

# Directed attention and nonparametric learning

Ian Dew-Becker and Charles G. Nathanson\*

September 19, 2017

## Abstract

We study an ambiguity-averse agent with uncertainty about income dynamics who chooses what aspects of the income process to learn about. The agent chooses to learn most about income dynamics at the very lowest frequencies, which have the greatest effect on utility. Deviations of consumption from the full-information benchmark are then largest at high frequencies, so consumption responds strongly to predictable changes in income in the short-run but is closer to a random walk in the long-run. Whereas ambiguity aversion typically leads agents to act as though shocks are more persistent than the truth, endogenous learning here eliminates that effect.

A growing literature studies economic behavior in the face of model uncertainty, while at the same time there is a large amount of recent work that studies optimal allocation of attention.<sup>1</sup> Those two areas are obviously related: the economy is highly complex, so people are unlikely to be able to understand all of it, and they must choose how to allocate their limited attention and information processing abilities. Furthermore, information acquisition is not free, so we would not necessarily expect people to be perfectly informed about everything.

Surprisingly, though, there is little or no research that studies the implications of directed attention in the face of model uncertainty.<sup>2</sup> The contribution of this paper is to study the behavior of an agent who allocates attention across different aspects of a model. We show that optimal learning about model features has important and interesting implications for behavior. On the one hand, it naturally leads to excess sensitivity of consumption to high-frequency components of income, as observed empirically. At the same time, though, we show that optimally directed attention drives the consumption policy closer to the optimum at lower frequencies than it would be if model uncertainty were purely exogenous.

---

\*Dew-Becker: Northwestern University and NBER; [ian.dewbecker@gmail.com](mailto:ian.dewbecker@gmail.com). Nathanson: Northwestern University; [nathanson@kellogg.northwestern.edu](mailto:nathanson@kellogg.northwestern.edu). We appreciate helpful comments from Ben Hebert, Peter Klibanoff, Konstantin Milbradt, Mikkel Plagborg-Møller, and seminar participants.

<sup>1</sup>On model uncertainty, most relevant for us is recent work on consumption under model uncertainty, e.g. Hansen, Sargent, and Tallarini (1999), Wang (2004, 2009), Luo (2008), Luo and Young (2010), but there is a large broader literature. See Sims (2003), Veldkamp (2011), and many citations therein for work on directed attention.

<sup>2</sup>There is substantial past work on directed learning (e.g. Van Nieuwerburgh and Veldkamp (2006), Peng and Xiong (2006), Veldkamp (2006), and Barron and Ni (2008)), but we are not aware of work that examines the choice of what part of a dynamic process to learn about.

More concretely, we study the problem of an ambiguity averse agent who is uncertain about the dynamics of an exogenous and untradeable income process. The key innovation compared to past work on ambiguity aversion is that the agent acquires information that can reduce the degree of ambiguity.

The agent's optimization has three phases. Conditional on a particular model of the world, the agent has standard Bayesian expected utility.<sup>3</sup> If the agent had a fully specified probability distribution over possible models for income, she could integrate across them, remaining a Bayesian. Similar to the literature on robust control (e.g. Hansen and Sargent (2007)), we argue that such behavior is implausible in the face of infinite-dimensional uncertainty about income dynamics. We therefore model the agent as ambiguity averse: among all sufficiently plausible models, she adopts the one under which her expected consumption utility is the smallest (this is the second phase of the optimization). This selection criterion ensures that the agent's consumption decisions are robust to uncertainty about the true model

The third phase of the optimization is most important. The agent allocates attention to different aspects of the income process, which allows her to endogenously limit the degree of ambiguity she faces. When the agent pays more attention to a particular aspect of income, such as its low-frequency behavior, she receives information about its true behavior along that dimension and the set of plausible models narrows.

We solve three phases of the optimization analytically and are therefore able to sharply establish our main result: the agent directs almost all attention to the behavior of income at the lowest frequencies (i.e. at long horizons). Because of this direction of attention, the agent's model is more accurate at lower frequencies than at higher frequencies. The agent learns in this manner because the low-frequency dynamics of income matter more for expected consumption utility than the high-frequency dynamics. By learning about low-frequency dynamics, the agent can deem many painful income processes implausible and therefore avoid selecting them in the second phase of her optimization.

Using this intuitive result, we derive two key implications for the agent's beliefs and behavior:

1. At short horizons, the agent's consumption growth is positively correlated with the *predictable* component of income growth. This comovement violates the permanent income hypothesis (Friedman (1957), Hall (1978)) but matches the extensive empirical evidence on the excess sensitivity of consumption to income (Jappelli and Pistaferri (2010), Kaplan and Violante (2014)). Because the agent fails to learn about the high-frequency characteristics of the income process, much of the predictable variation in income surprises her and therefore leads her to adjust her consumption.
2. The agent neither over- nor under-extrapolates current shocks to income when forecasting

---

<sup>3</sup>During this phase of the optimization, no dynamic learning about the model occurs. For boundedly rational models of dynamic learning, see Abel, Eberly, and Panageas (2007, 2013), Wang (2009), Bansal and Shaliastovich (2010), Hansen and Sargent (2010), Ju and Miao (2012), and Collin-Dufresne, Johannes, and Lochstoer (2015).

long-run future income. This lack of bias results from two offsetting forces in the agent’s optimization. Ambiguity aversion pushes her towards adopting an overly persistent model, as is typical in the literature on model selection under ambiguity aversion (Hansen and Sargent (2010, 2015), Bidder and Dew-Becker (2016)).<sup>4</sup> But the agent’s extra attention to the low-frequency behavior of income undoes this bias, as her knowledge attenuates her fears of persistent income processes.

What connects the two theoretical results is that high-frequency mistakes have minimal implications for lifetime utility, while low-frequency mistakes can have substantial effects. That idea has been suggested as an explanation for the excess sensitivity puzzle, and our model formalizes it.<sup>5</sup> People cannot achieve perfection, so they choose to make mistakes that are minimally costly.

After establishing these theoretical results, we explore them in a numerical example to quantify the importance of directed attention in guiding beliefs and behavior. As a comparison, we study the model adopted by the agent when she is restricted to allocate the same amount of attention to each part of the spectrum. This “fixed-attention” model provides the optimal statistical fit of the true spectrum and corresponds to the boundedly rational model studied by Fuster, Hebert, and Laibson (2011). In addition to confirming the theoretical results numerically, we find that attention allocated to low frequencies is 80 times higher under directed attention than fixed attention.

## 1 Environment and information

We begin by laying out preferences. We then describe the space of income processes, and finally the structure of the uncertainty that the agents face.

### 1.1 Preferences

We study agents who solve a consumption-savings problem under ambiguity aversion over model uncertainty. They face the following budget constraint.

**Assumption 1** *Financial wealth,  $W_t$ , follows the process*

$$W_t = RW_{t-1} + Y_t - C_t \tag{1}$$

where  $C_t$  is consumption,  $Y_t$  is income, and  $R$  is a fixed gross interest rate.

We denote possible income processes  $\hat{f}$ .

The agent’s preferences are represented by the following optimization.

---

<sup>4</sup>A bias towards belief in overly persistent processes is present also in the boundedly rational frameworks of Fuster, Hebert, and Laibson (2011) and Bordalo, Gennaioli, and Shleifer (2016).

