight steady state in such a case.) But there exists a continuum of steady states if and only if the matrix B in (1.23) has an eigenvalue equal to 1, and \bar{e} must also be the associated right eigenvector of B . We have seen above that B has such an eigenvalue if and only if (1.24) holds with equality.

Intuitively, the Taylor Principle guarantees determinacy of the REE dynamics, because perturbations of the expected future values of the elements of z result in a current TE value of z_t that is *closer* to zero than whatever is expected for the next period; hence REE forecasts $E_t z_{t+j}$ cannot be bounded for all j and also consistent with this contraction requirement, unless they are exactly zero for all j . But the fact that forecasts e_t different from zero give rise to TE values z_t that are closer to zero (on average) also implies that adaptive learning will move the forecasts closer to zero (on average), so that the learning dynamics eventually converge to the REE

³⁵In fact, when (1.24) is satisfied, both eigenvalues of A are inside the unit circle. When it fails, there is a real eigenvalue greater than 1.

³⁶If $A\bar{e} = \bar{e}$, then forecasts $e_t = \bar{e}$ each period lead to outcomes $\mathbf{x}_t = D\bar{e}$ each period. Then $\hat{\mathbf{x}}_t = D\bar{e}$ will be a perfect-foresight estimate each period, implying that the correct forecasts each period will be $CD\bar{e} = \bar{e}$.

forecasts.³⁷

The analysis above assumes a very simple kind of least-squares learning, in which the only contemplated forecasting rules are ones in which the forecasts are constants (estimates of the means of the various variables) and the same for all horizons. But Preston (2005) establishes the same conditions for learnability of the REE when people use forecasting rules of the form

$$\hat{E}_t x_{t+1} = \hat{a} + \hat{b}' \xi_t$$

for each variable x , and estimate the coefficients \hat{a}, \hat{b}' by regressing observations of x_{t-j} on ξ_{t-j-1} (for $0 \leq j \leq t-1$). Longer-horizon forecasts are then formed using

$$\hat{E}_t x_{t+k} = \hat{a} + \hat{b}' \hat{E}_t \xi_{t+k-1},$$

where forecasts of the future disturbances are based on estimated autoregressive models of each disturbance.³⁸

A policy inconsistent with the Taylor Principle still leads to instability, because perturbations of the estimated constant terms \hat{a} result in average values of the variables that differ from zero by an even greater amount, leading to explosive dynamics for the estimates as above. And on the other hand, conformity to the Taylor principle remains sufficient for stability of the learning dynamics, since the conditions under which estimates of the response coefficients \hat{b}' diverge are even more restrictive than those required for divergence of the estimates of the constant terms. Preston also shows that the Taylor Principle is necessary and sufficient for learnability of the MSV REE (1.25) using this approach, in the case that the disturbances are AR(1) processes (and hence Markovian). This result again follows because the key to convergence to the REE forecasting rule is the convergence of the estimates of the constant terms \hat{a} , and the mean dynamics of these estimates are unaffected by the stationary fluctuations in the disturbances.³⁹

³⁷See Woodford (2003, chap. 4, sec. 2.3) for further discussion.

³⁸See the discussion above in section 3.1.1 for further description of Preston's method of VAR-based forecasts.

³⁹Bullard and Mitra (2002) reach a similar conclusion, though on the basis of assumed TE dynamics derived by substituting subjective expectations for objective expectations in certain Euler equations of the REE model, rather than deriving the TE dynamics from infinite-horizon optimization under subjective expectations, as above. For comparison of this "Euler-equation approach" to modeling learning dynamics with the one used here, see Preston (2005) and Evans and McGough (2009).

3.2.2 The Possibility of a “Deflation Trap”

Thus under this approach we again conclude that a rule (1.21) that fails to conform to the Taylor Principle (1.24) makes the economy vulnerable to instability due to self-fulfilling fluctuations, though through a different mechanism than in section 1.5 above. The explosive dynamics of forecasts in the case of insufficient feedback from aggregate outcomes (especially for inflation) to the interest-rate target generalizes the informal argument of Friedman (1968) about the instability resulting from an interest-rate peg.⁴⁰

The problem is not necessarily avoided, however, by commitment to a rule that satisfies the Taylor principle. The reason is that the linear TE dynamics analyzed above cannot hold *globally*; in particular, policy cannot be described by (1.21) for all possible inflation rates and output gaps, because of the zero lower bound on the nominal interest rate.⁴¹ Even if the central bank follows (1.21) with coefficients satisfying (1.24) until the lower bound becomes a binding constraint, the altered response at low levels of inflation and output implies the existence of a second (deflationary) perfect-foresight steady state consistent with the policy rule, as discussed by Benhabib *et al.* (2001); and the insensitivity of the interest rate to variations in inflation and output once the lower bound binds implies that the learning dynamics will be unstable near the forecasting rule associated with the deflationary REE.⁴²

This is illustrated in Figure 1, where the global behavior of the ODE system corresponding to (3.22) is plotted now taking into account the zero lower bound (ZLB). Note that the third row of (3.19) can be written as $i_t = D'_i e_t$, if we average out the values of ξ_t (in order to describe the mean dynamics, which approximately characterize the asymptotic dynamics of our system). The asymptotic dynamics described by (3.22) are therefore consistent with the zero lower bound as long as $e(\tau)$

⁴⁰Howitt (1992) was the first attempt to formalize Friedman’s argument through an analysis of the convergence of learning dynamics to the REE.

⁴¹The other structural relations assumed above are merely local approximations to relations that should actually be nonlinear; but even if they are assumed to hold globally, the zero lower bound prevents (3.21) from holding globally.

⁴²Evans and Honkapohja (2010) and Benhabib *et al.* (2012) show this in the context of NK models with adaptive learning closely related to the one presented here.

remains in the region satisfying the inequality⁴³

$$\bar{r} + D'_i e \geq 0, \quad (3.23)$$

where $\bar{r} > 0$ is the steady-state real rate of interest.⁴⁴ Because both elements of D_i are positive under the sign assumptions stated above, the region in Figure 1 where the dynamics (3.22) apply is the region above and to the right of the line labeled *ZLB*, along which (3.23) holds with equality.

When (3.23) is violated, (1.21) must instead be replaced by

$$i_t = -\bar{r} < 0. \quad (3.24)$$

Solving the system consisting of (1.16)–(1.17) and (3.24), one obtains a linear solution of the form

$$x_t = \underline{x} + \underline{D}e_t + \underline{d}\xi_t \quad (3.25)$$

instead of (3.19). The complete TE solution for x_t is then given by (3.19) when (3.23) is satisfied, and (3.25) otherwise. (Note that this is a continuous, piecewise linear solution.) Repeating the derivation of (3.22), one finds that \dot{e} is given by (3.22) when (3.23) is satisfied (a region that includes a neighborhood of the origin), and instead by

$$\dot{e} = -(I - \underline{A})e + C\underline{x} \quad (3.26)$$

in the region where the inequality is reversed, where $\underline{A} \equiv C\underline{D}$. (Note that this makes \dot{e} a continuous, piecewise-linear function of e .) This is the system the trajectories of which are plotted in Figure 1.⁴⁵

One observes that the origin $e = 0$ (corresponding to the zero-inflation steady state) is a rest point of these dynamics, and locally stable under the ODE dynamics as discussed above. If the dynamics (3.22) applied globally (i.e., if the ZLB constraint were not an issue), this steady state would also be globally stable: the learning dynamics would converge to it asymptotically from all possible initial states of belief.

⁴³Recall that in the notation used here, i_t is the amount by which the nominal interest rate exceeds its steady-state value, so that the requirement for the nominal interest rate to be non-negative is $\bar{r} + i_t \geq 0$.

⁴⁴Because we log-linearize our equations around a stationary equilibrium with zero inflation, \bar{r} is also the steady-state nominal interest rate.

⁴⁵Analytical derivations of qualitative properties of this figure are given in the Appendix.

But the dynamics when the ZLB constraint binds are different; as a consequence, there is a second steady state in the region below the *ZLB* line, at

$$e = e^* \equiv (I - \underline{A})^{-1}C\underline{x},$$

corresponding to steady-state values

$$\pi^* = i^* = -\bar{r} < 0, \quad y^* = -(1 - \beta)\bar{r}/\kappa < 0.$$

Because $I - \underline{A}$ has two real eigenvalues, one positive and one negative, trajectories of (3.26) converge to e^* only from initial conditions along the line SM in the figure, the one-dimensional stable manifold. Trajectories above and to the right of this line eventually converge to the zero-inflation steady state, while those below and to the left of it diverge from e^* in the opposite direction, eventually being drawn into (and remaining forever in) the shaded region. Because the actual dynamics of inflation are stochastic (even for arbitrarily large t) rather than precisely equal to the approximating ODE dynamics, there is actually zero probability of convergence of the learning dynamics to the REE represented by e^* , even from initial conditions on the line SM ; the learning dynamics must diverge from e^* in one direction or the other.⁴⁶

One might think that the non-learnability of the deflationary REE (while the learning dynamics are instead locally convergent near the REE consistent with the central bank's inflation target) implies that one need not be concerned about the possibility of falling into a self-fulfilling “deflation trap” of the kind stressed by Benhabib *et al.*, on the basis of the REE analysis. But the divergent dynamics near the

⁴⁶Figure 2 of Evans and Honkapohja (2010) is qualitatively similar, though plotted in the plane of inflation and output expectations $(\hat{\pi}_t, \hat{y}_t)$. These authors obtain an autonomous differential equation system in the $\hat{\pi} - \hat{y}$ plane only by assuming that interest-rate forecasts are obtained from inflation and output forecasts using people's knowledge of the policy rule. This is not consistent with the assumption made here that interest rates (like all other variables) are forecasted on the basis of past observations of that variable. In particular, when trajectories in Figure 1 cross the ZLB line, the nominal interest rate becomes positive, and under the learning rule assumed here, a positive nominal interest rate must be expected in the future as well. Under the Evans-Honkapohja forecasting assumption, instead, people would continue for some time to forecast a zero nominal interest rate into the indefinite future, because inflation and output *expectations* (which lag behind actual inflation and output) would still be at levels that would imply an expectation that the zero lower bound should continue (indefinitely) to bind.

deflationary REE include the existence of trajectories that diverge in the direction of ever-lower levels of inflation and output (those in the shaded region of the figure),⁴⁷ as a result of which the learning dynamics do imply the possibility of a “deflation trap,” albeit not one that involves convergence to the deflationary REE emphasized by Benhabib *et al.* As Evans and Honkapohja (2010) and Benhabib *et al.* (2012) discuss, to the extent that expectations are necessarily formed in this backward-looking way, the only way out of such a “trap” is to use other tools of policy (such as fiscal stimulus⁴⁸) to raise inflation and/or output long enough for inflation and output expectations to return to the region in which the learning dynamics can be expected to converge toward the target REE without further artificial support. In this view, other tools of policy have an important stabilization role to play in deep crises, even if monetary policy alone suffices as a stabilization tool *except* when unusual shocks drive expectations far enough away from the target REE forecasting rule.

3.3 Learning Dynamics as a Source of Persistence

Much of the early literature on TE dynamics with learning was concerned with the question of asymptotic convergence to an REE; the positive prediction of interest was whether an REE (or which REE) should be reached, and hence observed in practice. But the learning dynamics themselves might also be regarded as a source of positive predictions. One such positive prediction of particular interest is the existence of persistent fluctuations resulting from the dynamics induced by evolving estimates of the coefficients of people’s forecasting rules.

⁴⁷In fact, one can show that all trajectories that begin in the region that is below both the ZLB line and the SM line converge eventually to the shaded region, where they remain forever, and diverge farther and farther from the deflationary steady state.

⁴⁸Even if fiscal policy expectations are Ricardian, and the forecasts e_t are purely backward-looking as assumed above, an increase in government purchases should increase output and inflation, by increasing the term g_t in (1.16). If a current increase in net government transfers does not reduce the present value of forecasted future net transfers — for example, because future primary surpluses are forecasted using an estimator like (3.14) — then an increase in net transfers will also increase output and inflation, by increasing the term b_{t+1} in (1.18), while (if anything) also increasing the forecasts $\hat{E}^i v_{t+1}^i$.

3.3.1 Constant-Gain Learning

To study the macroeconomic dynamics that result from learning, it is convenient to assume that the gain γ_t in (3.15) takes some constant value $0 < \gamma < 1$ for all t , so that convergence to the REE never occurs, even asymptotically. A “constant-gain” learning algorithm of this kind may be justified as making sense if people believe that the coefficients of the correct forecasting model may shift over time, and consequently place more weight on the most recent observations in their estimates of the current coefficients.⁴⁹ In this case, the predicted TE dynamics are time-invariant, and can be characterized in terms of predicted unconditional moments (variances, autocovariances, etc.).

The dynamics of forecasts are again given by (3.21), but now with the constant value γ substituted for γ_t ; the implied TE dynamics of other variables are then given by (3.19).⁵⁰ Equation (3.21) can alternatively be written in the form

$$e_t = \Lambda e_{t-1} + \lambda \xi_t, \quad (3.27)$$

where

$$\Lambda \equiv (1 - \gamma)[I - \gamma A]^{-1}, \quad \lambda \equiv \gamma[I - \gamma A]^{-1}a.$$

If policy satisfies (1.24), both eigenvalues of A are inside the unit circle, so that $[I - \gamma A]$ is invertible, Λ and λ are well-defined, and both eigenvalues of Λ are also inside the unit circle. Hence (3.27) defines stationary dynamics for the forecasts $\{e_t\}$ and consequently for the variables $\{\mathbf{x}_t\}$ as well.

Equation (3.21) implies that the forecasts will be serially correlated, even if the disturbances are i.i.d. More precisely, it implies that each element of e_t will be a linear combination of two first-order autoregressive processes (the innovations in which are generally correlated), with coefficients of serial correlation equal to the

⁴⁹See Sargent (1993) or Evans and Honkapohja (2001) for further discussion.

⁵⁰The type of adaptive learning dynamics considered here are again of a fairly simple kind, since people’s forecasting rules are assumed simply to forecast constant future values for each of the variables, and the only coefficients that must be learned are the estimated means of each variable. However, as Eusepi and Preston (2012b) discuss, even if one allows updating of the slope coefficients of a more complex linear regression model, in a local linear approximation to the implied TE dynamics, linearizing around the REE steady state, there are no additional dynamics resulting from the updating of the slope coefficients; the updating of the additional coefficients has only second-order effects on the TE dynamics.

two eigenvalues of Λ . If γ is small (estimates are based on a fairly long history), the eigenvalues of Λ will be near 1, and these processes will be highly persistent. It then follows from (3.19) that fluctuations in inflation and the output gap will have highly persistent components under the TE dynamics with learning. This contrasts sharply with the prediction of the REE analysis, according to which inflation and the output gap should both be serially uncorrelated if all fundamental disturbances are, as a consequence of (1.25).