<sup>5</sup>See Cochrane (1989), Eichenbaum (2011), and Kueng (2016) for discussions of the small utility costs of excess sensitivity to transitory income shocks.

**Assumption 2** Agents choose signal precision  $\tau$  and a consumption policy  $C^{policy}$  according to

$$\max_{\tau} E \left[ G \left( \max_{C^{policy}} \min_{\hat{f} \in F(x; \tau)} E \left[ \sum_{t=0}^{\infty} -\alpha^{-1} \beta^t \exp(-\alpha C_t) \mid \hat{f} \right] \right) \right], \quad (2)$$

where  $E$  denotes the expectation operator and  $G$  is a strictly increasing function that will be defined in equation (26) below.  $C^{policy}$  is a typical consumption rule mapping current wealth,  $W_t$ , and the income history,  $Y_0, \dots, Y_t$ , to consumption,  $C_t$ . The optimization is subject to the budget constraint (1) and a transversality condition.

The inner maxmin pair represents ambiguity averse preferences similar to those of Gilboa and Schmeidler (1989). The agent's aim is to maximize discounted utility over consumption, where  $\alpha$  is the coefficient of absolute risk aversion and  $\beta$  the time discount factor.

The source of uncertainty that the inner expectation applies to is the future realizations of income. Conditional on a (functional) parameter  $\hat{f}$ , agents calculate expectations over income realizations, and hence future consumption, using Bayes' rule.

The parameter  $\hat{f}$  is unknown. If people were Bayesian expected utility maximizers, they would choose the consumption policy to maximize expected utility under a probability measure for  $\hat{f}$ . That is, we would have

$$\max_{C^{policy}} \int E \left[ \sum_{t=0}^{\infty} -\alpha^{-1} \beta^t \exp(-\alpha C_t) \mid \hat{f} \right] d\Phi(\hat{f}) \quad (3)$$

where  $d\Phi(\hat{f})$  represents a probability density over models.  $\hat{f}$ , is a potentially infinite dimensional object. We therefore take the position that it is not reasonable to assume that people are able to fully articulate a probability distribution over all possible values of  $\hat{f}$  (the work of Hansen and Sargent (2007) on robust control is motivated similarly).

Instead, we model agents as ambiguity averse. They believe that  $\hat{f}$  may fall into a set  $F(x; \tau)$ , where  $x$  is a set of signals about the true model that they receive, which have precision  $\tau$ . The consumption policy is chosen to maximize expected utility with the understanding that nature will then select the least favorable value of  $\hat{f} \in F(x; \tau)$ .

The outer expectation is taken over possible realizations of the signals  $x$ . The agent chooses the signal precisions,  $\tau$ , to maximize the expected outcome of the ambiguity-averse consumption/savings problem. Intuitively, an agent that receives high-quality signals about income dynamics will have a smaller set  $F(x; \tau)$ , thus reducing the effects of ambiguity. The function  $G$  is applied to expected utility conditional on the signals for reasons that we discuss below.

Note that if there were no model uncertainty, so that the true model  $f$  is known, then the agent is solving a standard consumption-savings problem under CARA preferences:  $\max E \left[ \sum_{t=0}^{\infty} -\alpha^{-1} \beta^t \exp(-\alpha C_t) \right]$ . Our analysis therefore ignores wealth effects, but it is also more realistic than the assumption of quadratic utility over consumption used in Hansen, Sargent, and Tallarini (1999), among others.

Finally, it is also important to note that this is in certain regards a date-0 problem. Agents receive signals about the income process once. They then choose an optimal consumption policy that is meant to be robust to model uncertainty. We do not model how people update information about the income process ( $\hat{f}$ ) over time. That said, consumption is chosen fully dynamically in that the policy conditions on the past histories of wealth and income.

The remainder of this section defines more formally the various terms in assumption 2. We then proceed to solve the three optimization problems in section 3 and examine their implications in section 4.

## 1.2 Income

**Assumption 3** *Consumers face an exogenous and untradeable stochastic income stream,  $Y_t$ , that follows the process*

$$Y_t = a(L)Y_{t-1} + b_0\varepsilon_t \quad (4)$$

$$\varepsilon_t \sim i.i.d. N(0, 1) \quad (5)$$

where  $a(L)$  is a power series in the lag operator,  $L$ . We assume  $a(L)$  is such that  $Y$  is well behaved (in particular, has a spectrum that is positive and bounded).<sup>6</sup>

This is a standard baseline setup for time series models. While linearity and Gaussianity are certainly restrictive assumptions, they are in line with the past work we build on.

Much of our analysis will apply to the Wold representation,

$$Y_t = b(L)\varepsilon_t, \quad (6)$$

$$\text{where } b(L) \equiv \frac{b_0}{1 - La(L)}. \quad (7)$$

The coefficients in the power series  $b(L)$  are denoted  $b_j$  (i.e.  $b(L) = \sum_{j=0}^{\infty} b_j L^j$ ). Throughout the paper, we refer to models in the time domain in terms of  $b(L)$ . Since the distribution of  $\varepsilon_t$  is fixed,  $b(L)$  completely characterizes the statistical distribution of income. To be clear, though, the agent forecasts the future using only the past history of income. The  $\varepsilon_t$  are not directly observable.

Agents do not know the true income process. Alternative possible income processes are denoted  $\hat{b}$ .<sup>7</sup> Our focus is on uncertainty about the dynamics of income, rather than about the distribution of

---

<sup>6</sup>The assumption that  $Y_t$  is a linear Gaussian process is not necessary for most of the results. The critical assumptions about the true income process are that it is second-order stationary and that it has a spectral density that is finite and bounded away from zero. The distribution of the innovations is largely irrelevant (though it is important that it is fixed over time).

<sup>7</sup>Agents forecast with  $\hat{b}$  the same way they would with  $b$ ,

$$E_t \left[ Y_{t+n} \mid \hat{b} \right] = \sum_{k=0}^{\infty} \hat{b}_{n+k} \hat{b}(L)^{-1} Y_{t-k} \quad (8)$$

its innovations. The latter question is obviously also interesting, but our goal here is to understand how consumption responds to changes in income, and how well people understand the difference between permanent and transitory dynamics. An interpretation of our analysis is that it derives optimal attention to different aspects of income dynamics conditional on a choice having been made about how much attention to pay overall to dynamics versus the distribution of income innovations.

### 1.3 Signals about the spectrum of income

The key type of uncertainty that agents face is over the model that drives income. The definitions above are in the time domain, but our analysis examines a rotation, using the Fourier transform, into the frequency domain. The Fourier transform is used in time series analysis because it orthogonalizes problems. In our setting, one way for an agent to model income is to estimate its autocovariances ( $cov(Y_t, Y_{t-j})$ ). But estimates of autocovariances are in general correlated across lags, and one must impose complicated restrictions to guarantee positive definiteness, both of which substantially complicate the analysis.

Those issues do not arise in the frequency domain. In what we describe here, agents receive signals about features of the income process that are mutually orthogonal and always generate a positive definite covariance matrix for income. So a primary reason that we model agents as learning in the frequency domain is that such learning represents acquiring information about fundamentally independent aspects of the income process. We will show, though, that there is a direct link between estimation in the frequency and time domains, so there is no deep restriction on behavior; rather, the transform is used yield analytic and economically interpretable solutions.

The next subsection defines the frequency transform. We then describe precisely how agents acquire information about income dynamics. Last, we discuss the relationship between the information scheme used here and those studied in models of rational inattention.

#### 1.3.1 The Fourier transform and its distribution

The spectral density of the income process is defined as the Fourier transform of the autocovariances and thus also fully represents income dynamics:

$$\exp f(\omega) \equiv \sum_{j=-\infty}^{\infty} \cos(\omega j) \text{cov}(Y_t, Y_{t-j}). \quad (9)$$

(the notation  $f$  is used for the log spectrum because that is what maps most directly into utility). As  $f$  is periodic, we restrict attention to the domain  $\omega \in [0, \pi]$ . There is a one-to-one mapping between the spectral density and the autocovariances since they are a Fourier transform pair. One simple interpretation of the spectrum is that it represents a variance decomposition for income in

terms of fluctuations at different frequencies:

$$\text{var}(Y_t) = \frac{1}{\pi} \int_0^\pi \exp(f(\omega)) d\omega. \quad (10)$$

$\exp(f(\omega))$  measures the contribution of fluctuations in income at frequency  $\omega$  to the total variance of income. The relative magnitude of  $f$  across frequencies determines the extent to which variation in income is driven by low- versus high-frequency fluctuations. An AR(1) process with an autocorrelation near 1 has a spectrum whose mass is isolated at low frequencies, whereas a process that features reversals, such as  $x_t = \varepsilon_t - (1/2)\varepsilon_{t-1}$ , has a spectrum with mass concentrated at high frequencies (those near  $\pi$ ).

As with  $b$  and  $\hat{b}$ ,  $f$  is the true log spectral density of income, while alternative hypothetical spectra are denoted  $\hat{f}$ . The agent can construct forecasts of future income based on  $\hat{f}$  since there is a unique  $\hat{b}$  associated with each  $\hat{f}$  (Priestley (1981) section 10.1).

The key feature of the spectrum for our purposes is that sample estimates of it are asymptotically uncorrelated across frequencies:

**Lemma 1** *Denote the sample autocovariance of a mean-zero time series*

$$\tilde{v}_j \equiv (T-j)^{-1} \sum_{t=1}^{T-j} Y_{t+j} Y_t, \quad (11)$$

and define  $\exp(\tilde{f}(\omega))$  to be the sample analog of (9). Then as  $T \rightarrow \infty$ ,

$$E[\tilde{f}(\omega) - f(\omega)] \rightarrow -\gamma \quad (12)$$

$$\text{cov}(\tilde{f}(\omega_1) - f(\omega_1), \tilde{f}(\omega_2) - f(\omega_2)) \rightarrow \frac{\pi^2}{6} 1\{\omega_1 = \omega_2\} \quad (13)$$

for  $\omega \in [0, \pi]$ , where  $\gamma$  is Euler's constant and  $1\{\cdot\}$  is the indicator function.

**Proof.** This result follows directly from Brillinger (1981) theorem 5.2.6. ■

Lemma 1 shows why we study income dynamics in the frequency domain: the sample spectrum yields an estimate of the true spectrum with errors that are asymptotically uncorrelated across frequencies. In what follows, we assume people use the asymptotic approximation and treat the spectral estimates as truly uncorrelated across frequencies.

Intuitively, estimating the sample spectrum is not fundamentally different from estimating sample autocovariances – they are linear transformations of each other and hence contain the same information. So in estimating the spectrum, people are effectively learning about the properties of income by directly calculating its autocovariances.

### 1.3.2 Information acquisition method

Information acquisition happens prior to the consumption policy being chosen; this step can be thought of as happening behind a veil of ignorance, before realizations of the agent's own income history have occurred. We assume that people gather information by estimating the log spectrum on observed income histories driven by the same model that will drive their own income. For example, a person who knows they will become an economist might investigate the income histories that older economists have received. These observations and calculations are costly, though, so we assume that people are limited in the number of measurements of spectra (equivalently, calculations of autocovariances) that they can undertake.

More specifically, we assume that there is a large dataset available that reports the income histories of many people, all of whom have the same parameters determining their income processes (i.e. the same  $f$  and hence  $b$ ), but different realizations (different  $\varepsilon$ 's). To obtain information about the log spectrum of income at some frequency  $\omega_j$ , the agent calculates the sample spectrum for  $\tau_j$  of those income histories.

Combining the central limit theorem with lemma 1 above, the mean of  $\tilde{f}(\omega_j) + \gamma$  across many income histories is asymptotically normally distributed with mean  $f(\omega_j)$  and variance  $\frac{\pi^2}{6}\tau_j^{-1}$  (we ignore the  $\frac{\pi^2}{6}$  scaling in what follows). To represent limited information capacity, we assume that the sum of the total calculations that an agent may perform across all frequencies is limited. That is,  $\sum_j \tau_j$  is limited (section 4.5 generalizes the constraint to allow for differential information costs across frequencies).<sup>8</sup>

For technical reasons (to avoid infinite information flows, for example), we assume that the agent gains information on the spectrum on the uniform discretization of  $[0, \pi]$  given by  $\omega_j = \pi j/n$  for  $j \in \{1, \dots, n\}$  and we take  $n$  as large. We scale the variances by  $d\omega \equiv \pi/n$  so that they can be interpreted as the information density at each point.

**Assumption 4** *The agent receives signals  $\{x(\omega_j)\}_{j=1, \dots, n}$  that are distributed as*

$$x(\omega_j) \sim N\left(f(\omega_j), \tau(\omega_j)^{-1}/d\omega\right) \quad (14)$$

where the errors are uncorrelated across frequencies. In choosing the precision of their signals, agents face the constraint

$$\sum_{j=1}^n \tau(\omega_j) d\omega \leq \bar{\tau}. \quad (15)$$

If  $\tau$  differs across frequencies, that means that the agent calculated the sample spectrum for more income histories at some frequencies than others. That is, they have many income histories

---

<sup>8</sup>Note that in the interest of conserving effort, an agent can calculate the spectrum most efficiently by simply calculating the squared Fourier transform of the income history itself – that is numerically equivalent to taking the Fourier transform of the sample autocovariances. The constraint on  $\sum_j \tau_j$  is then a constraint on the number of Fourier transforms the agent may calculate.



to examine, but for frequencies that they learn less about, they only estimate the spectrum using a small number of them, and the effort saved is allocated elsewhere.

A natural statistical benchmark is for the agent to instead use equal information at all frequencies – i.e. to calculate the entire spectrum for each income history – making no attention allocation decision. In that case, the function  $\tau$  is constant across frequencies.

The estimation scheme that we endow agents with is the standard nonparametric method for estimating the spectrum. If agents knew the specific parametric form of the income process, e.g. that it was an ARMA( $p, q$ ) for known  $p$  and  $q$ , then they could estimate it more efficiently through maximum likelihood. We instead leave their model uncertainty unstructured, which makes nonparametric analysis based on the periodogram most natural. The lack of structure is in fact precisely what delivers the independence across frequencies that will make our analysis tractable.