If the fundamental disturbances are instead themselves serially correlated, then persistent fluctuations in inflation and output are possible even under the REE dynamics. However, empirical New Keynesian models, such as those of Christiano *et al.* (2005) or Smets and Wouters (2007), generally find it necessary to introduce additional sources of persistence (indexation of prices and wages to past inflation, adjustment costs for expenditure), of debatable microeconomic realism, in order to fit the kind of persistence that is actually observed.⁵¹ Learning dynamics provide an alternative potential source of intrinsic persistence, and some studies (*e.g.*, Milani, 2005, 2007, 2011; Slobodyan and Wouters, 2009) find that there is less need for *ad hoc* structural persistence in econometric models that assume least-squares learning rather than rational expectations.⁵²

3.3.2 Consequences for Policy Evaluation

The additional dynamics resulting from learning can change one's conclusions regarding the relative desirability of alternative monetary policy rules, even with respect to comparisons among rules that do not imply explosive learning dynamics. As a simple example, suppose that the central bank's short-run inflation target depends linearly on the cost-push shock,

$$\pi_t^* = \phi_u u_t, \tag{3.28}$$

⁵¹See Woodford (2003, chap. 5) for discussion of the reasons for this.

⁵²Eusepi and Preston (2011) similarly find that the introduction of learning dynamics introduces a new channel for the propagation of the effects of technology shocks in an otherwise standard real business cycle model, and argue that the model with learning produces fluctuations more similar to observed business cycles. Even larger departures from REE business-cycle dynamics are predicted in the case of a model of learning that involves discrete switching between simple forecasting models of dramatically different character, as proposed by DeGrauwe (2010).

for some $0 \leq \phi_u \leq 1$, and i_t is adjusted each period as necessary in order to ensure that $\pi_t = \pi_t^*$. The required interest rate can be determined, as a function of current disturbances and expectations, from equations (1.16) and (1.17); these equations also indicate the implied evolution of the output gap.

Under the further assumption of RE beliefs, the model predicts that

$$y_t = -(1 - \phi_u)\kappa^{-1}u_t, \quad (3.29)$$

and hence that the unconditional variances of inflation and of the output gap will be

$$\text{var}(\pi) = \phi_u^2 \sigma_u^2, \quad \text{var}(y) = (1 - \phi_u)^2 \kappa^{-2} \sigma_u^2,$$

where σ_u^2 is the variance of the cost-push shock. It follows that for all ϕ_u in this interval, increasing ϕ_u increases the volatility of equilibrium inflation, but *reduces* the volatility of the equilibrium output gap. If policy is concerned to minimize some weighted average of the two variances, the optimal choice of ϕ_u will be somewhere between the two extremes, at a point that depends on the relative weight on the two stabilization objectives.

If we instead assume adaptive learning of the kind specified by (3.15), substitution of the policy rule (3.28) into the TE relation (1.17) implies that the output gap each period will be given by

$$y_t = -(1 - \phi_u)\kappa^{-1}u_t - (1 - \alpha)\beta\kappa^{-1}\hat{p}_t^*, \quad (3.30)$$

where \hat{p}_t^* is the common forecast at date t of the value of p_{t+1}^* . The latter forecast will be given by (3.17); if the estimates of the means of each of the variables evolve in accordance with (3.15), for some $0 < \gamma < 1$, this implies that

$$\hat{p}_t^* = (1 - \gamma)\hat{p}_{t-1}^* + \gamma[(1 - \alpha\beta)^{-1}\pi_t + \xi y_t + \mu_t].$$

Substituting (3.28) and (3.30) for π_t and y_t respectively in this expression yields a law of motion for the forecast of the form

$$\hat{p}_t^* = \rho\hat{p}_{t-1}^* + \phi_u\psi u_t, \quad (3.31)$$

where

$$0 < \rho \equiv \frac{(1 - \gamma)(1 - \alpha\beta)}{1 - (1 - \gamma)\alpha\beta} < 1, \quad \psi \equiv \frac{\gamma}{\alpha[1 - (1 - \gamma)\alpha\beta]} > 0$$

both are independent of the choice of ϕ_u .

Equation (3.31) implies that if $\phi_u > 0$, a positive cost-push shock immediately raises the forecast \hat{p}_t^* , and the forecast continues to be higher in subsequent periods as well (to an extent that decreases exponentially over time). Comparing (3.30) with the REE prediction (3.29), we see that with adaptive learning, the output reduction in the period of the shock is greater than would occur under rational expectations; moreover, the negative effect on the output gap persists, rather than being limited to the period of the shock. For both reasons, a given value of $\phi_u > 0$ does not reduce the predicted variance of the output gap as much as is predicted by the REE analysis, though it continues to increase the predicted variance of inflation by the same amount. Thus the trade-off between inflation stabilization and output-gap stabilization is steeper in the case of learning: less reduction in the variance of the output gap is achieved by a given increase in the variance of inflation. The implication, as argued by Orphanides and Williams (2005), is that under given preferences with regard to inflation and output-gap stability, it will be optimal to choose a lower value of ϕ_u (maintaining tighter control of inflation) when one recognizes that people must learn to forecast macroeconomic conditions, relative to what one would conclude from the REE analysis.⁵³

4 TE Dynamics with Nearly Correct Beliefs

There is, however, another way of ensuring that one's model's predictions do not depend on a supposition that people will fail to notice patterns in the data that should actually be easily discerned. These alternative approaches are based not on an explicit specification of the procedure used to look for such patterns, as in the case of econometric learning models, but rather on a direct requirement that probability beliefs — however obtained — not be too different from the true probabilities (according to one's model). Approaches of this kind propose no model of how people reason to the probability beliefs that they hold, but instead focus on defining the respects in which subjective beliefs should reasonably be expected to be similar to

⁵³Orphanides and Williams consider a one-parameter family of policies similar to the one considered here, but in the context of a simpler model of the way in which expectations affect aggregate supply. For implications of learning dynamics for the optimal choice of a policy rule within more complex families of candidate policies, see Gaspar *et al.* (2011).

objective probabilities, and the other respects in which one might expect more variation in subjective beliefs. In this section, I discuss two examples of how this might be done: the “rational belief equilibria” of Mordecai Kurz and coauthors (Kurz, 1994, 1997, 2012; Kurz and Motolesse, 2011), and the “near-rational expectations” proposed by Woodford (2010) and explored further in Adam and Woodford (2012).

Before describing these approaches, it is important to note that a hypothesis that beliefs are “nearly correct” does *not* imply that they are nearly the same as (any possible) REE beliefs. The extent to which beliefs are correct depends on their conformity with the actual TE dynamics, which may differ greatly from REE dynamics, and not their conformity with REE predictions. This difference is emphasized in particular by Kurz (2012), who emphasizes the possibility of sizeable aggregate fluctuations even when the magnitude of exogenous disturbances to “fundamentals” is much smaller than must be postulated to account for the fluctuations using DSGE models that assume rational expectations.

4.1 “Rational Belief Equilibria”

Kurz (1994) proposes a relaxation of the rational expectations hypothesis in which the probability beliefs of decisionmakers are required to imply model-consistent values for some data moments, but not for all of the data moments that are relevant to their forecasts and hence to their decisions. Certain quantities (including conventional macroeconomic aggregates, such as rate of growth of GDP or the CPI) are assumed to be objectively measurable, and as a consequence everyone is assumed to agree about the current and past values of these variables. The postulate of “rational beliefs” (RB) then requires that in any stationary equilibrium (a “rational belief equilibrium,” RBE) consistent with some time-invariant policy, everyone must also agree about all of the unconditional first and second moments⁵⁴ of these objectively measurable variables, and assign values to these moments that coincide with the predictions of the model about this particular RBE.⁵⁵

⁵⁴By “all second moments” I mean to include all covariances between leads and lags of the various variables.

⁵⁵In fact, Kurz (1994) proposes the stronger postulate that the subjective assessment of the unconditional joint distribution of the objectively measurable variables must coincide with their model-implied unconditional distribution. In the case of a linear model with additive Gaussian disturbances, of the kind used in applications such as Kurz (2012), and in the example presented

But these variables are not the only ones on the basis of which individuals form their forecasts; there are also subjective variables (“belief states”) about which they need not agree. A given decisionmaker is assumed to have coherent probability beliefs about the joint distribution of her own belief states and the objectively measurable variables, on the basis of which the belief states modify her forecasts of the future paths of the objectively measurable variables; but *these* data moments need not be ones about which others agree, and the probability beliefs of an individual need not coincide in this respect with the predictions of the model. It is in this latter respect that the RB postulate is weaker than RE. Insofar as people are assumed to learn the correct values of some data moments but not others, the RBE concept is a cousin of the RPE concept discussed in section 3.1.

The content of the RB postulate, as well as the sense in which it is weaker than RE, is best illustrated using an example. Suppose that the natural rate of output is the sum of two components,

$$Y_t^n = \bar{Y}_t^n + \xi_{2t}, \quad (4.1)$$

where the permanent component \bar{Y}_t^n evolves as a random walk,

$$\bar{Y}_t^n = \bar{Y}_{t-1}^n + \xi_{1t}, \quad (4.2)$$

with $\{\xi_{1t}\}$ an i.i.d. innovation distributed as $N(0, \sigma_1^2)$, and the transitory component ξ_{2t} is another i.i.d. innovation, distributed as $N(0, \sigma_2^2)$ and independent of the permanent shocks. If the process $\{Y_t^n\}$ is objectively measurable but its permanent and transitory components are not, and no other objectively measurable variables provide information about this decomposition, then an optimal estimate of the permanent component (or optimal forecast of the long-run level) of Y_t^n at any time t is given by an exponentially-weighted moving average⁵⁶

$$\bar{Y}_t \equiv (1 - \lambda) \sum_{j=0}^{\infty} \lambda^j Y_{t-j}^n,$$

where the smoothing factor $0 < \lambda < 1$ is given by

$$\lambda = \frac{2}{2 + q + \sqrt{q^2 + 4q}}, \quad q \equiv \frac{\sigma_1^2}{\sigma_2^2}.$$

below, identity of the two unconditional distributions is equivalent to identity of the complete list of first and second moments.

⁵⁶This corresponds to the Bayesian posterior mean, or minimum-mean-squared-error forecast, using a Kalman filter (Harvey, 1989, chap. 4), as originally derived by Muth (1960).

Conditional only on objectively measurable data, then, an optimal forecast of the future natural rate at any horizon $k \geq 0$ will be given by

$$\bar{E}_t Y_{t+k}^n = \bar{Y}_t. \quad (4.3)$$

Suppose, however, that in addition to the objectively measurable data, each individual price-setter j has a *subjective estimate* of the permanent component, that I shall denote z_t^j . If each price-setter correctly understands the laws of motion (4.1)–(4.2), this implies that subjective forecasts will be given by

$$\hat{E}_t^j Y_{t+k}^n = z_t^j, \quad (4.4)$$

rather than by (4.3). Note that each individual’s beliefs are described by a completely specified, internally consistent probability measure, that is moreover consistent with the true first and second moments of all objectively measurable data; for example, these beliefs imply correct (model-consistent) values for the unconditional moments $E[\Delta Y_t^n]$ and $\text{cov}(\Delta Y_t^n, \Delta Y_{t-k}^n)$ for all k , as these can be derived from (4.1)–(4.2).⁵⁷

But individual beliefs about the statistics of the subjective belief state z_t^j and its co-movement with objectively measurable data need *not* coincide with the beliefs of others, or with the way that the model describes the evolution of these variables. In our example, while each individual j believes that $z_t^j = \bar{Y}_t^n$, it does not follow that z_t^j must take the same numerical value for all j . Moreover, even if, as in Kurz and Motolesse (2011) and Kurz (2012), we suppose that the population distribution of subjective beliefs, and hence the population mean $Z_t \equiv \int z_t^j dj$, are objectively measurable data, it does not follow that z_t^j must equal Z_t for each individual. Individuals can be aware that their personal estimate z_t^j differs from the average estimate, without any internal inconsistency of their beliefs. The RB postulate requires that people all have model-consistent beliefs about unconditional moments such as $E[Y_t^n - Z_t]$, $E[\Delta Z_t]$, $\text{cov}(\Delta Y_{t+k}^n, Y_t^n - Z_t)$, and so on. But awareness of these moments and observation of Z_t does not give a price-setter any reason to doubt the validity of her forecast (4.4), given her belief in the laws of motion (4.1)–(4.2) and her belief

⁵⁷Kurz does not refer to these common beliefs as “correct” beliefs about unconditional moments, but only as “empirical frequencies.” However, in applications such as Kurz and Motolesse (2011) or Kurz (2012), the calculations used to explain or predict data are carried out under the assumption that the empirical frequencies correspond to model-implied unconditional moments, under time-invariant stochastic processes for the various disturbances specified in the model.

that z_t^j is an accurate (but personal) observation of the value of \bar{Y}_t^n . (This simply requires each individual to believe that others' personal assessments of the value of the permanent component are erroneous, even though she understands that, like her, they each believe their personal assessments to be correct.⁵⁸)

As one possible example of how this makes possible an additional source of aggregate fluctuations, suppose that people's subjective assessments are given by $z_t^j = \bar{Y}_t + \nu_t^j$, where ν_t^j is a random term that evolves independently of all "fundamental" variables, including both the permanent and transitory components of Y_t^n . Thus, since \bar{Y}_t is objectively measurable, the subjective state z_t^j reflects no additional information about the future evolution of the natural rate (or any other fundamentals). Moreover, suppose that the errors ν_t^j are correlated across individuals, so that the aggregate error $\nu_t \equiv \int \nu_t^j dj$ is not equal to zero. Then because ν_t represents an error in the average estimate of a variable that is relevant for pricing decisions (the error in the average estimate $\int \hat{E}_t^j [\bar{Y}_t^n - \bar{Y}_t] dj$), it will affect the determination of endogenous aggregate variables, such as output and inflation; and variation in ν_t will be an additional source of variability in these variables, in addition to the random variation in fundamentals such as $\{\xi_{1t}, \xi_{2t}\}$.

To illustrate the effects of fluctuations in the aggregate belief state on endogenous variables, let the monetary policy rule be specified by a *target criterion*: that is, the central bank adjusts its instrument as necessary in order to ensure that the linear relationship

$$\pi_t + \phi(Y_t - \bar{Y}_t) = 0 \tag{4.5}$$

holds at all times.⁵⁹ This represents a form of "flexible inflation targeting" (Svensson, 1999), where the concept of the output gap in the central bank's target criterion

⁵⁸Thus this equilibrium concept allows a much wider range of possible specifications of belief dynamics than a rational-expectations model with "private information," of the kind considered by Rondina and Walker (2012). In the latter paper, all individuals are assumed to agree about the joint distribution of all publicly *or privately* observable variables, though individuals do not observe other individuals' private signals about the separate components of the aggregate disturbance. In Kurz's work, instead, people "agree to disagree." It is therefore not necessary to suppose that there is anything secret about individuals' subjective beliefs; it is only the basis for accepting these subjective assessments as *correct* that is not shared.

⁵⁹The state-contingent path for the interest rate i_t required in order for (4.5) to hold each period will depend on the subjective expectations of both price-setters and consumers. I suppose here that the central bank observes average expectations when setting i_t , and so can implement a reaction function that makes (4.5) a necessary consequence of the TE relations that determine π_t and y_t ,

is output relative to the central bank's estimate of *long-run* potential, rather than relative to the current natural rate of output. As a simple example in which belief fluctuations provide an independent source of aggregate variability, suppose that ν_t evolves as an AR(1) process,

$$\nu_t = \rho\nu_{t-1} + \epsilon_t, \quad (4.6)$$

where $0 \geq \rho < 1$, and $\{\epsilon_t\}$ is an i.i.d. innovation, with distribution $N(0, \sigma_\epsilon^2)$ independent of all fundamental states. We can then solve for the equilibrium dynamics of inflation and output implied by the TE relations, (4.5) and the above assumptions about subjective expectations, using the method of undetermined coefficients.

Let us conjecture beliefs on the part of each price-setter j of the form

$$\hat{E}_t^j p_{t+1}^{*j} = \gamma_1(\bar{Y}_t - z_t^j) + \gamma_2(\bar{Y}_t - Z_t), \quad (4.7)$$

for coefficients γ_1, γ_2 that remain to be determined. Average beliefs are then given by

$$\int \hat{E}_t^j p_{t+1}^{*j} dj = \gamma(\bar{Y}_t - Z_t),$$

where $\gamma \equiv \gamma_1 + \gamma_2$. Let us also suppose for simplicity that the cost-push shock u_t is equal to zero at all times.⁶⁰ The TE relation (1.17) then implies that inflation and the output gap must be given by

$$\pi_t = \frac{\kappa\phi}{\kappa + \phi}(\bar{Y}_t - Y_t^n) - \frac{(1 - \alpha)\beta\gamma\phi}{\kappa + \phi}\nu_t, \quad (4.8)$$

$$y_t = \frac{\phi}{\kappa + \phi}(\bar{Y}_t - Y_t^n) + \frac{(1 - \alpha)\beta\gamma}{\kappa + \phi}\nu_t. \quad (4.9)$$

Since equations (4.8)–(4.9) are relationships among objectively measurable variables,⁶¹ the RB postulate requires that the subjective probability beliefs of each price-setter be consistent with them. These relations, together with the laws of motion (4.1)–(4.2) and (4.6) for the exogenous aggregate state variables, further imply

regardless of what those subjective expectations may be. Thus (4.5) can be treated as an equilibrium relation in solving for the equilibrium dynamics under a given hypothesis about expectations.

⁶⁰Note that even under this assumption, the model implies the existence of equilibrium fluctuations in inflation and in the output gap, owing to the discrepancy between the concept of potential output (\bar{Y}_t) used in the central bank's target criterion (4.5) and the one (Y_t^n) that shifts the AS relation (1.17).

⁶¹Recall that $\nu_t = Z_t - \bar{Y}_t$, and the average belief state Z_t is assumed to be objectively measurable.

that correct forecasts of future inflation and output are given by

$$E_t \pi_{t+k} = \frac{\kappa \phi}{\kappa + \phi} \lambda^k (\bar{Y}_t - \bar{Y}_t^n) - \frac{(1 - \alpha) \beta \gamma \phi}{\kappa + \phi} \rho^k \nu_t,$$

$$E_t y_{t+k} = \frac{\phi}{\kappa + \phi} \lambda^k (\bar{Y}_t - \bar{Y}_t^n) + \frac{(1 - \alpha) \beta \gamma}{\kappa + \phi} \rho^k \nu_t$$

for any horizon $k \geq 1$, where $E_t[\cdot]$ refers to the expectation conditional on the history of all exogenous states up through period t , including the (unobserved) value of \bar{Y}_t^n . If price-setters are assumed to correctly understand these laws of motion,⁶² then their subjective forecasts of future inflation and output gaps must conform to these equations as well, but with the value of \bar{Y}_t^n replaced by each individual's subjective estimate of this state. Thus for any price-setter j , each of the forecasts is a linear function of $\bar{Y}_t - z_t^j$ and $\nu_t \equiv Z_t - \bar{Y}_t$.⁶³

Substituting these subjective forecasts into definition (1.12), we can obtain an expression for $\hat{E}_t^j p_{t+1}^{*j}$ as a linear function of $\bar{Y}_t - z_t^j$ and $\bar{Y}_t - Z_t$, as conjectured in (4.7). We now however have expressions for the coefficients γ_1, γ_2 (given in the Appendix) as functions of the assumed value of γ . Requiring the implied values of these coefficients to equal their conjectured values yields two linear equations to solve for the unknown coefficients γ_1, γ_2 . As shown in the Appendix, our sign assumptions on parameters imply the existence of a unique solution, with $\gamma > 0$.

We thus obtain TE dynamics consistent with the RB postulate in which fluctuations in the aggregate belief state ν_t cause random variations in inflation and output.

⁶²The RB postulate requires that price-setters all correctly understand the autocorrelation function of the objectively measurable process $\{\nu_t\}$, but it does not require that they agree that an unbiased forecast of ν_{t+k} at time t depends only on the current value ν_t ; they may have subjective judgments about the likely future path of the aggregate belief state that they believe are more accurate than the forecast that could be made on the basis of objectively measurable data alone. Here I make the more restrictive assumption that no one believes that they have additional insight into the future evolution of any exogenous states *except* for believing in their personal estimates of the permanent component \bar{Y}_t^n .

⁶³Kurz and Motolese (2011) say that “those who believe the economy is stationary” will necessarily forecast using the “empirical measure” — that is, using only the information contained in the history of objectively measurable variables, and so make forecasts such as (4.3). But in fact, the subjective probability beliefs specified here imply that $\{\pi_t, y_t, \bar{Y}_t - \bar{Y}_t^n, Y_t^n - \bar{Y}_t^n, \bar{Y}_t - Z_t, \Delta \bar{Y}_t^n\}$ are jointly stationary processes; and the same is true of the RBE beliefs specified in the applications proposed by Kurz and Motolese (2011) and Kurz (2012). The crucial issue is actually not stationarity, but whether variables other than objectively measurable variables are also used in forecasting.

It is instructive to compare this solution with the REE dynamics under policy rule (4.5). We may assume as above that each individual observes a personal state variable z_t^j (a “gut feeling,” if one likes), that is distributed as assumed above; but under the RE hypothesis, each individual must correctly understand the joint distribution of z_t^j and all other variables. This would mean correctly understanding that z_t^j contains no information that is useful for predicting the future path of the natural rate of output (given that \bar{Y}_t is independently observable), and similarly that Z_t is uninformative. RE forecasts of all variables would then correspond simply to the expectations of those variables conditional on the observed history of the natural rate of output; thus, for example, the common forecast of the future natural rate of output would be given by (4.3). It is shown in the Appendix that under policy rule (4.5), the unique stationary REE is one in which inflation and the output gap are given by equations (4.8)–(4.9), but with $\gamma = 0$, whereas $\gamma > 0$ in the RBE discussed above. Thus the RBE beliefs do not change the response of inflation or output to exogenous fluctuations in the natural rate of output, but result in increased variability of both inflation and the output gap for any value of ϕ , relative to the REE prediction, owing to the existence of fluctuations unrelated to any changes in fundamentals, but due purely to variation in the aggregate belief state.

The above simple calculation may make it appear that the RBE hypothesis makes definite quantitative predictions about the evolution of endogenous variables under a given policy rule, but this is actually not true; the RBE constructed above is only one possible example of TE dynamics consistent with the RB postulate under the assumed policy rule. First of all, there is an RBE of the kind assumed above for any specification of the serial correlation coefficient ρ and of the innovation variance σ_ϵ^2 for the process $\{\nu_t\}$. Moreover, there is no reason why $\{\nu_t\}$ must be an AR(1) process; this allowed us to verify the conjecture that subjective forecasts were of the form (4.7), but we might equally well have assumed a more complex process for $\{\nu_t\}$, and still solved for an RBE, in which however, subjective forecasts would be correspondingly more complex. Thus if we allow $\{\nu_t\}$ to be any process in some larger parametric family, we can obtain a multi-parameter family of RBE associated with the given policy (4.5). But even this understates the multiplicity of possible RBE. For it was not necessary to have assumed that people believe that they have an additional (personal) awareness of the decomposition of Y_t^n into permanent and transitory components, but no additional personal insight into the economy’s future

evolution of any *other* sort. Allowing for other types of subjective beliefs (that need not be correlated with actual outcomes, according to one's model, in the way that people believe that they are) would further expand the set of RBE solutions consistent with a given policy rule.

Kurz and coauthors argue that the more flexible relationship between the evolution of exogenous fundamentals and that of endogenous variables allowed by this relaxation of the RE hypothesis can make sense of some of the empirical difficulties faced by RE models. For example, Kurz and Motolese (2011) discuss RBE of an asset-pricing model in which there is a risky asset in fixed supply and an exogenously given riskless rate of return (independent of the quantity invested in the riskless asset). The dividend on the risky asset is an exogenous process, about the future evolution of which individual investors believe they have additional personal insight, beyond the information contained in the past history of the dividend, just as in the case of subjective forecasts of the natural rate of output in the above example. In the RBE, variations in the aggregate belief state become an additional source of variation in equilibrium asset prices, and in particular result in a time-varying risk premium, of the kind that is found to be empirically important in many asset markets.

Kurz and Motolese estimate the parameters of their model using data on term premia associated with federal funds futures and Treasury bills, and find that allowance for the more flexible class of equilibria allows the data to be fit better; their best-fitting RBE implies that more than half of the measured risk premia are due to fluctuations in the aggregate belief state. This suggests that the kind of additional flexibility allowed by the concept of an RBE may be of empirical relevance. At the same time, because the predictions of the more general theory are much less specific, it is not obvious that findings such as those of Kurz and Motolese can be regarded as confirming a specific theoretical view of the nature of subjective beliefs.⁶⁴

Relaxation of the RE hypothesis also has potential consequences for policy design;

⁶⁴While the sets of possible RBE discussed in papers such as Kurz and Motolese (2011) and Kurz (2012) involve only a few free parameters, this is not because the RB postulate alone allows one to derive such specific conclusions — a large number of additional (theoretically unmotivated) assumptions are made as well, in order to obtain equilibria of a particular form. Moreover, in neither of these applied papers are all of the restrictions implied by the RB postulate imposed; the solutions proposed as possible accounts of actual data are actually examples of an even weaker version of the RBE concept, and the proposed restrictions on the stochastic processes characterizing the data are mainly adopted for convenience rather than following from a conception of rationality.

as illustrated by the above example, the degree of macroeconomic stability guaranteed by commitment to a given policy rule may not be as great as a mere analysis of the REE dynamics consistent with it would suggest. This raises the possibility that alternative rules might provide more robust approaches to stabilization, even if they do not lead to a superior REE. While there will not be a unique RBE consistent with a given policy rule, or even a unique RBE associated with a given restricted state space (as in the analysis of MSV-REE above), it may be compare the data moments associated with the entire range of possible RBE, for alternative parameterizations of a policy rule. Kurz (2012) undertakes an illustrative analysis of this kind of the consequences of alternative central-bank reaction functions in the context of a New Keynesian model closely related to the one presented here. However, the comparisons undertaken consider only certain parametric classes of RBE, and it is unclear why attention should be restricted to these specific types of equilibria. This seems an important limitation of the Kurz approach for purposes of policy analysis. The alternative approach presented next instead allows a clear delineation of the set of equilibria consistent with a given policy rule.

4.2 “Near-Rational Expectations”

Rather than distinguishing *a priori* between data moments that individuals should correctly assess and those that they may not, depending on the nature of the variable in question, the assumption of “near-rational expectations” in Woodford (2010) instead defines a set of probability beliefs that are close enough to the predictions of one’s model to be plausibly held by decisionmakers in such a situation, on purely statistical grounds. Essentially, an alternative probability distribution is “close” to the predicted probabilities of outcomes in a given equilibrium if the alternative distribution represents a sample distribution of outcomes that could be observed in some finite number of repetitions of the equilibrium. This requires, for example, that “near-rational” subjective expectations assign zero probability to all outcomes that occur with zero probability in equilibrium.

This means that each agent’s subjective probability measure over possible paths for all variables must be *absolutely continuous* with respect to the equilibrium probability measure. This in turn implies there must exist a scalar stochastic process $\{m_t\}$ for each agent — the agent’s *belief distortion factor* — with $m_t \geq 0$, $E_t m_{t+1} = 1$ at

all times, such that the agent’s subjective one-period-ahead forecast of any variable X_{t+1} is given by

$$\hat{E}_t X_{t+1} = E_t[m_{t+1} X_{t+1}], \quad (4.10)$$

where $E_t[\cdot]$ indicates the conditional expectation under the true (model-implied) probabilities in the particular equilibrium. Thus a value $m_{t+1} > 1$ in a particular state of the world at date $t + 1$ implies that, conditional on reaching the predecessor state at date t , the agent exaggerates the probability of reaching this state relative to the correct equilibrium probability. Internal consistency of individual probability beliefs then implies that longer-horizon subjective forecasts are correspondingly given by

$$\hat{E}_t X_{t+j} = E_t[m_{t+1} \cdots m_{t+j} X_{t+j}].$$

The degree of discrepancy between subjective and objective probability beliefs can then be measured by the degree to which the distortion factor $\{m_t\}$ differs from a constant factor, equal to 1 in all states. A measure of the degree of discrepancy in one-period-ahead beliefs (looking forward from any period t) with appealing properties is the *relative entropy* between the subjective and objective conditional probabilities

$$R_t \equiv E_t[m_{t+1} \log m_{t+1}]. \quad (4.11)$$

This is non-negative convex function of the belief distortion factor that achieves its minimum possible value of zero if and only if $m_{t+1} = 1$ almost surely (the RE case). Moreover, the probability of observing a sample frequency distribution for the possible outcomes at date $t + 1$ that is close to any given subjective measure, in the case of a large (but finite) number of independent draws from the equilibrium probability measure, is (in the case of a large enough number of draws) a decreasing function of the relative entropy of the subjective measure.⁶⁵ Hence subjective beliefs under which R_t is small (though positive) each period are ones that could plausibly be maintained even by an agent with considerable experience of typical equilibrium outcomes.