### 1.3.3 Relationship with rational inattention

Rational inattention provides an alternative and equally important interpretation of the information structure. It is possible that complete information about the spectrum of income is available, but agents have trouble processing it. Then the noise in the signals  $x(\omega_j)$  represents cognitive errors that people make in interpreting the available information. The frequencies at which  $\tau$  is larger are the ones the agent pays the most attention to.

In terms of the literature, the signal structure we analyze is highly similar to that in Kacperczyk, Van Nieuwerburgh, and Veldkamp (2016) in that agents receive signals with normally distributed error and they are constrained by the total precision of the signals. This constraint is most natural when each independent observation of the spectrum is equally costly to obtain. Sims (2003) proposes an alternative constraint based on information flow or entropy. In our setting, the total entropy of the signals is  $\sum_{j=1}^n \log(\tau(\omega_j) d\omega)$ , so high-precision signals are relatively less costly under an entropy constraint.

That said, our setting is more restricted than the fully general rational inattention specification: the information the agents acquire is always independent across frequencies (which is motivated by the fact that statistical estimates of the spectrum are independent across frequencies) and the errors are Gaussian (motivated by the central limit theorem). In the most general form of the models that Sims (2003) studies, those restrictions need not hold. However, they are commonly imposed elsewhere, as in Kacperczyk, Van Nieuwerburgh, and Veldkamp (2016).

## 1.4 Priors and model plausibility

Agents measure the plausibility of models, and define the set they worry about,  $F(x; \tau)$ , based on their signals and a prior. Given that the model space is infinite-dimensional, it is difficult to imagine that a person would have a fully defined prior, though. People likely cannot place a formal probability on every possible model, or even necessarily express a view about the relative likelihood of all possible pairs of models. That fact motivates our use of ambiguity aversion, and it leads us

to specify prior beliefs as loosely as possible.

We assume agents believe that the log spectrum is likely to be smooth in the sense that its differences across frequencies have limited variation. The smoothness prior is a belief in simplicity: agents believe that spectra typically are smooth across frequencies, rather than fluctuating wildly. Following Shiller (1973), Akaike (1979), and Kitagawa and Gersch (1984, 1996), the prior is represented by a penalty on variability that is appended to the likelihood of the data.<sup>9</sup> Given assumption 4, the penalized log likelihood of the data given a model  $\hat{f}$  is

$$PL(x | \hat{f}, \tau) = \underbrace{-\frac{1}{2} \sum_{j=1}^n \left( x(\omega_j) - \hat{f}(\omega_j) \right)^2 \tau(\omega_j) d\omega}_{\text{Data likelihood}} - \underbrace{\frac{\lambda}{2} \sum_{j=2}^n \left( \frac{\hat{f}(\omega_j) - \hat{f}(\omega_{j-1})}{d\omega} \right)^2 d\omega}_{\text{Roughness penalty}}. \quad (16)$$

$PL(x | \hat{f}, \tau)$  depends on two factors: the log likelihood for normally distributed data and a term encoding the belief in smoothness. Models are viewed as less plausible when they are rougher or more complicated in the sense of having a larger average squared derivative. The most plausible models have perfectly flat spectra – white noise – while the least plausible have highly variable spectra.<sup>10,11</sup>

The parameter  $\lambda$  controls the strength of prior. For any fixed  $\lambda$ , as the signal precision grows large, the smoothness penalty becomes irrelevant. One reason we include the smoothness prior is that without it,  $\hat{f} = x$  is the maximum-likelihood estimate, which would imply that  $\hat{f}$  has infinite variation and would yield an inconsistent estimate of  $f$ , even as  $n \rightarrow \infty$  (Wahba (1980)).

The smoothness prior also implies that when people have weak signals, they use simple and smooth models. Complexity here only arises when people have a wealth of information. When signals are more precise, so that  $\tau$  is large relative to  $\lambda$ , the roughness penalty is relatively less important and the agent will consider more complex models.

The penalized likelihood leads to the following assumption:

**Assumption 5** *Nature chooses the income process from the set*

$$F(x; \tau) = \left\{ \hat{f} : PL(x | \hat{f}, \tau) \geq \bar{L} \right\} \quad (17)$$

In addition to the roughness penalty, we also assume that agents are able to express a prior mean over possible models. In the absence of any information about the world, they believe the

<sup>9</sup>The smoothness prior is often explicitly justified as a belief in simplicity. In Shiller (1973), which is the first application of such a prior, a justification is that “[i]n most applications...the researcher will feel that...the lag coefficients should trace out a ‘smooth’ or ‘simple’ curve.” While Shiller’s (1973) smoothness prior is stated in the time domain, those in Akaike (1979) and Kitagawa and Gersch (1985, 1989) are specified in the frequency domain in a manner almost identical to ours.

<sup>10</sup>That white noise is treated as the most plausible is also sensible from an information theoretic perspective since Gaussian white noise has the greatest Shannon entropy among all time series processes with a given variance.

<sup>11</sup>An alternative way to penalize complexity in models would use the coefficients of the ARMA representation. We will see below, though, that the smoothness prior we impose here also ends up imposing smoothness on the AR and MA coefficients.

average spectrum is flat at  $\bar{f}$ . This assumption is introduced so that it is possible for the agent to calculate expectations for  $\hat{f}$  prior to observing signals (i.e. the outer expectation operator in assumption 2).

At this point, all the basic terms in the preferences in assumption 2 have been defined. The next section examines the three optimizations.

## 2 Solution

All three optimizations in the preferences – the consumption policy, nature’s choice of a model, and the information decision – are analytically solvable. The solution is itself an important contribution of the paper. There is little work that obtains closed-form solutions for optimal consumption under model uncertainty and rational inattention, and the fact that the model can be solved when model uncertainty and attention are themselves endogenous is even more surprising. We analyze the three pieces of the optimization in turn.

### 2.1 Optimal consumption conditional on a model

The minimax theorem implies that the inner maximization and minimization in the preferences (2) can be reversed. Intuitively, the operations represent a zero-sum game with a pure strategy Nash equilibrium. We first solve for the optimal consumption given a model.

**Lemma 2** *The consumption policy that maximizes expected utility conditional on some  $\hat{b}$  is*

$$C_t = (R - 1) (W_{t-1} + \hat{z}(L) \hat{\varepsilon}_t) - \frac{\alpha}{2} R^{-1} (1 - R^{-1}) \hat{b} (R^{-1})^2 - \alpha^{-1} \frac{\log \beta R}{R - 1}, \quad (18)$$

where  $z(L)$  is a lag polynomial with coefficients

$$\hat{z}_j = \sum_{k=j}^{\infty} R^{-(k-j)} \hat{b}_k. \quad (19)$$

Expected utility from consumption is then

$$\begin{aligned} & \max_{C^{policy}} E \left[ -\alpha^{-1} \sum_{t=0}^{\infty} \beta^t \exp(-\alpha C_t) \mid \hat{b} \right] \\ &= \frac{-\alpha^{-1}}{1 - \beta} \exp \left( \frac{\alpha^2}{2} R^{-1} (1 - R^{-1}) \hat{b} (R^{-1})^2 + \log \frac{(1 - \beta)}{1 - R} + \frac{\log \beta R}{R - 1} \right). \end{aligned} \quad (20)$$

Lemma 2 provides two useful results. First, it shows that we obtain a standard consumption function: agents consume the annuity value of financial plus human wealth,  $(R - 1) (W_{t-1} + \hat{z}(L) \hat{\varepsilon}_t)$ , minus a precautionary saving term  $\frac{\alpha}{2} R^{-1} (1 - R^{-1}) \hat{b} (R^{-1})^2$ . The behavior of consumption thus depends on beliefs about income dynamics through two channels. First,  $\hat{b}$  affects the riskiness of

the income stream, and hence the amount of precautionary savings agents desire to hold. Second, and more importantly, current consumption depends on beliefs about future income, which are driven by  $\hat{b}$ . When  $\hat{b}$  implies that income shocks are more persistent ( $\sum_{k=j}^{\infty} R^{-(k-j)} \hat{b}_k$  is larger) consumption responds more strongly to the shocks. These are all standard results. Deviations of the behavior of the agents in our model from the standard permanent-income predictions are caused by deviations of their model,  $\hat{b}$ , from the truth.