Woodford (2010) accordingly defines an equilibrium with “near-rational expectations” (NRE) as a situation in which each agent optimizes on the basis of internally consistent probability beliefs for which R_t is sufficiently small each period, when calculated with respect to the equilibrium probability measure describing the outcomes that result (in each possible state of the world) from their collective choices. Note

⁶⁵See, for example, Cover and Thomas (2006).

that the equilibrium measure will not generally be the REE measure, because people act on the basis of non-REE beliefs; hence NRE equilibrium outcomes need not be near the REE outcomes in order for beliefs to be “near-rational.”

In the context of the NK model described above, an NRE equilibrium (NREE) corresponds to stochastic processes $\{\bar{v}_t^i\}$ and distortion factors $\{m_t^i\}$ for each household, and processes $\{p_t^{*j}\}$ and distortion factors $\{m_t^j\}$ for each firm, such that (1.10) and (1.15) hold each period when subjective forecasts are given by (4.10) for each agent, and the distortion factors imply that the relative entropy for each agent remains within some bound. If, for example, we assume a monetary policy rule of the form (1.21) and restrict attention to the special case of common subjective probability beliefs for all agents, then an NREE corresponds to a vector stochastic process $\{z_t\}$ and distortion factor $\{m_t\}$, such that

$$z_t = B E_t[m_{t+1}z_{t+1}] + b\xi_t \quad (4.12)$$

holds each period (where B and b are again the matrices in (1.22)), and the relative entropy (4.11) implied by the distortion factor satisfies the specified bound.

4.3 Robustly Optimal Policy

For any positive upper bound on the allowable relative entropy, the set of NREE consistent with a given policy rule will be large. How, then, can this kind of theory provide a basis for selection of a particular policy rule? Woodford (2010) proposes that one choose the policy that implies the highest possible *lower bound* for one’s welfare objective (or lowest possible upper bound for one’s loss function), across the entire set of NREE consistent with the rule, under some specified bound on the allowable size of belief distortions. Such a “maximin” approach to policy choice is in the spirit of the “robust control” approach to dealing with model uncertainty advocated by Hansen and Sargent (2008).

This approach requires one to determine, for any candidate policy rule, the “worst-case” belief distortion process, that implies an NREE that is the worst possible for the policymaker’s welfare objective, subject to the bound on the size of belief distortions that are contemplated. As an example, suppose that the objective of policy is to minimize a discounted loss function of the form

$$E_0 \sum_{t=0}^{\infty} \beta^t [\pi_t^2 + \lambda(y_t - y^*)^2], \quad (4.13)$$

for some relative weight $\lambda > 0$ on output-gap stabilization, and some optimal output gap $y^* > 0$, as in Clarida *et al.* (1999). Here the expectation $E_0[\cdot]$ used to define the objective refers to the probability beliefs of the policymaker, which need not be shared by others. And let us again consider policy commitments in the simple family (3.28).

If for simplicity we restrict attention to equilibria in which belief distortions are common to all agents, (1.11) and (1.20) imply that

$$\pi_t = \kappa y_t + u_t + \beta E_t[m_{t+1}\pi_{t+1}]$$

each period, which is just the NRE generalization of the “New Keynesian Phillips curve” assumed by Clarida *et al.* Substituting (3.28) for the path of inflation, this implies that the output gap must satisfy

$$y_t = -\kappa^{-1} \{(1 - \phi_u)u_t + \beta\phi_u E_t[m_{t+1}u_{t+1}]\} \quad (4.14)$$

in the NREE corresponding to any distortion process $\{m_t\}$.

Since the path of inflation is independent of belief distortions under a policy commitment of the hypothesized type, the belief distortions that maximize (4.13) involve a choice of the one-period-ahead distortion factors $\{m_{t+1}\}$ looking forward from any date t so as to maximize $(y_t - y^*)^2$ subject to an upper bound

$$R_t \leq \bar{R}, \quad (4.15)$$

where R_t is defined in (4.11), and the size of $\bar{R} > 0$ indicates the allowable degree of departure from model-consistency. Hence the factors $\{m_{t+1}\}$ should be chosen to maximize

$$\frac{1}{2}(y_t - y^*)^2 + \theta_t E_t[m_{t+1} \log m_{t+1}]$$

subject to the constraint that $E_t m_{t+1} = 1$, where y_t is given by (4.14) and θ_t is a Lagrange multiplier associated with the constraint (4.15).

If u_{t+1} is i.i.d. $\mathcal{N}(0, \sigma_u^2)$, the solution to this problem is easily shown to involve a state-contingent distortion factor

$$\log m_{t+1} = \alpha + \gamma u_{t+1},$$

where

$$\alpha = -\bar{R}, \quad \gamma = \pm \frac{(2\bar{R})^{1/2}}{\sigma_u}.$$

The positive root for γ is optimal if $y_t < y^*$ (the most common case), while the negative root is optimal if $y_t > y^*$. (When $y_t < y^*$, the policymaker's tradeoff is made even more painful by an increase in inflation expectations, that shift the short-run Phillips curve in a way that increases the tension between the goals of keeping inflation near zero and the output gap near y^* ; and expected inflation is increased if people exaggerate the likelihood of positive cost-push shocks. If $y_t > y^*$, instead, the worst-case belief distortions would be ones that *reduce* inflation expectations, by exaggerating the likelihood of negative cost-push shocks.)

The worst-case beliefs then imply

$$E_t[m_{t+1}u_{t+1}] = \gamma\sigma_u^2 = \pm(2\bar{R})^{1/2}\sigma_u,$$

taking care to select the root that implies the largest gap between y_t and y^* . It follows that

$$|y_t - y^*| = |y^* + (1 - \phi_u)\kappa^{-1}u_t| + \phi_u\beta\kappa^{-1}(2\bar{R})^{1/2}\sigma_u \quad (4.16)$$

for all realizations of u_t . Equation (4.16) shows that increasing ϕ_u reduces the sensitivity of $y_t - y^*$ to cost-push shocks (as in the RE analysis), but at the cost of making it possible for the absolute value of the gap in the *absence* of any cost-push shock to be larger, as a result of belief distortions.

The upper bound for (4.13) in the case of any policy in the simple family (3.28) is then given by

$$(1 - \beta)^{-1}[L_\pi + \lambda L_y^{pess}],$$

where

$$L_\pi \equiv E[\pi_t^2] = \phi_u^2\sigma_u^2$$

is the same function of ϕ_u as in the RE analysis, and

$$L_y^{pess} \equiv E[(y_t - y^*)^2]$$

is evaluated under the worst-case belief distortions; in both expressions the expectation is over possible realizations of u_t . One observes that $\partial L_\pi / \partial \phi_u = 2\phi_u\sigma_u^2$, regardless of the degree of concern for robustness; but one finds that allowance for belief distortions ($\bar{R} > 0$) makes the value of the derivative $\partial L_y^{pess} / \partial \phi_u$ less negative (or more positive) at each value of ϕ_u . Hence the value of ϕ_u at which the marginal reduction in expected losses with respect to output-gap stabilization no longer outweighs the marginal increase in expected losses with respect to inflation stabilization will be

reached at a lower value of ϕ_u under the worst-case beliefs (when $\bar{R} > 0$) than under the RE analysis. In fact, the upper bound for expected losses may be minimized at $\phi_u = 0$, whereas complete inflation stabilization is never optimal under the RE analysis; and the optimal ϕ_u remains bounded away from 1, no matter how large the weight λ on the output-gap stabilization objective may be, whereas the optimal ϕ_u approaches 1 as $\lambda \rightarrow \infty$ in the RE analysis.

The “robustly optimal” policy within this simple family thus involves greater stability of inflation in the face of cost-push shocks than would be optimal if one could be sure that people would have model-consistent expectations. This is similar to the conclusion obtained in section 3.3.2 when analyzing alternative policies under the assumption of adaptive learning, and again the basic reason is that variations of inflation in response to cost-push shocks make it too easy for people to mis-estimate the average future rate of inflation, causing undesirable instability in the short-run Phillips curve tradeoff.

The reasons for inflation expectations to be insufficiently well-anchored are somewhat different in the two cases: in the learning analysis, it was assumed that inflation expectations *necessarily* drift in response to certain observations of inflation outcomes, and a large value of ϕ_u was dangerous because it increased the frequency of occurrence of observations that would lead to significant expectational errors; here, instead, no precise prediction is made about what expectations must be, but a large value of ϕ_u is dangerous because it allows more significant expectational errors to be consistent with the assumed bound on relative entropy. Yet ultimately, the problematic feature of the large- ϕ_u policy is the same in both cases: it makes sample paths in which average observed inflation differs significantly from the actual long-run inflation target (in particular, paths in which the sample average is significantly *higher*) occur too frequently.

Woodford (2010) extends the above analysis by considering a much more flexible family of policies, in which the short-term inflation target π_t^* is an arbitrary linear function of the history of cost-push shocks, and shows how the optimal commitment of this form differs from the optimal commitment in the RE analysis of Clarida *et al.* (1999). As in the simpler exercise above, the robustly optimal policy commitment involves a lower amplitude of inflation surprises in response to cost-push shocks; Woodford shows that it also involves a *greater* degree of commitment to subsequent

reversal of any effects on the price level of past cost-push shocks.⁶⁶ Kwon and Miao (2012) show how a similar method can be used to characterize the robustly optimal policy commitment for a broad class of linear-quadratic policy problems, and generalize the results of Woodford (2010) to the cases of persistent cost-push shocks and of a more general form of aggregate-supply relation that incorporates intrinsic inflation inertia.

Adam and Woodford (2012) further extend the analysis of Woodford (2010), considering policy commitments that are not necessarily expressed in terms of inflation targets that depend only on the history of exogenous disturbances. They find that the conclusions mentioned above continue to hold, as a description of how inflation must be expected to evolve in response to cost-push shocks *under the worst-case beliefs*, even if a robustly optimal policy commitment (within the more general family) need not require inflation to evolve this way *regardless* of the nature of belief distortions. In the analysis of Adam and Woodford, there is not a uniquely defined policy rule that is robustly optimal; instead, there exists a large class of policy rules that all imply the same dynamics under the worst-case belief distortions, and hence achieve the same minimum upper bound for the loss function, though they may be associated with different TE dynamics under other kinds of distorted beliefs that are also consistent with the relative-entropy bound.⁶⁷

Among the robustly optimal policy rules is one that involves commitment to a target criterion: the central bank uses its policy instrument to ensure that the joint evolution of inflation and output satisfy a linear relationship of the form

$$\pi_t + \phi_s(\pi_t - E_{t-1}\pi_t) + \phi_y(y_t - y_{t-1}) = 0 \quad (4.17)$$

⁶⁶This kind of robust policy problem is compared to a alternative ways of introducing robustness to uncertainty about the correctness of model equations into an optimal monetary stabilization policy problem, in the context of the same linear-quadratic New Keynesian framework used here, in Hansen and Sargent (2012).

⁶⁷It should be recalled that also under the RE analysis, the optimal policy commitment is not uniquely defined. Instead, the *optimal REE dynamics* are uniquely defined, while there are many different policy rules that can achieve these dynamics as a determinate equilibrium outcome; the rules differ in the behavior that they prescribe out of equilibrium, though the policy instrument evolves in the same way *in equilibrium* under each of them. Under the robust policy analysis, the different robustly optimal rules also differ in the sets of possible equilibrium outcomes associated with them, since a given rule does not imply a determinate equilibrium, except under a particular specification of the belief distortions.

each period, where ϕ_s, ϕ_y are both positive coefficients, that depend both on model parameters and on the relative weight λ assumed in the objective (4.13),⁶⁸ and $E_{t-1}\pi_t$ indicates the *policymaker's* forecast of inflation a period earlier. Here the coefficient $\phi_s > 0$ multiplying the inflation surprise results from the concern for robustness, and this coefficient is larger the greater the concern for robustness (as measured by the relative-entropy bound). The presence of this term reduces the extent to which a “cost-push shock” should be allowed to cause a surprise change in the rate of inflation, since it requires the surprise reduction in the output gap to be $(1 + \phi_s)/\phi_y$ times as large as the surprise increase in inflation, rather than only $1/\phi_y$ times as large, as under the optimal commitment assuming rational expectations (Woodford, 2003, chap. 7). Thus one again concludes, as in the analysis of robustness to adaptive learning dynamics in section 3.3.2, that ensuring greater robustness to potential (modest) departures from fully model-consistent expectations requires one to adjust the relative weights on inflation and the output gap in a monetary policy rule, in the direction of stronger relative responses to fluctuations in the rate of inflation.

5 Conclusion

This review has been able to illustrate only a few of the possible methods of macroeconomic analysis that depart in one way or another from the complete requirements of the rational expectations hypothesis. Rather than presenting all of the possible specifications of expectations or reviewing all of the conclusions obtained using them in particular models, I have sought only to compare broad classes of approaches, that differ in the respects in which they maintain or depart from particular aspects of the knowledge assumptions maintained in the RE literature. Even this brief overview has shown that there is a considerable range of alternative approaches, leading to different conclusions about a variety of issues.

It may be asked how macroeconomic analysis can be possible with such a wide range of candidate assumptions. One answer would be that empirical studies should

⁶⁸Adam and Woodford also show how to characterize welfare-maximizing policy, when welfare is defined by the expected utility (under the policymaker's expectations) of a representative household. There is again a robustly optimal target criterion of the form (4.17), the coefficients of which now depend purely on model parameters.

be undertaken to determine which of these possible specifications of subjective expectations best describe observed behavior. A few studies of that kind already exist, but the empirical literature remains at a fairly early stage. Much of the early work on the alternatives surveyed here has been undertaken in order to clarify or criticize the conceptual foundations of rational expectations equilibrium, rather than to provide a positive analysis of observed phenomena; further empirical applications are much to be desired.

Nonetheless, it is probably a mistake to suppose that empirical investigations should identify a single model of expectations that can be judged to have been historically valid, and that can then be treated as the way in which expectations *must* be formed in the future, for purposes of counterfactual policy analyses. It is more reasonable, in my view, to search for policies that should be *robust* to a variety of possible specifications of expectations. Of course, it is not possible (and probably would not be desirable, even if feasible) to demand that a policy be robust to *all* possible views of the world; it is therefore important that macroeconomists continue to seek greater certainty about which models of the economy are more accurate. But one need not settle upon a single model specification before policy analysis is possible.

Indeed, the approaches discussed in sections 2 and 4 above seek to define *classes* of reasonable specifications of expectations under a given policy regime, rather than a single correct specification; and even in the case of the econometric learning models discussed in section 3, the identification of a best-fitting learning rule for some historical data set would better be taken as providing evidence about the *types* of learning rules that should be allowed for in a robustness analysis, rather than as identifying a “true” learning rule that can be relied upon in the future. If macroeconomic analysis is approached in this spirit, then awareness of a variety of arguably reasonable specifications should contribute to the robustness of the conclusions reached, rather than preventing any policy recommendations from being given.

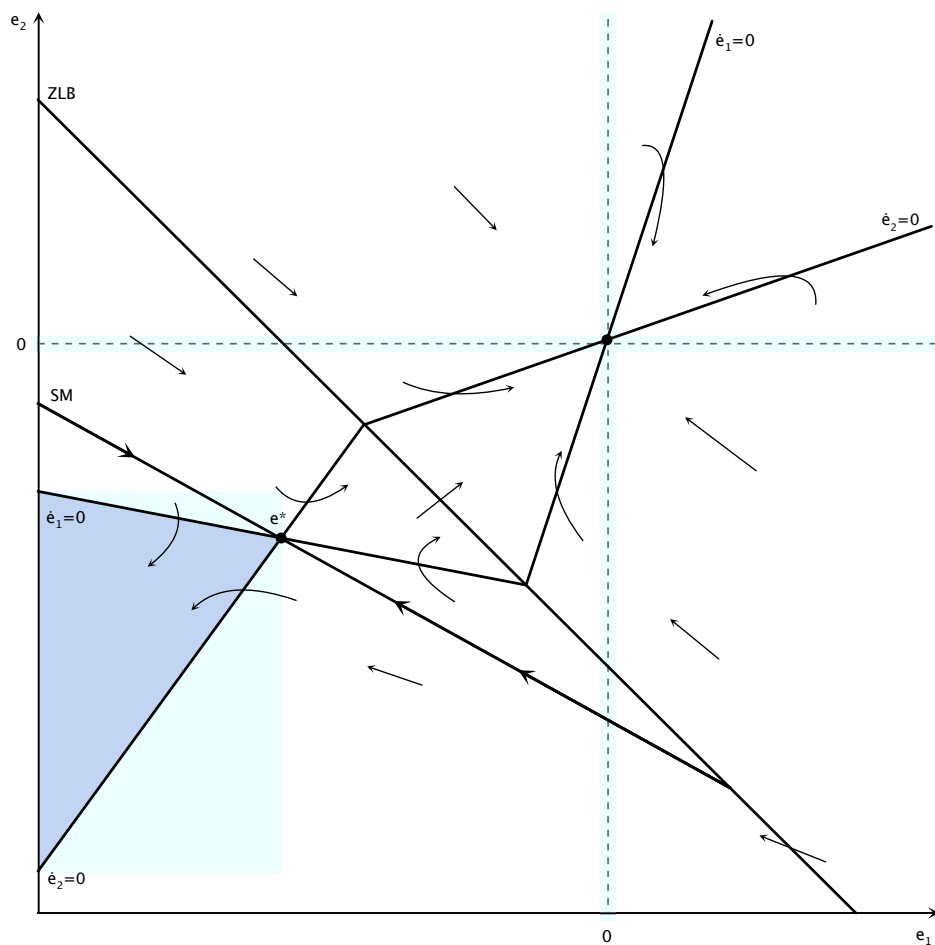


Figure 1: ODE trajectories that approximate asymptotic learning dynamics, when interest-rate policy is constrained by the zero lower bound. The points at which the constraint binds are those below and to the left of the line ZLB. The steady state in which the inflation target is achieved (corresponding to the origin) is a locally stable rest point of the ODE dynamics, but there is also a second steady state (point e^*) at which the ZLB constraint binds, with stable manifold SM.

References

- [1] Adam, Klaus, and Michael Woodford, “Robustly Optimal Monetary Policy in a Microfounded New Keynesian Model,” *Journal of Monetary Economics* 59: 468-487 (2012).
- [2] Bassetto, Marco, “A Game-Theoretic View of the Fiscal Theory of the Price Level,” *Econometrica* 70: 2167-2195 (2002).
- [3] Benhabib, Jess, George W. Evans, and Seppo Honkapohja, “Liquidity Traps and Expectation Dynamics: Fiscal Stimulus or Fiscal Austerity?” NBER Working Paper no. 18114, May 2012.
- [4] Benhabib, Jess, Stephanie Schmitt-Grohé, and Martin Uribe, “The Perils of Taylor Rules,” *Journal of Economic Theory* 96: 40-69 (2001).
- [5] Bernheim, Douglas, “Rationalizable Strategic Behavior,” *Econometrica* 52: 1007-1028 (1984).
- [6] Blanchard, Olivier J., and Charles Kahn, “The Solution of Linear Difference Equations under Rational Expectations,” *Econometrica* 48: 1305-1311 (1980).
- [7] Branch, William A., “Restricted Perceptions Equilibria and Learning in Macroeconomics,” in D. Colander, ed., *Post Walrasian Macroeconomics: Beyond the Dynamic Stochastic General Equilibrium Model*, Cambridge: Cambridge University Press, 2004.
- [8] Bray, Margaret, “Learning, Estimation and Stability of Rational Expectations,” *Journal of Economic Theory* 26: 318-339 (1982).
- [9] Bullard, James, and Kaushik Mitra, “Learning About Monetary Policy Rules,” *Journal of Monetary Economics* 49: 1105-1129 (2002).
- [10] Christiano, Lawrence J., Martin Eichenbaum, and Charles L. Evans, “Nominal Rigidities and the Dynamic Effects of a Shock to Monetary Policy,” *Journal of Political Economy* 113: 1-45 (2005).
- [11] Clarida, Richard, Jordi Gali, and Mark Gertler, “The Science of Monetary Policy,” *Journal of Economic Literature* 37: 1661-1707 (1999).

- [12] Clarida, Richard, Jordi Gali, and Mark Gertler, “Monetary Policy Rules and Macroeconomic Stability: Evidence and Some Theory,” *Quarterly Journal of Economics* 115: 147-180 (2000).
- [13] Cover, Thomas M., and Joy A. Thomas, *Elements of Information Theory*, New York: Wiley-Interscience, 2d ed., 2006.
- [14] Deaton, Angus, *Understanding Consumption*, Oxford: Oxford University Press, 1992.
- [15] De Grauwe, Paul, “Top-Down Versus Bottom-Up Macroeconomics,” *CESifo Economic Studies* 56: 465-497 (2010).
- [16] Eusepi, Stefano, and Bruce Preston, “Expectations, Learning and Business Cycle Fluctuations,” *American Economic Review* 101: 2844-2872 (2011).
- [17] Eusepi, Stefano, and Bruce Preston, “Debt, Policy Uncertainty, and Expectations Stabilization,” *Journal of the European Economic Association*, forthcoming 2012a.
- [18] Eusepi, Stefano, and Bruce Preston, “Fiscal Foundations of Inflation: Imperfect Knowledge,” unpublished, Federal Reserve Bank of New York, September 2012b.
- [19] Evans, George W., and Seppo Honkapohja, *Learning and Expectations in Macroeconomics*, Princeton: Princeton University Press, 2001.
- [20] Evans, George W., and Seppo Honkapohja, “Expectations and the Stability Problem for Optimal Monetary Policies,” *Review of Economic Studies* 70: 807-824 (2003).
- [21] Evans, George W., and Seppo Honkapohja, “Monetary Policy, Expectations and Commitment,” *Scandinavian Journal of Economics* 108: 15-38 (2006).
- [22] Evans, George W., and Seppo Honkapohja, “Learning and Macroeconomics,” *Annual Review of Economics* 1: 421-451 (2009).
- [23] Evans, George W., and Seppo Honkapohja, “Expectations, Deflation Traps, and Macroeconomic Policy,” in D. Cobham, *et al.*, eds., *Twenty Years of Inflation Targeting: Lessons Learned and Future Prospects*, Cambridge: Cambridge University Press, 2010.

- [24] Evans, George W., Seppo Honkapohja, and Kaushik Mitra, “Does Ricardian Equivalence Hold when Expectations are Not Rational?” *Journal of Money, Credit and Banking*, forthcoming 2012.
- [25] Evans, George W., and Bruce McGough, “Learning to Optimize,” unpublished, University of Oregon, March 2009.
- [26] Evans, George W., and Garey Ramey, “Expectation Calculation and Macroeconomic Dynamics,” *American Economic Review* 82: 207-224 (1992).
- [27] Friedman, Milton, “The Role of Monetary Policy,” *American Economic Review* 58: 1-17 (1968).
- [28] Fuster, Andreas, David Laibson, and Brock Mendel, “Natural Expectations and Macroeconomic Fluctuations,” *Journal of Economic Perspectives* 24: 67-84 (2010).
- [29] Fuster, Andreas, Benjamin Hebert, and David Laibson, “Natural Expectations, Macroeconomic Dynamics, and Asset Pricing,” *NBER Macroeconomics Annual* 26: 1-48 (2011).
- [30] Fuster, Andreas, Benjamin Hebert, and David Laibson, “Investment Dynamics with Natural Expectations,” *International Journal of Central Banking* 8: 243-265 (2012).
- [31] Galí, Jordi, *Monetary Policy, Inflation, and the Business Cycle: An Introduction to the New Keynesian Framework*, Princeton: Princeton University Press, 2008.
- [32] Gaspar, Vitor, Frank Smets, and David Vestin, “Inflation Expectations, Adaptive Learning and Optimal Monetary Policy,” in B.M. Friedman and M. Woodford, eds., *Handbook of Monetary Economics*, volume 3B, Amsterdam: Elsevier, 2011.
- [33] Grandmont, Jean-Michel, “Temporary General Equilibrium Theory,” *Econometrica* 45: 535-572 (1977).
- [34] Grandmont, Jean-Michel, ed., *Temporary Equilibrium: Selected Readings*, Boston: Academic Press, 1988.

- [35] Guesnerie, Roger, “An Exploration of the Eductive Justifications of the Rational Expectations Hypothesis,” *American Economic Review* 82: 1254-1278 (1992).
- [36] Guesnerie, Roger, *Assessing Rational Expectations 2: Eductive Stability in Economics*, Cambridge, MA: MIT Press, 2005.
- [37] Guesnerie, Roger, “Macroeconomic and Monetary Policies from the Eductive Viewpoint,” in K. Schmidt-Hebbel and C. Walsh, eds., *Monetary Policy Under Uncertainty and Learning*, Santiago: Central Bank of Chile, 2008.
- [38] Hansen, Lars Peter, and Thomas J. Sargent, *Robustness*, Princeton: Princeton University Press, 2008.
- [39] Hansen, Lars Peter, and Thomas J. Sargent, “Three Types of Ambiguity,” unpublished, University of Chicago, July 2012.
- [40] Harvey, Andrew C., *Forecasting, Structural Time Series Models and the Kalman Filter*, Cambridge: Cambridge University Press, 1989.
- [41] Hicks, John, *Value and Capital*, Oxford: Clarendon Press, 1939.
- [42] Howitt, Peter, “Interest-Rate Control and Nonconvergence to Rational Expectations,” *Journal of Political Economy* 100: 776-800 (1992).
- [43] Kurz, Mordecai, “On the Structure and Diversity of Rational Beliefs,” *Economic Theory* 4: 877-900 (1994).
- [44] Kurz, Mordecai, ed., *Endogenous Economic Fluctuations: Studies in the Theory of Rational Belief*, Berlin: Springer, 1997.
- [45] Kurz, Mordecai, “A New Keynesian Model with Diverse Beliefs,” unpublished, Stanford University, February 2012.
- [46] Kurz, Mordecai, and Maurizio Motolese, “Diverse Beliefs and Time Variability of Risk Premia,” *Economic Theory* 47: 293-335 (2011).
- [47] Kwon, Hyosung, and Jianjun Miao, “Woodford’s Approach to Robust Policy Analysis in a Linear-Quadratic Framework,” unpublished, Boston University, December 2012.

- [48] Lindahl, Erik, “Theory of Money and Capital,” London: Allen and Unwin, 1939.
- [49] Lubik, Thomas A., and Frank Schorfheide, “Testing for Indeterminacy: An Application to U.S. Monetary Policy,” *American Economic Review* 94: 190-217 (2004).
- [50] Marcet, Albert, and Thomas J. Sargent, “Convergence of Least Squares Learning Mechanisms in Self-Referential Linear Stochastic Models,” *Journal of Economic Theory* 48: 337-368 (1989).
- [51] Maskin, Eric, and Jean Tirole, “Markov Perfect Equilibrium: I. Observable Actions,” *Journal of Economic Theory* 100: 191-219 (2001).
- [52] McCallum, Bennett T., “On Nonuniqueness in Linear Rational Expectations Models: An Attempt at Perspective,” *Journal of Monetary Economics* 11: 134-168 (1983).
- [53] Milani, Fabio, “Learning, Monetary Policy Rules, and Macroeconomic Stability,” *Journal of Economic Dynamics and Control* 32: 3148-3165 (2005).
- [54] Milani, Fabio, “Expectations, Learning and Macroeconomic Persistence,” *Journal of Monetary Economics* 54: 2065-2082 (2007).
- [55] Milani, Fabio, “Expectation Shocks and Learning as Drivers of the Business Cycle,” *Economic Journal* 121: 379-401 (2011).
- [56] Muth, John F., “Optimal Properties of Exponentially Weighted Forecasts,” *Journal of the American Statistical Association* 55: 299-305 (1960).
- [57] Orphanides, Athanasios, and John C. Williams, “Imperfect Knowledge, Inflation Expectations and Monetary Policy,” in B.S. Bernanke and M. Woodford, eds., *The Inflation Targeting Debate*, Chicago: University of Chicago Press, 2005.
- [58] Pearce, David, “Rationalizable Strategic Behavior and the Problem of Perfection,” *Econometrica* 52: 1029-1050 (1984).
- [59] Phelps, Edmund S., “The Trouble with ‘Rational Expectations’ and the Problem of Inflation Stabilization,” in R. Frydman and E.S. Phelps, eds., *Individual Forecasting and Aggregate Outcomes: “Rational Expectations” Reconsidered*, Cambridge: Cambridge University Press, 1983.