Lemma 2 also characterizes optimized expected utility from consumption for a given income process  $\hat{b}$ . The only term that differs across models is  $\hat{b}(R^{-1})^2$ , which measures the variance of innovations to permanent income, and hence the variance of consumption growth. Utility is lower when the variance of consumption growth is higher.

The information structure laid out in the previous section refers entirely to the log spectrum, but utility is derived in lemma 2 terms of the lag polynomial  $\hat{b}$ . We link the two through the following novel result.

**Lemma 3** *For a log spectrum  $\hat{f}$  that is bounded from above and below, where  $\hat{b}(L)$  is the associated Wold representation,*

$$\log \hat{b}(R^{-1})^2 = \frac{1}{\pi} \int_0^{\pi} Z(\omega) \hat{f}(\omega) d\omega, \quad (21)$$

$$\text{where } Z(\omega) \equiv 1 + 2 \sum_{j=1}^{\infty} \cos(\omega j) R^{-j}. \quad (22)$$

Lemma 3 gives us a powerful result:  $\log \hat{b}(R^{-1})^2$ , the statistic that determines expected utility from consumption conditional on a model, is linear in the log spectrum. This result is the key novel mathematical innovation in the paper that will allow us to solve the model analytically, and it is likely useful in other contexts, since it is general result for NPV innovations.<sup>12</sup>

Lemma 3 shows that utility is always decreasing in  $\hat{f}$ , which implies utility is decreasing in the variance of income growth. Moreover, though, utility depends on different frequencies differently, according to the function  $Z$ . The frequency domain is useful here for showing how risk at different frequencies drives utility. The left-hand panel of figure 1 plots  $Z$  for an annual calibration with  $R = 1.025$ .  $Z(\omega) > 0$  for all  $\omega$ , it is bounded from above for  $R > 1$ , reaching its maximum at  $\omega = 0$ , and it is decreasing on  $(0, \pi)$ . When  $R = 1$ ,  $Z$  is equivalent to the Dirac delta function. The mass of  $Z$  primarily lies on extremely low frequencies. So what matters for the agent's utility, through

---

<sup>12</sup>Lemma 3 does not appear to have been previously noted in the literature, and we are not aware of any direct derivation from known results. It is a generalization of the Szegő-Kolmogorov formula for the innovation variance of a time series. Specifically,  $\hat{b}(0)^2$  is the innovation variance, which the Szegő-Kolmogorov formula says is the geometric mean of the spectrum. The equation  $\hat{b}(1)^2 = \exp \int \delta(\omega) f(\omega) d\omega$  for the Dirac delta function  $\delta$  is also well known. So lemma 3 fills in  $\hat{b}(x)^2$  for  $x$  between 0 and 1.

The innovation variance for the NPV of a time series arises naturally in many economic settings, such as the consumption/savings problem here, equilibrium macroeconomic models (Hansen and Sargent (1980, 1981)), models with generalized recursive preferences (Bidder and Dew-Becker (2016); Dew-Becker and Giglio (2016); Dew-Becker (2016)), the q theory of investment, and Calvo-type price setting. The appendix provides a proof.

$\hat{b}(R^{-1})^2$ , is the variance of the most persistent components of income. Transitory fluctuations in income do not pass into consumption. Rather, permanent income shocks change human wealth, and thus consumption, hence reducing utility. These characteristics of the utility function and  $Z$  are robust features of the model, as they do not depend on any sort of detailed calibration – the only parameter affecting  $Z$  is the gross interest rate.<sup>13</sup>

## 2.2 Nature’s minimization

Since  $\hat{b}(R^{-1})^2$  is the only term in (20) that differs across models, nature’s minimization problem in (2) is equivalent to choosing  $\hat{f}$  from the set  $F(x; \tau)$  to maximize  $\int_0^\pi Z(\kappa) \hat{f}(\kappa) d\kappa$ . Nature’s Lagrangian is

$$\min_{\hat{f} \in F(x; \tau)} -\frac{1}{\pi} \int_0^\pi Z(\omega) \hat{f}(\omega) d\omega - \psi PL(x | \hat{f}, \tau), \quad (23)$$

where  $\psi$  is a Lagrange multiplier. We refer to the model that achieves the minimum in (23) as  $f^w(\omega; x, \tau)$ .

It is straightforward to solve for  $f^w$  from the first-order conditions for the nature’s optimization. The solution can be obtained most easily by creating vectors (in boldface) of the form  $\mathbf{f}^w(x; \tau) \equiv [f^w(\omega_1; x, \tau), \dots, f^w(\omega_n; x, \tau)]'$  (recall that the frequencies  $\omega_j = \pi j/n$  are the uniform discretization of the interval  $[0, \pi]$  on which the agent receives signals and that we think of  $n$  as large). We define  $diag(\cdot)$  to be an operator that creates a matrix with its argument on the main diagonal and zeros elsewhere.

**Proposition 1** *The model that solves (23) is*

$$\mathbf{f}^w(x; \tau) = (I_{n \times n} - \lambda diag(\tau^{-1}) D)^{-1} (\psi diag(\tau^{-1}) \mathbf{Z} + \mathbf{x}) \quad (24)$$

where  $I_{n \times n}$  is an  $n \times n$  identity matrix and  $D$  is a differencing matrix of the form

$$D \equiv \begin{bmatrix} -1 & 1 & 0 & \cdots & 0 \\ 1 & -2 & 1 & & \\ 0 & 1 & -2 & \ddots & \vdots \\ \vdots & & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & 1 & -1 \end{bmatrix} d\omega^{-2}. \quad (25)$$

$f^w$  is a linear function of  $x$  and  $Z$ . The worst-case spectrum is higher at frequencies where  $\tau$  is smaller – there is more uncertainty about the spectrum – and where  $Z$  is larger – increases in the

<sup>13</sup>The analysis so far has assumed income is stationary. That assumption has no effects on our results. In the presence of a unit root, the analysis applies to the first difference of income. If  $\hat{g}(L)$  is the Wold representation for the first difference of income, then  $\hat{b}(R^{-1}) = \hat{g}(R^{-1}) / (1 - R^{-1})$ . The agent then can calculate  $\log \hat{b}(R^{-1})^2$  by using Lemma 3 applied to the log spectrum of income *growth* and subtracting  $\log(1 - R^{-1})$ . The loading of utility on frequencies for the level of income is the same as for the first difference.

spectrum are more painful. Similarly, when  $\psi$  is larger, so that the agent is more ambiguity-averse, the worst-case model tilts more in the direction of  $Z$ .