- [60] Preston, Bruce, “Learning about Monetary Policy Rules when Long-Horizon Expectations Matter,” *International Journal of Central Banking* 1: 81-126 (2005).
- [61] Rondina, Giacomo, and Todd B. Walker, “Information Equilibria in Dynamic Economies with Dispersed Information,” with Todd B. Walker, unpublished, U.C. San Diego, March 2012.
- [62] Sargent, Thomas J., *Bounded Rationality in Macroeconomics*, Oxford: Oxford University Press, 1993.
- [63] Sargent, Thomas J., “Evolution and Intelligent Design,” *American Economic Review* 98: 5-37 (2008).
- [64] Slobodyan, Sergey, and Raf Wouters, “Learning in an Estimated Medium-Scale DSGE Model,” CERGE-EI working paper no. 396, November 2009.
- [65] Smets, Frank, and Raf Wouters, “Shocks and Frictions in US Business Cycles: A Bayesian DSGE Approach,” *American Economic Review* 97: 586-606 (2007).
- [66] Svensson, Lars E.O., “Inflation Targeting as a Monetary Policy Rule,” *Journal of Monetary Economics* 43: 607-654 (1999).
- [67] Taylor, John B., “Discretion Versus Policy Rules in Practice,” *Carnegie-Rochester Conference Series on Public Policy* 39: 195-214 (1993).
- [68] Walsh, Carl E., *Monetary Theory and Policy*, Cambridge, MA: MIT Press, 3d ed., 2010.
- [69] Woodford, Michael, “Fiscal Requirements for Price Stability,” *Journal of Money, Credit and Banking* 33: 669-728 (2001).
- [70] Woodford, Michael, *Interest and Prices: Foundations of a Theory of Monetary Policy*, Princeton: Princeton University Press, 2003.
- [71] Woodford, Michael, “Robustly Optimal Monetary Policy with Near-Rational Expectations,” *American Economic Review* 100: 274-303 (2010).

A Appendix: Details of Calculations

This appendix provides additional details of several derivations referred to in the main text.

A.1 Derivation of the Temporary Equilibrium Conditions

The equilibrium conditions given in the text represent linearized versions of the conditions required for a temporary equilibrium in a model described more fully here. The economy is made up of a continuum of identical infinite-lived households, indexed by $i \in [0, 1]$. In the plan that each household i formulates in period t , it seeks to maximize its estimate (based on subjective probabilities) of the discounted sum of utilities in the remaining periods of its life,

$$\hat{E}_t^i \sum_{T=t}^{\infty} \beta^{T-t} [u(C_T^i; \xi_T) - v(H_T^i; \xi_T)]. \quad (\text{A.1})$$

Here C_t^i is a Dixit-Stiglitz (or CES) aggregate of the household's purchases of differentiated consumer goods,

$$C_t^i \equiv \left[\int_{j=0}^1 c_t^i(j)^{\frac{\theta-1}{\theta}} dj \right]^{\frac{\theta}{\theta-1}}, \quad (\text{A.2})$$

where $c_t^i(j)$ is the quantity purchased of good j and $\theta > 1$ is the elasticity of substitution among different goods; H_t^i is hours worked by the household in period t ; and ξ_t is a vector of exogenous disturbances that includes possible disturbances to both the urgency of immediate consumption and the disutility of labor (that need not be correlated, since ξ_t is a vector).

Given the assumption of a single financial asset, a one-period riskless nominal bond, the household's bond holdings evolve according to

$$B_{t+1}^i = (1 + i_t) \left[B_t^i + W_t H_t + \int_{j=0}^1 \Pi_t(j) dj - \int_{j=0}^1 p_t(j) c_t^i(j) dj - T_t \right], \quad (\text{A.3})$$

where B_t^i is the nominal value at maturity of the bonds carried into period t , W_t is the nominal wage, Π_t^j is nominal profits of firm j (distributed in equal shares to the households who own the firms), $p_t(j)$ is the price of good j , and T_t is the net nominal (lump-sum) tax obligation (assumed to be equal for all households). (Note that here I write H_t for hours worked by the household, because of the assumption that available hours of work are allocated equally to each household.) Each firm's profits are given by

$$\Pi_t(j) = p_t(j) y_t(j) - W_t H_t(j),$$

where $y_t(j)$ is the quantity produced and sold of good j and $H_t(j)$ is labor hired by firm j . Integrating this over firms and noting that

$$\int_{j=0}^1 H_t(j) dj \equiv H_t,$$

we see that each household's income other than from its bond holdings must equal

$$W_t H_t + \int_{j=0}^1 \Pi_t(j) dj = \int_{j=0}^1 p_t(j) y_t(j) dj.$$

Finally, the form of the Dixit-Stiglitz aggregator (A.2) implies that in the case of an optimal allocation of household expenditure across differentiated goods, total expenditure will equal

$$\int_{j=0}^1 p_t(j) c_t^i(j) dj = P_t C_t^i,$$

where

$$P_t \equiv \left[\int_{j=0}^1 p_t(j)^{1-\theta} dj \right]^{\frac{1}{1-\theta}}$$

is the Dixit-Stiglitz price index. As this holds for all purchasers of goods (including the government, it is assumed), total sales revenues will similarly equal

$$\int_{j=0}^1 p_t(j) y_t(j) dj = P_t Y_t,$$

where Y_t is total demand for the composite good defined in (A.2). Substituting these expressions into (A.3), the law of motion for bond holdings can be written more simply as

$$B_{t+1}^i = (1 + i_t) [B_t^i + P_t Y_t - P_t C_t^i - T_t]. \quad (\text{A.4})$$

The household's consumption plan can then be formulated purely as a choice of a planned state-contingent evolution for $\{C_T^i\}_{T=t}^\infty$, and the household's perceived intertemporal budget constraint (i.e., the set of plans that are believed to be feasible) depends only on the household's initial wealth B_t^i and the expected evolution of the variables $\{Y_T, P_T, i_T, T_T\}_{T=t}^\infty$.

The household's planning problem at date t is then the choice of state-contingent paths $\{C_T^i, B_{T+1}^i\}_{T=t}^\infty$ consistent with (A.4) at all dates $T \geq t$ (together with a bound on how negative bond holdings can be asymptotically, to rule out Ponzi schemes) so as to maximize (A.1), given the household's initial wealth B_t^i and its expectations regarding the evolution of the variables $\{Y_T, P_T, i_T, T_T\}_{T=t}^\infty$ outside its control. (Note that the evolution of $\{H_T\}_{T=t}^\infty$ is also outside the household's control, under the labor market institutions assumed in the model; but the household's consumption-planning

problem is independent of its expectations about how much it will be working.) While we are at present agnostic about the nature of households' expectations about the future evolution of the variables outside their individual control (and do not assume, in general, that households are necessarily aware of any of the structural relations that determine such variables), we assume that households correctly understand the constraints (A.4) that define their own problem, and accordingly that each chooses a plan that solves the problem just stated, under some internally consistent subjective expectations about the evolution of the variables outside its control.

A.1.1 Subjectively Optimal Expenditure

The first-order conditions for the problem just defined are

$$u_C(C_T^i; \xi_T) = \beta(1 + i_T) \hat{E}_T^i [\Pi_{T+1}^{-1} u_C(C_{T+1}^i; \xi_{T+1})] \quad (\text{A.5})$$

for any state of the world that might be reached at any date $T \geq t$, where $\Pi_{T+1} \equiv P_{T+1}/P_T$ is the gross rate of inflation. The household's subjectively optimal plan is then a pair of processes $\{C_T^i, B_{T+1}^i\}_{T=t}^\infty$ satisfying (A.4) and (A.5) at all dates $T \geq t$, together with bounds on the asymptotic growth of net financial wealth (a transversality condition) that guarantee both consistency with the borrowing limit and no inefficient overaccumulation of wealth. We can approximately characterize the optimal plan, looking forward from any date t , by linearizing the conditions that implicitly define the optimal plan around the steady-state values of the endogenous variables that represent a solution in the case of no random fluctuations in the exogenous states.

Assuming constant values $\xi_T = \bar{\xi}$ for all of the exogenous disturbances and that expectations about the economy's deterministic evolution are correct (i.e., perfect foresight), the structural relations (A.4) and (A.5) are consistent with a stationary solution in which $\Pi_T = 1$ (a zero steady-state inflation rate), $i_T = \beta^{-1} - 1 > 0$, $b_T^i \equiv B_T^i/P_{T-1} = \bar{b}$, $Y_T = \bar{Y}$, $C_T^i = \bar{C}$, and $\tau_T \equiv T_T/P_T = \bar{\tau}$ with certainty for all $T \geq t$, under the assumption that (i) the household's initial financial wealth is given by $b_t^i = \bar{b}$,⁶⁹ (ii) the value of \bar{C} satisfies

$$\bar{C} = (1 - \beta)\bar{b} + \bar{T} - \bar{\tau};$$

and (iii) these steady-state values are also consistent with the remaining model structural relations (discussed further below), and in particular with the specifications of monetary and fiscal policy. We now wish to solve for a solution to the structural relations (A.4) and (A.5) near this steady-state solution, in the case of values for the exogenous disturbances that remain near enough to the values $\bar{\xi}$ for all $T \geq t$; a level of initial financial wealth b_t^i near enough to \bar{b} ; and subjective expectations about

⁶⁹Note that b_t^i is defined as B_t^i/P_{t-1} rather than B_t^i/P_t so that it is a predetermined state variable.

the future evolution of all variables that are *close enough* to being correct (i.e., near enough to model-consistency).⁷⁰

We obtain a local linear approximation to the perturbed solution by solving local linear approximations to the structural equations (A.4) and (A.5). In terms of the perturbations

$$\begin{aligned}\hat{b}_t^i &\equiv \frac{b_t^i - \bar{b}}{\bar{Y}}, & \hat{i}_t &\equiv \log(1 + i_t) - \log(\beta^{-1}), & \pi_t &\equiv \log \Pi_t, \\ \hat{Y}_t &\equiv \frac{Y_t - \bar{Y}}{\bar{Y}}, & \hat{c}_t^i &\equiv \frac{C_t^i - \bar{C}}{\bar{Y}}, & \hat{\tau}_t &\equiv \frac{\tau_t - \bar{\tau}}{\bar{Y}},\end{aligned}$$

a linear approximation to (A.4) can be written in the form

$$\hat{b}_{t+1}^i = s_b(\hat{i}_t - \beta^{-1}\pi_t) + \beta^{-1}(\hat{b}_t^i + \hat{Y}_t - \hat{c}_t^i - \hat{\tau}_t), \quad (\text{A.6})$$

where $s_b \equiv \bar{b}/\bar{Y}$. A local linear approximation to the marginal utility of expenditure can similarly be written in the form

$$\log u_C(C_t^i; \xi_t) = \log u_C(\bar{C}; \bar{\xi}) - \sigma^{-1}(\hat{c}_t^i - \bar{c}_t), \quad (\text{A.7})$$

where $\sigma > 0$ is a parameter proportional to the intertemporal elasticity of substitution of expenditure, and \bar{c}_t is an exogenous disturbance (a shock to the urgency of private expenditure), indicating the shift in the size of \hat{c}_t^i required in order to maintain a constant marginal utility of expenditure. A local linear approximation to (A.5) then can be written in the form

$$\hat{c}_t^i - \bar{c}_t = -\sigma(\hat{i}_t - \hat{E}_t^i \pi_{t+1}) + \hat{E}_t^i [\hat{c}_{t+1}^i - \bar{c}_{t+1}]. \quad (\text{A.8})$$

One can show that there is a unique solution to equations (A.6) and (A.8) for periods $T \geq t$ consistent with the transversality condition

$$\lim_{T \rightarrow \infty} \beta^{T-t} \hat{E}_t^i \hat{b}_T^i = 0.$$

Rearranging the terms in (A.6), we can alternatively write

$$\hat{b}_t^i = -(\hat{Y}_t - \hat{\tau}_t) - s_b(\beta \hat{i}_t - \pi_t) + \bar{c}_t + (\hat{c}_t^i - \bar{c}_t) + \beta \hat{b}_{t+1}^i.$$

This can then be “solved forward” to obtain

$$\hat{b}_t^i = - \sum_{T=t}^{\infty} \hat{E}_t^i \left\{ (\hat{Y}_T - \hat{\tau}_T) + s_b(\beta \hat{i}_T - \pi_T) + \bar{c}_T + (\hat{c}_T^i - \bar{c}_T) \right\} \quad (\text{A.9})$$

⁷⁰Because of the local approximation that is relied upon after this point, there is a sense in which all of the analysis in the paper assumes “near-rational expectations.” Nonetheless, while the linearized temporary-equilibrium relations are relied upon throughout the text, bounds on the possible expectational errors are not assumed except when additional restrictions on expectations are explicitly introduced, as for example in section 4.2 of the text.

But (A.8) implies that for any $T > t$,

$$\hat{E}_t^i[\hat{c}_T^i - \bar{c}_T] = (\hat{c}_t^i - \bar{c}_t) + \sigma \sum_{s=t}^{T-1} \hat{E}_t^i[\hat{i}_s - \pi_{s+1}],$$

so that

$$\sum_{T=t}^{\infty} \beta^{T-t} \hat{E}_t^i[\hat{c}_T^i - \bar{c}_T] = (1 - \beta)^{-1} (\hat{c}_t^i - \bar{c}_t) + (1 - \beta)^{-1} \sigma \sum_{T=t}^{\infty} \beta^{T+1-t} \hat{E}_t^i[\hat{i}_T - \pi_{T+1}].$$

Substituting this last result into (A.9), we obtain

$$\begin{aligned} \hat{b}_t^i &= - \sum_{T=t}^{\infty} \hat{E}_t^i \left\{ (\hat{Y}_T - \hat{\tau}_T) + s_b(\beta \hat{i}_T - \pi_T) + \bar{c}_T + (1 - \beta)^{-1} \beta \sigma (\hat{i}_T - \pi_{T+1}) \right\} + (1 - \beta)^{-1} (\hat{c}_t^i - \bar{c}_t) \\ &= - \sum_{T=t}^{\infty} \hat{E}_t^i \left\{ (\hat{Y}_T - \hat{\tau}_T) + s_b(\beta \hat{i}_T - \pi_T) + (1 - \beta)^{-1} \beta \sigma (\hat{i}_T - \pi_{T+1}) + (1 - \beta)^{-1} \beta (\bar{c}_{T+1} - \bar{c}_T) \right\} \\ &\quad + (1 - \beta)^{-1} \hat{c}_t^i. \end{aligned}$$

This equation can be solved for the value of \hat{c}_t^i under the subjectively optimal plan. This yields equation (1.1) in the text. (Note, however, that hats are omitted from all variables in the text. That is, b_t^i in the text refers to the perturbation variable here denoted by \hat{b}_t^i , and so on.)

A.1.2 Labor Supply and Wage Determination

If we instead write the law of motion (A.3) for a household's financial wealth in real terms, we observe that the household's intertemporal budget set looking forward from any date t depends only on the value of the household's real period t wealth

$$a_t^i \equiv \frac{B_t^i + W_t H_t^i}{P_t}$$

(prior to profit distributions, taxes and transfer payments) and the expected paths of the variables $\{Y_T, i_T, \int_j \Pi_T(j) dj, \tau_T\}_{T=t}^{\infty}$ and $\{\Pi_T\}_{T=t+1}^{\infty}$ outside the household's control. Let $V_t^i(a_t^i)$ denote the household's subjective evaluation in period t of the maximum attainable value of the continuation utility (A.1) given the value of a_t^i , where the time subscript indicates the dependence of this function on expectations at t regarding the variables outside the household's control. Using the envelope theorem, we observe that the derivative of this value function will equal

$$V_t^{i'}(a_t^i) = u_C(C_t^i; \xi_t) \tag{A.10}$$

where C_t^i is the household's optimal period t expenditure, characterized above.

Under the assumed labor market institution, then, the union supplies H_t hours of work from each household in period t so as to maximize the value of

$$\int_0^1 V_t^i(b_t^i/\Pi_t + w_t H_t) di - v(H_t; \xi_t),$$

in the case of any real wage $w_t \equiv W_t/P_t$. The first-order condition for optimal aggregate labor supply is then

$$\int_0^1 V_t^{i'}(a_t^i) di \cdot w_t = v_H(H_t; \xi_t).$$

Using (A.10), this can alternatively be written in the form

$$\int_0^1 u_C(C_t^i; \xi_t) di \cdot w_t = v_H(H_t; \xi_t). \quad (\text{A.11})$$

We can then log-linearize relation (A.11) around the same steady-state values as above. Using (A.7), the local linear approximation takes the form

$$\hat{w}_t - \sigma^{-1} \int_0^1 (\hat{c}_t^i - \bar{c}_t) di = \nu_t, \quad (\text{A.12})$$

introducing the notation

$$\hat{w}_t \equiv \log(w_t/\bar{w}), \quad \nu_t \equiv \log v_H(H_t; \xi_t) - \log v_H(\bar{H}; \bar{\xi})$$

where \bar{w} is the steady-state real wage. Finally, letting $C_t \equiv \int_0^1 C_t^i di$ denote aggregate household expenditure and introducing the notation

$$\hat{c}_t \equiv \log(C_t/\bar{C}),$$

we observe that to a linear approximation

$$\hat{c}_t = \int_0^1 \hat{c}_t(i) di.$$

Substituting this into (A.12), we obtain

$$\hat{w}_t - \sigma^{-1}(\hat{c}_t - \bar{c}_t) = \nu_t.$$

This is the log-linear wage equation given in section 1.3 of the text, except that once again hats are omitted in the notation used in the text.

A.2 Non-uniqueness of Rationalizable Equilibrium

Here we present additional details of the calculations involved in establishing the non-uniqueness of rationalizable equilibrium (or, failure of “eductive stability”) in the example of Guesnerie (2008), as stated in section 2.2 of the text. Under the assumption of a monetary policy rule of the form $i_t = \phi_\pi \pi_t$, the TE dynamics must satisfy

$$\bar{v}_t^i = (1 - \beta\phi_\pi)\bar{v}_t + \beta\hat{E}_t^i \bar{v}_{t+1}^i \quad (\text{A.13})$$

for all i , as explained in the text.⁷¹ In a rationalizable TE, not only must this hold at date t , but everyone must expect anyone else to expect anyone else ... to expect it to hold at any future date. We wish to consider whether the RE equilibrium in which $\bar{v}_t^i = 0$ for all i and all t constitutes the unique bounded process $\{\bar{v}_t^i\}$ consistent with common knowledge of (A.13). In the case that $|\phi_\pi| < 1$, the solution is obvious not unique, because in this case, there is not even a unique RE equilibrium, as is well known; and all REE are also rationalizable TE. Here we show that even when $\phi_\pi > 1$, so that $\bar{v}_t^i = 0$ is the unique bounded REE, it is possible to have a large multiplicity of bounded rationalizable TE.

Common knowledge that (A.13) holds for all i and all t implies that the hierarchy of beliefs at any date t must satisfy

$$\begin{aligned} \hat{E}_t^{i_1} \hat{E}_{t_1}^{i_2} \cdots \hat{E}_{t+j_{n-1}}^{i_n} \bar{v}_{t+j_n}^{i_n} &= (1 - \beta\phi_\pi) \hat{E}_t^{i_1} \hat{E}_{t_1}^{i_2} \cdots \hat{E}_{t+j_{n-1}}^{i_n} \bar{v}_{t+j_n} \\ &+ \beta \hat{E}_t^{i_1} \hat{E}_{t_1}^{i_2} \cdots \hat{E}_{t+j_{n-1}}^{i_n} \bar{v}_{t+j_{n+1}}^{i_n}, \end{aligned} \quad (\text{A.14})$$

$$\begin{aligned} \hat{E}_t^{i_1} \hat{E}_{t_1}^{i_2} \cdots \hat{E}_{t+j_{n-1}}^{i_n} \bar{v}_{t+j_n}^{i_{n+1}} &= (1 - \beta\phi_\pi) \hat{E}_t^{i_1} \hat{E}_{t_1}^{i_2} \cdots \hat{E}_{t+j_{n-1}}^{i_n} \bar{v}_{t+j_n} \\ &+ \beta \hat{E}_t^{i_1} \hat{E}_{t_1}^{i_2} \cdots \hat{E}_{t+j_{n-1}}^{i_n} \hat{E}_{t+j_n}^{i_{n+1}} \bar{v}_{t+j_{n+1}}^{i_{n+1}}, \end{aligned} \quad (\text{A.15})$$

where $t < t_1 < \cdots < t_{n+1}$ is any sequence of dates beginning with t , and $i_1 \neq i_2 \neq \cdots \neq i_{n+1}$ is any sequence of households. At the same time, any hierarchy of beliefs satisfying (A.14)–(A.15) will describe a rationalizable TE. The hierarchy of beliefs regarding the future paths of inflation and the interest rate is derivable from the hierarchy of beliefs about the $\{\bar{v}_t^i\}$, using the assumption that both equation (1.8) in the text and the policy rule are common knowledge.

One possible solution to (A.13), (A.14) and (A.15) is given by $\bar{v}_t^i = \epsilon$, and

$$\hat{E}_t^{i_1} \hat{E}_{t_1}^{i_2} \cdots \hat{E}_{t+j_{n-1}}^{i_n} \bar{v}_{t+j_n}^{i_n} = (-\mu)^{1-n} \phi_\pi \epsilon,$$

$$\hat{E}_t^{i_1} \hat{E}_{t_1}^{i_2} \cdots \hat{E}_{t+j_{n-1}}^{i_n} \bar{v}_{t+j_n}^{i_{n+1}} = (-\mu)^{-n} \epsilon$$

⁷¹See equation (2.2) in the text.

for any sequences of households and dates of the kind assumed above, where ϵ is an arbitrary real number and

$$\mu \equiv \frac{\beta\phi_\pi - 1}{(1 - \beta)\phi_\pi}.$$

These beliefs satisfy all of the requirements for rationalizability for any real number ϵ , as can be verified by substituting the candidate solution into equations (A.13), (A.14) and (A.15) and verifying that each condition is satisfied.

Moreover, if $1/2 < \beta < 1$ and $\phi_\pi > (2\beta - 1)^{-1} > 1$ is satisfied, as stated in the text, then $\mu > 1$, and forecasts of all orders satisfy a uniform bound. There is then (at least) a continuum of uniformly bounded rationalizable TE. Moreover, ϵ may represent the realization of a “sunspot” event unrelated to fundamentals, so that there are seen to exist bounded sunspot equilibria, despite the fact that monetary policy satisfies the Taylor Principle.

A.3 Restricted Perception Equilibrium

Here we explain further details of the example of a restricted perception equilibrium in which Ricardian Equivalence fails. The policy regime assumed in the example of section 3.1.2 implies that the dynamics of b_t and s_t are given by

$$b_{t+1} = \beta^{-1}(b_t - s_t) \tag{A.16}$$

$$s_t = \phi_b b_t + \epsilon_t^s \tag{A.17}$$

where $\{\epsilon_t^s\}$ is an i.i.d. random variable with mean zero and variance σ^2 . Substitution of (A.17) into (A.16) yields the univariate law of motion

$$b_{t+1} = \rho b_t - \beta^{-1}\epsilon_t^s,$$

where $\rho \equiv \beta^{-1}(1 - \phi_b)$. Under the assumption that

$$0 < 1 - \beta < \phi_b < 1, \tag{A.18}$$

$0 < \rho < 1$ and this law of motion implies that $\{b_t\}$ is stationary AR(1) process with positive serial correlation.

It follows from this law of motion that the unconditional variance of the stationary process is given by

$$\mathbb{E}[b^2] = \frac{\beta^{-1}\sigma^2}{1 - \rho^2}. \tag{A.19}$$

It then follows from (A.17) that

$$\mathbb{E}[sb] = \phi_b \mathbb{E}[b^2], \quad \mathbb{E}[s^2] = \phi_b^2 \mathbb{E}[b^2] + \sigma^2, \tag{A.20}$$

and from (A.16) that

$$\mathbb{E}[b_{t+1}s_t] = (\rho - \beta) \mathbb{E}[b^2]. \quad (\text{A.21})$$

The linear projection is then defined as $P_t[b_{t+1}] = \Lambda_b s_t$, where

$$\Lambda_b \equiv \frac{\mathbb{E}[b_{t+1}s_t]}{\mathbb{E}[s_t^2]} \quad (\text{A.22})$$

is the OLS regression coefficient. It follows from (A.20), (A.21) and this definition that

$$\Lambda_b = \frac{(\rho - \beta)\mathbb{E}[b^2]}{\phi_b^2\mathbb{E}[b^2] + \sigma^2} < \frac{\rho - \beta}{\phi_b^2} < \frac{1 - \beta}{\phi_b^2} < \frac{1}{\phi_b}, \quad (\text{A.23})$$

where the last inequality relies upon (A.18).

The linear equation (3.12) in the text, to solve for ψ_v , can be written in the form

$$A(\psi_v) = 0,$$

where the function $A(\psi_v)$ is defined as the left-hand side of (3.12) minus the right-hand side. This is a linear function of the form

$$A(\psi_v) = a\psi_v + b,$$

where

$$a = 1 - \phi_b\Lambda_b > 0$$

as a consequence of (A.23). There will thus be a unique value of ψ_v for which $A(\psi_v) = 0$. Moreover, we observe that

$$A(\beta^{-1} - 1) = (\beta^{-1} - 1) - (1 - \beta)(\beta^{-1} - 1)\Lambda_b > (\beta^{-1} - 1) \left(1 - \frac{1 - \beta}{\phi_b}\right) > 0,$$

where the first inequality follows from (A.23) and the second from (A.18). Then since $A(\psi_v)$ is an increasing function, the zero of the function must occur for a value of ψ_v less than this. Hence the unique solution satisfies $\psi_v < \beta^{-1} - 1$, as asserted in the text.

A.4 Phase Dynamics Shown in Figure 1

When the zero lower bound (ZLB) constraint does not bind, the non-stochastic part of the solution given in equation (3.19) of the text can be written explicitly as

$$\mathbf{x}_t \equiv \begin{bmatrix} \pi_t \\ y_t \\ i_t \end{bmatrix} = \frac{1}{\Delta} \begin{bmatrix} \kappa & (1 - \alpha)\beta(1 + \sigma\phi_y) \\ 1 & -(1 - \alpha)\beta\sigma\phi_\pi \\ \kappa\phi_\pi + \phi_y & (1 - \alpha)\beta\phi_\pi \end{bmatrix} \begin{bmatrix} e_{1t} \\ e_{2t} \end{bmatrix} \equiv De_t,$$

where

$$\Delta \equiv 1 + \kappa\sigma\phi_\pi + \sigma\phi_y > 0,$$

allowing us to identify the matrix D . Using the elements of the matrix C indicated by the coefficients of the linear equations (3.16)–(3.17) in the text, we then find that

$$I - A \equiv I - CD = \frac{1}{\Delta} \begin{bmatrix} \frac{\sigma}{1-\beta}[\kappa(\phi_\pi - 1) + \phi_y] & \frac{\sigma}{1-\beta}(1-\alpha)\beta[(\phi_\pi - 1) - \sigma\phi_y] \\ -\frac{\kappa}{(1-\alpha)(1-\alpha\beta)} & \frac{\kappa\sigma\phi_\pi + (1-\beta)(1+\sigma\phi_y)}{1-\alpha\beta} \end{bmatrix}.$$

This matrix defines the linear ODE system given by (3.22) in the text, for the part of the plane where the ZLB constraint does not bind.

The two rows of the matrix $M \equiv I - A$ give the coefficients of the linear equations that define the loci $\dot{e}_1 = 0$ and $\dot{e}_2 = 0$ respectively. Let the elements of M be denoted $\{m_{ij}\}$. Then since $m_{21} < 0, m_{22} > 0$, the locus $\dot{e}_2 = 0$ is necessarily an upward-sloping line, as shown in the figure. Furthermore, $\dot{e}_2 < 0$ for all points in the phase plane above this line, while $\dot{e}_2 > 0$ for all points below it.

Because m_{12} cannot be signed in general, the locus $\dot{e}_1 = 0$ may have either a positive or negative slope (or may be perfectly vertical). However, in the case of response coefficients satisfying

$$\phi_\pi + \frac{1-\beta}{\kappa}\phi_y > 1, \tag{A.24}$$

i.e., consistent with the Taylor Principle, $m_{11} > 0$, which implies that $\dot{e}_1 > 0$ for all points to the left of this locus, while $\dot{e}_1 < 0$ for all points to its right. In addition, assumption (A.24) implies that

$$\frac{m_{22}}{m_{21}} - \frac{m_{12}}{m_{11}} = -\frac{(1-\alpha)\Delta}{\kappa(\phi_\pi - 1) + \phi_y} \left[\phi_\pi + \frac{1-\beta}{\kappa}\phi_y - 1 \right] < 0,$$

so that the inverse slope de_1/de_2 must be greater (more positive) for the $\dot{e}_2 = 0$ locus than for the $\dot{e}_1 = 0$ locus. Hence the relative slopes of the two loci in the unconstrained region are necessarily as shown in Figure 1.

It also follows from the above discussion that at all points above the $\dot{e}_2 = 0$ locus and to the right of the $\dot{e}_1 = 0$ locus, trajectories must all move down and to the left; while at all point above the $\dot{e}_2 = 0$ locus and to the left of the $\dot{e}_1 = 0$ locus, they must all move down and to the right; and so on. Hence the phase dynamics in the unconstrained region (the region above the line labeled “ZLB”) must be as shown in the figure. These sign restrictions suffice to imply that in an open set around the zero-inflation steady state (point $e = 0$ in the figure), all trajectories converge asymptotically to that steady state (i.e., the zero-inflation steady state is locally asymptotically stable under the ODE dynamics). If there were no ZLB constraint,

so that equation (3.22) in the text applied globally, this steady state (which would then be the unique rest point of the ODE system) would also be *globally* stable.