Before analyzing the implications of proposition 1 in detail, we first solve the agent's optimal  $\tau$  to help frame the effects of information choice on consumption behavior.

### 2.3 Optimal information choice

Consistent with nature's decision, we assume that agents choose information so as to minimize  $E \left[ b^w (R^{-1})^2 \right] = \frac{1}{\pi} \int_0^\pi Z(\omega) E[f^w(\omega; x, \tau)] d\omega$ . The reason we choose that particular function to minimize, beyond the fact that it is the measure of risk that drives utility, is that it can be calculated under the incomplete prior that agents have. Recall that agents do not have a fully specified prior over  $f$ ; they only have a mean and a belief in smoothness. But since  $f^w$  is linear in  $x$ , its mean is well defined under those beliefs. So assuming agents minimize  $b^w (R^{-1})^2$  allows us to specify the least informative prior possible – one that has only a mean.

The assumption that information is chosen in order to minimize the expected riskiness of income immediately yields a functional form for  $G$ :

#### Assumption 6

$$G(x) \equiv \log \log \left( -\alpha (1-R) (\beta R)^{1-R} x \right) + \log \left( 2\alpha^{-2} R (1-R^{-1})^{-1} \right) \quad (26)$$

While that function looks complicated, its purpose is that it yields a simple reduced-form version of the preferences (2):

$$\max_{\tau} E \left[ \min_{\hat{f} \in F(x; \tau)} -\frac{1}{\pi} \int_0^\pi Z(\omega) \hat{f}(\omega) d\omega \right] = \max_{\tau} E \left[ \min_{\hat{b} \in B(x; \tau)} -\hat{b} (R^{-1})^2 \right] \quad (27)$$

where  $B$  denotes the set of valid  $\hat{b}$  induced by  $F$ .<sup>14</sup> That problem has a simple solution.

**Proposition 2** *The optimal information policy under the preferences (2) (and (27)) is*

$$\tau^*(\omega_j) = \underbrace{\theta^{-1/2}}_{\text{Shadow cost of info.}} \times \underbrace{\psi^{1/2}}_{\text{Ambiguity aversion}} \times \underbrace{Z(\omega_j)}_{\text{Utility weights}} \quad (28)$$

where  $\theta$  is the Lagrange multiplier on the information constraint,  $\sum_j \tau(\omega_j) \leq \bar{\tau}$ .

Recall that the function  $Z$  measures how the level of the log spectrum,  $f$ , affects utility. Agents optimally gather information exactly in proportion to  $Z$ , learning the most about the frequencies that are most important for utility. In terms of the adversarial game with nature, the agent chooses precision to constrain nature most at the frequencies that are potentially most painful. Since  $\tau$

<sup>14</sup>Technically,  $B(x; \tau) = \left\{ \hat{b} : PL \left( x \mid \log \left| \hat{b}(e^{i\omega}) \right|^2 \right) \geq \bar{L} \right\}$  subject to the additional restriction that  $\hat{b}$  be a Wold representation of some process.

also controls the potential complexity of  $f^w$ , the agent's choice of  $\tau^*$  implies that models are most complex where  $Z$  is highest – very low frequencies.

While  $Z$  controls the shape of  $\tau^*$ ,  $\theta$  and  $\psi$  determine its scale. An increase in the available precision  $\bar{\tau}$  lowers the Lagrange multiplier  $\theta$ , leading to more precision at all frequencies.  $\psi$  determines the extent to which nature is constrained by the penalized likelihood, i.e. how ambiguity-averse people are. Holding the shadow cost  $\theta$  of precision constant, a decrease in ambiguity-aversion through  $\psi$  lowers the chosen precisions.

To see the implication of proposition 2 for noise in the signals at each frequency, the right-hand panel of figure 1 plots  $Z(\omega)^{-1} \propto \tau^*(\omega)^{-1}$ . The variance of the signals that the agents receive is a simple function of frequency, rising smoothly as the frequency increases (it is an affine function of  $\cos(\omega)$ ).

Lemma 2 and propositions 1 and 2 give the complete analytic solution to the model. The remainder of the paper analyzes the implications of the solution for the types of models that agents optimally use and how those choices affect observable consumption behavior.

### 3 Behavior of the model agents use

We have two relevant cases for  $\tau$ . The utility-optimal information policy,  $\tau^*(\omega)$ , says that it is proportional to  $Z(\omega)$ , while the statistical benchmark is to set  $\tau(\omega)$  to equal a constant. We focus on two key results for  $f^w$  under those policies:

1. **Optimal learning eliminates excessive extrapolation:** Without an optimal information ( $\tau$ ) policy, the worst-case model displays excessive persistence compared to the truth – people over-extrapolate shocks. But under optimal information ( $\tau^*$ ), that bias disappears.
2. **Agents make mistakes primarily about the transitory component of income:** Under the optimal policy, agents use models that tend to deviate from the truth more at high than at low frequencies.

This section derives those results theoretically and examines them in numerical simulations of the model.

#### 3.1 Optimal learning eliminates excessive extrapolation

Taking an expansion around an infinite level of precision, the appendix derives the following first-order approximation in the continuous limit of the problem ( $d\omega \rightarrow 0$ ) for arbitrary  $\tau$ :

$$E[f^w(\omega; x, \tau) - f | f] \approx \psi\tau(\omega)^{-1} Z(\omega) + \lambda\tau(\omega)^{-1} f''(\omega). \quad (29)$$

Equation (29) yields our first important result. In the statistical benchmark case where  $\tau$  is constant across frequencies,  $f^w$  is biased in the direction of  $Z(\omega)$ . Recall from figure 1 that  $Z$

is large at low frequencies and close to zero elsewhere. So under the statistical benchmark, the worst-case model has excessively high power at low frequencies, which means that it implies is more persistent than the truth ( $f$ ). That result is almost exactly what is obtained in Bidder and Dew-Becker (2016), and is closely related to results in Hansen and Sargent (2010, 2016). Intuitively, since highly persistent models lead to the lowest utility, the agent naturally fears them.

Equation (29) also yields the more important part of the result, though, which is that under the optimal policy,  $\tau^*$ , there is no systematic bias towards either under- or over-extrapolation. Specifically, we have under the optimal policy

$$E[f^w(\omega; x, \tau^*) - f | f] \approx \psi^{1/2} \theta^{1/2} + \lambda \tau^*(\omega)^{-1} f''(\omega). \quad (30)$$

Since  $\tau^*(\omega) \propto Z(\omega)$ , the frequencies that are most important for utility are also the ones that the agent learns the most about, thus constraining the worst-case model. The proportionality completely cancels  $Z$  out of the bias, leaving just a constant.

When  $f^w$  deviates from  $f$  by only a constant, the two models have identical autocorrelations and differ only in the conditional variances. For example (ignoring the effects of  $f''$  for the moment; i.e. for small  $\lambda$ ), if income follows an AR(1) process with persistence  $\rho$ , then  $E[f^w]$  is the log spectrum for an AR(1) also with persistence  $\rho$ , but with innovations that have a greater variance.