The equation of the locus of points at which the ZLB constraint just binds (line *ZLB* in the figure) is of the form

$$a'e = -\bar{r}\Delta < 0,$$

where

$$a' \equiv [\kappa\phi_\pi + \phi_y \quad (1 - \alpha)\beta\phi_\pi].$$

(This follows from equation (3.23) in the text, and our explicit solution above for the matrix D .) Since $a_1, a_2 > 0$, this equation describes a downward-sloping straight line, located below and to the left of the point $e = 0$, as shown in Figure 1. The points above and to the right of this line constitute the values of e for which the ZLB constraint will not bind, while below and to the left of the line, the constraint is strictly binding (and hence the ODE system given by (3.22) in the text does not apply).

Then since the $\dot{e}_2 = 0$ locus is upward-sloping, as established above, it must intersect the line *ZLB* at a point below and to the left of the point $e = 0$, as shown in the figure. The locus $\dot{e}_1 = 0$ is not necessarily upward-sloping (as drawn in the figure), but under assumption (A.24), one can show that

$$\frac{a_2}{a_1} - \frac{m_{12}}{m_{11}} = -\frac{(1 - \alpha)\beta\phi_y\Delta}{[\kappa(\phi_\pi - 1) + \phi_y](\kappa\phi_\pi + \phi_y)} > 0.$$

This implies that the inverse slope of the line *ZLB* is necessarily more negative than that of the locus $\dot{e}_1 = 0$, so that the two lines also must intersect, as shown. It further follows from our conclusion above about the relative slopes of the $\dot{e}_1 = 0$ and $\dot{e}_2 = 0$ loci in the unconstrained region that the point at which the locus $\dot{e}_1 = 0$ intersects the line *ZLB* must lie below and to the right of the point at which the locus $\dot{e}_2 = 0$ intersects this line, as shown in the figure. Hence the qualitative dynamics in the unconstrained region are as shown in the figure.

In the constrained region, the dynamics are instead defined by the alternative linear ODE system specified in equation (3.26) of the text, where the matrix \underline{D} is obtained by substituting $\phi_\pi = \phi_y = 0$ into the expression given above for the matrix D . This system can alternatively be written in the form

$$\dot{e} = -\underline{M}(e - e^*), \tag{A.25}$$

where the matrix \underline{M} is obtained by substituting $\phi_\pi = \phi_y = 0$ into the expression given above for the matrix M .

The elements of the matrix \underline{M} then define the coefficients of the equations corresponding to the loci $\dot{e}_1 = 0$ and $\dot{e}_2 = 0$ in the constrained region. Since $\underline{m}_{11}, \underline{m}_{12} < 0$,

the $\dot{e}_1 = 0$ locus must be upward-sloping; and similarly, since $m_{21} < 0$ and $m_{22} > 0$, the $\dot{e}_2 = 0$ locus must be downward-sloping, as shown in Figure 1. Then given the relative positions of the points at which these loci intersect the line ZLB , discussed above, the two loci necessarily intersect at a point e^* in the interior of the constrained region, as also shown in the figure. (As discussed in the text, this point represents a second, deflationary steady state.)

The signs of the elements of \underline{M} also imply that in the region above the $\dot{e}_1 = 0$ locus but below the $\dot{e}_2 = 0$ locus, all trajectories must move up and to the right, and so on. In particular, they imply that all trajectories starting in the grey region of the figure must move down and to the left. Hence all trajectories starting in this region remain trapped in it forever.

The signs of the elements of \underline{M} also imply that the determinant of the matrix is negative, implying that the matrix must have two real eigenvalues, one positive and one negative. It then follows from standard results regarding linear ODE systems that the stable manifold of the system defined by equation (A.25) is one-dimensional: it corresponds to a line (the line SM shown in Figure 1) passing through point e^* along which trajectories converge asymptotically to point e^* starting from any point on this line. Instead, all trajectories starting from points off this line diverge from the line, and hence diverge from the steady state e^* , which is thus locally unstable under the dynamics defined by (A.25).

Moreover, at any points in the constrained region that are below the line SM but above the $\dot{e}_2 = 0$ locus, $\dot{e}_2 < 0$, so that trajectories starting at any such point must eventually enter the grey region, unless they leave the constrained region (i.e., they cross the line ZLB , in which case equation (A.25) would no longer apply). Similarly, trajectories starting at any point below the line SM but to the right of the $\dot{e}_1 = 0$ locus must eventually enter the grey region, unless they leave the constrained region. Thus all trajectories beginning in the constrained region and below the line SM must eventually enter and be permanently trapped in the grey region, unless they leave the constrained region before entering the grey region.

It remains to determine whether trajectories beginning in the constrained region below the line SM can ever leave the constrained region. To answer this, we need to examine whether trajectories point into or out of the constrained region, at points along the line ZLB that defines the boundary of the region. The line ZLB consists of all points e such that

$$a'e = -\bar{r}\Delta; \tag{A.26}$$

the constrained region consists of all points e for which $a'e$ is more negative than this. Thus at any point on the line ZLB , the trajectory points out of the constrained region if $a'\dot{e} > 0$, and into the constrained region if instead $a'\dot{e} < 0$. At points between the point where the $\dot{e}_2 = 0$ locus intersects the line ZLB and the point where the $\dot{e}_1 = 0$ locus intersects it, both elements of \dot{e} are positive, and hence $a'\dot{e} > 0$. Hence

the trajectories all point outside the constrained region along this interval, as shown in Figure 1.

Consider instead the point where the line SM intersects the line ZLB , if such a point exists. The line SM consists of all points e of the form

$$e = e^* + \alpha u, \tag{A.27}$$

for arbitrary (positive or negative) values of the coefficient α , where u is the right eigenvector of the matrix \underline{M} corresponding to the positive eigenvalue. Thus

$$\underline{M}(e - e^*) = \lambda(e - e^*) \tag{A.28}$$

for some $\lambda > 0$. Then at any point on the line SM , we must have

$$\begin{aligned} a'\dot{e} &= -a'\underline{M}(e - e^*) \\ &= -\lambda a'(e - e^*), \end{aligned}$$

using (A.25) and (A.28). Because e^* is in the interior of the constrained region, $a'(e - e^*) > 0$ for all points on the line ZLB defined by (A.26). Hence at a point that is *both* on the line SM and on the line ZLB , we must have $a'\dot{e} < 0$, and the trajectory of the ODE system through this point points into the constrained region, as shown in Figure 1.

Now suppose that the stable manifold SM intersects the line ZLB at a point below and to the right of the point where the $\dot{e}_1 = 0$ locus intersects ZLB , as shown in the figure. Because $a'\dot{e} = -a'\underline{M}(e - e^*)$ is a linear function of e , the fact that it is positive at the point of intersection with the $\dot{e}_1 = 0$ locus but negative at the point of intersection with SM implies that it must also take an even more negative value at all points on the line ZLB that are below and to the right of the intersection with SM . Hence the ODE trajectories point into the constrained region at all such points, as shown in Figure 1, and it is not possible for a trajectory that begins in the constrained region and below the line SM to ever leave the constrained region.

Alternatively, suppose that the stable manifold intersects the line ZLB at a point above and to the left of the point where the $\dot{e}_2 = 0$ locus intersects ZLB . In this case, the fact that $a'\dot{e} > 0$ at the point of intersection with the $\dot{e}_2 = 0$ locus while $a'\dot{e} < 0$ at the point of intersection with SM implies that we must have $a'\dot{e} < 0$ at all points on the line ZLB above and to the left of the intersection with SM . Hence the ODE trajectories again point into the constrained region at such points, and again it is not possible for a trajectory that begins in the constrained region and below SM to ever leave the constrained region.

Finally, suppose that the stable manifold SM is exactly parallel to the line ZLB (the case in which $a'u = 0$). In this case, SM never intersects ZLB , and trajectories that begin below SM can never approach the boundary of the constrained region.

Thus in all possible cases, we conclude that trajectories that begin in the constrained region and below the line SM remain forever in the constrained region, as stated in the text (and illustrated in Figure 1). It follows that all such trajectories must eventually enter the grey region and remain trapped there forever, as stated in the text.

A.5 An Example of Rational Belief Equilibria

Here we provide additional details of the calculations referred to in the example discussed in section 4.1 of the text. It follows from the definition of the subjective variable p_t^{*j} in section 1.3 of the text that

$$p_t^{*j} = \sum_{k=0}^{\infty} (\alpha\beta)^k \hat{E}_t^j \left[\pi_{t+k} + \frac{\alpha}{1-\alpha} \kappa y_{t+k} \right].$$

(This can be obtained, for example, by “solving forward” equation (1.15) in the text, and then using equation (1.11) in the text to substitute for the $\hat{E}_t^j p_{t+k}^{*j}$ terms in terms of $\hat{E}_t^j \pi_{t+k}$.) This in turn implies that

$$\hat{E}_t^j p_{t+1}^{*j} = \sum_{k=0}^{\infty} (\alpha\beta)^k \hat{E}_t^j \left[\pi_{t+k+1} + \frac{\alpha}{1-\alpha} \kappa y_{t+k+1} \right]. \quad (\text{A.29})$$

Then substituting the expressions given in the text for the subjective expectations $\hat{E}_t^j \pi_{t+k}$ and $\hat{E}_t^j y_{t+k}$, we obtain a solution for $\hat{E}_t^j p_{t+1}^{*j}$ of the form conjectured in equation (4.7), with coefficients

$$\gamma_1 = \frac{1}{1-\alpha} \frac{\kappa\phi}{\kappa + \phi} \frac{\lambda}{1-\alpha\beta\lambda}, \quad (\text{A.30})$$

$$\gamma_2 = g \cdot \gamma, \quad g \equiv \frac{(1-\alpha)\phi - \alpha\kappa}{\kappa + \phi} \frac{\beta\rho}{1-\alpha\beta\rho}. \quad (\text{A.31})$$

This is not, however, a complete solution, as the expression obtained for γ_2 is a function of γ , the sum $\gamma_1 + \gamma_2$.

The solutions (A.30)–(A.31) imply that

$$\gamma = \gamma_1 + g\gamma, \quad (\text{A.32})$$

where γ_1 and g are both explicit functions of the model parameters given above. We further note that under our assumptions that $0 < \alpha, \beta, \lambda < 1$, $0 \leq \rho < 1$, and $\kappa, \phi > 0$, the coefficients of the linear equation (A.32) satisfy $\gamma_1 > 0$ and $g < 1$. Hence the equation has a unique solution, given by

$$\gamma = \frac{\gamma_1}{1-g} > 0. \quad (\text{A.33})$$

Substituting this solution into (A.30) and (A.31), we obtain unique solutions for γ_1 and γ_2 .

We thus obtain a unique solution for the TE dynamics under beliefs of the postulated form. The complete system of equations describing these dynamics consists of equations (4.1)–(4.2) in the text, specifying the exogenous evolution of the natural rate of output Y_t^n ; the equation

$$\bar{Y}_t = \lambda \bar{Y}_{t-1} + (1 - \lambda) Y_t^n \quad (\text{A.34})$$

for the dynamics of the central bank's estimate of the permanent component; equation (4.6) in the text for the evolution of the aggregate belief state ν_t ; and equations (4.8)–(4.9) in the text, giving the TE dynamics of inflation and the output gap, where the value of γ in these equations is given by (A.33).

Note that this does *not* mean that there is a unique equilibrium in this model consistent with the rational belief postulate. The assumption, in this derivation, that the aggregate belief state follows dynamics of the form (4.6) is an arbitrary one — not only our assumption that ν_t is an AR(1) process, but that it evolves independently of the fundamental disturbances — chosen purely to illustrate the kind of calculations that are involved in verifying the existence of an RBE. The assumption that each firm j regards its subjective belief state z_t^j as an accurate observation of the current value of \bar{Y}_t^n , rather than simply a random state that provides *information about* Y_t^n , is also an extreme special case. By relaxing either or both of these assumptions, we could construct a very large class of alternative RBE for this model, so that the predictions about such matters as the serial correlation of observed fluctuations in inflation are not nearly as sharp as the calculation above might suggest.

It may be useful to compare this example of RBE fluctuations to the REE dynamics implied by the assumed policy rule. Under the RE hypothesis, equation (1.20) implies that⁷²

$$p_t^* = \frac{\kappa}{1 - \alpha} y_t + \beta E_t p_{t+1}^*,$$

or alternatively (using equation (1.11) in the text), that

$$\begin{aligned} \pi_t &= \kappa y_t + \beta E_t \pi_{t+1} \\ &= \kappa (Y_t - \bar{Y}_t) + \kappa (\bar{Y}_t - Y_t^n) + \beta E_t \pi_{t+1}. \end{aligned} \quad (\text{A.35})$$

Then using equation (4.5) in the text to substitute for the $Y_t - \bar{Y}_t$ term, we obtain

$$\pi_t = \frac{\kappa \phi}{\kappa + \phi} (\bar{Y}_t - Y_t^n) + \frac{\beta \phi}{\kappa + \phi} E_t \pi_{t+1}$$

⁷²Recall that in this discussion, we assume that the cost-push disturbance u_t is equal to zero at all times.

as an equation that the REE dynamics of inflation must satisfy.

Since $\beta\phi/(\kappa + \phi) < 1$, this equation can be “solved forward” to yield

$$\pi_t = \frac{\kappa\phi}{\kappa + \phi} \sum_{k=0}^{\infty} \left(\frac{\beta\phi}{\kappa + \phi} \right)^k E_t[\bar{Y}_{t+k} - Y_{t+k}^n]. \quad (\text{A.36})$$

But since the RE forecasts of the future natural rate of output are given by equation (4.3) in the text, one has

$$E_t[\bar{Y}_{t+k} - Y_{t+k}^n] = 0$$

for all $k \geq 1$, so that (A.36) reduces to

$$pi_t = \frac{\kappa\phi}{\kappa + \phi} (\bar{Y}_t - Y_t^n). \quad (\text{A.37})$$

The implied solution for the output gap is then

$$y_t = \frac{\phi}{\kappa + \phi} (\bar{Y}_t - Y_t^n). \quad (\text{A.38})$$

Comparing these equations with equations (4.8)–(4.9) in the text for the RBE dynamics of inflation and the output gap, one sees that the REE dynamics are given by the same equations, but with the value $\gamma = 0$ instead of the positive value given in (A.33). Since the exogenous dynamics of the variables $\bar{Y}_t - Y_t^n$ and ν_t are the same in either case, this allows us to directly compare the RBE dynamics with the REE dynamics. As discussed in the text, we see that the responses of both π_t and y_t to the variations in exogenous fundamentals $\bar{Y}_t - Y_t^n$ are the same in both equilibria; the difference is that in the RBE dynamics, variations in the aggregate belief state ν_t (independent of the fluctuations in fundamentals) cause additional variation in π_t and y_t , causing each variable to be more volatile in the RBE than it would be under the REE dynamics.