Equation (30) is a key result of the paper. It shows that endogenous learning can completely eliminate overextrapolation. Intuitively, ambiguity averse agents tend to focus on models with excessive persistence because they are associated with low utility. But that fact also causes them to obtain the most information about those frequencies, thus entirely canceling out the effect of ambiguity.

This result stands in conflict with recent work that argues that ambiguity aversion and information processing constraints lead to overextrapolation. What we show here is that when people are able to choose what aspects of income to learn about, they naturally focus on the low frequencies, since those are most important for utility. But it is precisely that focus that then eliminates any bias towards excessive extrapolation.

### 3.1.1 Numerical example

To make the result above more concrete, we consider a simple numerical example. Suppose income is truly i.i.d. over time,  $Y_t = \varepsilon_t$ , so that the true model has zero persistence. Since  $f''(\omega) = 0$ , the second term in equations (29) and (30) is equal to zero. The left-hand panel of figure 5 plots the true (flat) log spectrum  $f(\omega)$  along with the mean worst-case spectra under the optimal information policy  $\tau^*$  and for the statistical benchmark in which  $\tau$  is constant across frequencies (the calibration is set so that they have equal total precision:  $\sum_j \tau^*(\omega_j) = \sum_j \tau$ ), which we denote with  $\bar{f}_*^w$  and  $\bar{f}_F^w$ , respectively. The figure shows that  $\bar{f}_*^w$  is shifted up by a constant compared to  $f$ , while  $\bar{f}_F^w$  actually has a significantly different shape, with a peak at low frequencies indicating persistence in income.



The right-hand panel of figure 5 plots the impulse response functions (the  $b$ 's) associated with the three models. Since income is truly i.i.d.,  $b_j = 0$  for  $j \geq 1$  under the true model. Under the optimal information policy with model uncertainty, the only thing that changes on average is that  $b_0$  becomes larger – people fear a higher variance, but they do not on average act as though income actually has any persistence. Under the statistical benchmark, though, there is clearly persistence in income: the impulse response is consistently positive after the initial impact. Figure 5 thus illustrates our first basic result. While ambiguity aversion and model uncertainty can often drive agents to act as though income is excessively persistent, that result is delicate: it disappears when people can allocate attention and information acquisition optimally.

### 3.2 Agents make mistakes about the transitory component of income

The primary mistakes in the agent's worst-case model come from the term involving  $f''(\omega)$ . That part of the formula is driven by the agent's smoothness prior. In the face of noisy data, agents estimate the spectrum of income by smoothing information across frequencies. Since  $f^w(\omega; x, \tau)$  is a convex combination of the data  $x$  local to  $\omega$ , it is biased upward when  $f'' > 0$  and downward when  $f'' < 0$ . Intuitively, if there is a narrow peak in  $f$ , a simple model will tend to smooth the peak out, and thus be biased downward.

In that sense, the agents also have a bias towards simplicity: they use models with smaller variations across frequencies when they have less information.<sup>15</sup> When the true spectrum is in fact complex, in the sense that it has local peaks and troughs, the worst-case model will tend to make mistakes in smoothing those peaks out. So the errors appear exactly where  $f''$  is large.

Since the optimal information policy gives the agents noisier signals about the spectrum at high frequencies, that is also where they make the largest smoothing errors. In (30),  $f''(\omega)$  is multiplied by  $\tau^*(\omega)^{-1}$ . So when precision is high, the term is scaled down and the worst-case spectrum tracks the true spectrum closely. But when  $\tau^*$  is small – at high frequencies – agents do more smoothing across frequencies and make larger mistakes.

#### 3.2.1 Numerical example

To illustrate the errors caused by smoothing, we now consider a richer numerical example with multiple peaks in the spectrum that we view as more realistic. The left-hand and middle panels of figure 5 plot the log spectrum of the data-generating process for income, while the right-hand panel plots the impulse response of income to a shock,  $\varepsilon$ . The calibration is chosen to have both high-and low-frequency components. The high-frequency piece – which generates the middle peak in the spectrum – is driven by the fact that a component of the shocks to income reverts: when income rises higher by \$1 today, it is lower on average by 50 cents over the next three periods. That behavior can be caused by forces that shift income over time but have little effect on total lifetime

<sup>15</sup>That intuition can be formalized. It is possible to show that correlations in the estimated spectrum,  $f^w(\omega; x, \tau)$ , are higher across frequencies, implying that complexity is lower, in regions where  $\tau$  is smaller.

income. For example, many people overpay taxes during the year and then receive refunds (e.g. Souleles (1999)). The low-frequency component of income – the left-hand peak in the spectrum – comes from the fact that the impulse response is persistently positive in the later periods following a shock. This represents a persistent component in income growth, and could come from variation over time in the average growth rate of the economy or the performance of one’s employer.<sup>16</sup>

We examine two specifications for  $\tau$ : the first is the optimum derived above,  $\tau^*$ , which is proportional to  $Z(\omega)$ ; the second specification is the statistical benchmark that sets  $\tau(\omega)$  to be constant at the mean of  $\tau^*$ :

$$\tau^F(\omega_j) = \tau^F \equiv n^{-1} \sum_{i=1}^n \tau^*(\omega_j). \quad (31)$$

As in the previous example, the choice of the mean for  $\tau^F$  implies that it has the exact same information cost as  $\tau^*$ . Note, though, that since precision is the inverse of variance, the average variance of the errors across frequencies is in fact much smaller under  $\tau^F$  than under  $\tau^*$ .

Figure 5 plots  $\bar{f}_*^w$  and  $\bar{f}_F^w$  for the two-peak calibration. The two log spectra are rather different from each other and the true model.  $\bar{f}_*^w$  matches  $f$  very well at the lowest frequencies, but it does a poor job of matching the middle-frequency peak in  $\bar{f}$  and also deviates substantially at higher frequencies.  $\bar{f}_F^w$  has the opposite behavior: it matches the middle-frequency peak and high-frequency behavior well, and in fact matches  $f$  well at almost all frequencies, but it fits relatively poorly at low frequencies. That is exactly what the formulas predict: optimal learning,  $\tau^*$ , causes models to be relatively more accurate at low than high frequencies. Overall, though,  $\bar{f}_F^w$  has a much better fit than  $\bar{f}_*^w$ , with a root mean squared error that is 42 percent smaller, due to the fact that  $\bar{f}_F^w$  spreads information evenly across frequencies.

The right-hand panel of figure 5 plots the lag polynomials,  $b$ ,  $\bar{b}_*^w$ , and  $\bar{b}_F^w$ , associated with the log spectra  $f$ ,  $\bar{f}_*^w$ , and  $\bar{f}_F^w$ , respectively.  $\bar{b}_*^w$  fails to match the short-run mean-reversion in the income process, while the lag polynomial for the suboptimal information policy,  $\bar{b}_F^w$ , does not, as predicted by the analytic results. The figure shows that the greater smoothness of  $\bar{f}_*^w$  also translates into smoothness in the associated lag polynomial, and in particular errors in the transitory behavior of income. But at longer lags, the figure shows that the optimal policy performs better, giving a closer fit to the persistent component of the impulse response function. Since it is the long-run part that determines human wealth, and hence optimal consumption, it is optimal from an expected utility perspective for agents to use models that fit the persistent component at the cost of missing the transitory dynamics.

---

<sup>16</sup>Technically, the impulse response function for income is equal to  $[1, -0.15, -0.3, -0.15, 0, \dots]$  plus  $0.04 \exp(-0.05j)$ . It is then scaled so that the standard deviation of consumption growth is 1.56 percent (when initial consumption is equal to 1).

As discussed above,  $n$  is intended to be taken as large – it is only used to avoid infinities – so we set it to 4000.  $\beta = 0.975$  to represent an annual calibration, and  $R = \beta^{-1}$  for simplicity.  $\bar{\tau}$ ,  $\lambda$ , and  $\psi$  are chosen in order to ensure that the agents make non-trivial mistakes in modeling consumption and that the behavior is visibly different across the two policies for  $\tau$ .  $\psi = 10^{-4}$ ;  $\lambda = 0.00075$ ;  $\bar{\tau} = 405.83$ . The parameterization is meant to illustrate the qualitative behavior of the model rather than match specific quantitative data.

## 4 Implications for observable consumption behavior

We now explore the implications of the results in the previous section for the observable behavior of consumption.

### 4.1 Consumption function

The consumption function from (18) implies that consumption growth follows

$$\Delta C_t = (1 - R^{-1}) b^w (R^{-1}) \varepsilon_{t+1}^w + \frac{\alpha}{2} (1 - R^{-1})^2 b^w (R^{-1})^2 + \alpha^{-1} \log \beta R \quad (32)$$

$$\text{where } \varepsilon_{t+1}^w \equiv b^w (L)^{-1} Y_t, \quad (33)$$

$\Delta \equiv 1 - L$  is the first-difference operator, and  $b^w(L)$  is the Wold representation associated with the worst-case model  $f^w$ . In the case where agents use the true model, so that  $b^w = b$  (i.e. under complete information), the filtered shocks,  $\varepsilon^w$  are equal to the true shocks,  $\varepsilon$ , and consumption follows a random walk with innovations equal to the innovation in the annuity value of the NPV of future income,  $(1 - R^{-1}) b (R^{-1}) \varepsilon_{t+1}$ . When the agent uses a model that differs from the truth, though,  $\varepsilon_{t+1}^w$  is no longer an i.i.d. process and consumption growth is no longer uncorrelated over time. That is, the agent's estimated shocks,  $\varepsilon^w$ , are in general serially correlated, which leads to (suboptimal) serial correlation in consumption growth.

To better understand the implications of the worst-case model for the behavior of consumption growth, we can write the log spectrum of consumption growth as

$$f_{\Delta C}^w(\omega; x, \tau) = \log \left( (1 - R^{-1})^2 b^w (R^{-1}; x, \tau)^2 \right) + f(\omega) - f^w(\omega; x, \tau). \quad (34)$$

Just like the spectrum of income,  $f_{\Delta C}^w$  represents a variance decomposition, measuring what types of fluctuations drive the overall variance of consumption growth. When the agent knows the true model,  $f_{\Delta C}^w$  is perfectly flat, which means that consumption growth is uncorrelated over time and the level of consumption is a random walk. But in general the agent does not know the true model. For example, if the true spectral density has a peak at some frequency but the worst-case spectrum does not, then  $f_{\Delta C}^w$  will inherit the same peak through the term  $f(\omega) - f^w(\omega; x, \tau)$ . That is, features of the income spectrum that the agent "ignores" in the sense that they do not appear in  $f^w$  are passed through to the spectrum of consumption growth.

Using (34), we can immediately map the results in the previous subsections into the spectrum of consumption growth. Specifically, for general information policies and for the optimal policy, we have

$$E[f_{\Delta C}^w(\omega; x, \tau) | f] \approx E \log \left( (1 - R^{-1})^2 b^w (R^{-1}; x, \tau)^2 \right) - \psi_{\tau}(\omega)^{-1} Z(\omega) - \lambda_{\tau}(\omega)^{-1} f''(\omega) \quad (35)$$

$$E[f_{\Delta C}^w(\omega; x, \tau^*) | f] \approx E \log \left( (1 - R^{-1})^2 b^w (R^{-1}; x, \tau^*)^2 \right) - \psi^{1/2} \theta^{1/2} - \lambda_{\tau^*}(\omega)^{-1} f''(\omega). \quad (36)$$

Again, the information policies differ in two key ways. First, comparing the terms  $\psi\tau(\omega)^{-1}Z(\omega)$  and  $\psi^{1/2}\theta^{1/2}$ , there are no systematic deviations of consumption growth from white noise under the optimal information policy. Under other policies, though, since people overextrapolate income shocks, consumption is actually mean reverting in the long-run – there is a trough in  $f_{\Delta C}^w$ . Intuitively, overextrapolation causes people to consume more than they can afford (more than human wealth) following positive shocks. Eventually, then, they must reduce consumption, causing long-run mean reversion. So the observable prediction of the model is that we actually *should not* observe long-run mean reversion in consumption growth. By the same token, people should also not underreact to shocks (as under rational inattention), which would lead to long-run persistence in consumption growth.

The second class of mistakes is the smoothing errors due to the term  $\lambda\tau(\omega)^{-1}f''(\omega)$ . This term says essentially that variation in the spectrum of income that the agent is not aware of passes directly into consumption growth. When  $f''$  is negative, for example, there is a local peak in the spectrum of income, and the spectrum of consumption growth then is also relatively high. Again, these errors are scaled by the precision of signals. The model predicts that consumption should track income relatively more closely – in the sense that their autocorrelations or impulse-responses are the same – at high than low frequencies. That is, transitory variation in income, such as the shifts in income over time studied by Souleles (1999), is predicted to pass directly into consumption. We illustrate that behavior below in a numerical example.

Compared to the behavior under the standard setup with no model uncertainty, our model generates, through limited information, excessive sensitivity of consumption to high-frequency shocks to income. This result is not obtained, though, by appealing to some sort of irrationality; rather, it arises simply from people optimally choosing to focus their attention on low frequencies. Endogenous attention leads to our second difference from the literature, which is that unlike other recent work on model uncertainty (Fuster, Hebert, and Laibson (2012), Bidder and Dew-Becker (2016), and Hansen and Sargent (2016)), the model does *not* predict excessive extrapolation of shocks. The model predicts excess sensitivity to transitory variation in income, but in fact the *correct* sensitivity to the permanent component.

The model also has rather different predictions from rational inattention over state variables (as opposed to rational inattention over model specifications), which suggest that they could be tested against each other empirically. As discussed by Sims (2003), the most prominent prediction of rational inattention is delayed reaction to shocks, due to the fact that people observe the shocks imperfectly. If income rises permanently, Sims (2003) shows that in general people will take a number of periods to fully realize that such a shock has occurred, meaning that consumption responds slowly to permanent shocks to income. Here, on the other hand, agents respond rapidly to permanent shocks because it is precisely the low-frequency part of income that they understand best.

Sims (2003) and Luo (2008) show that rational inattention can also generate excess sensitivity of consumption to income shocks, but the effects are calibration-specific and may be quantitatively

small (e.g. see the simulations in Sims (2003)). Intuitively, excess sensitivity arises because agents are not able to distinguish permanent from transitory shocks. So to obtain high-frequency mistakes, the rational inattention model must also predict low-frequency mistakes. In our model, though, the prediction of optimal information acquisition is in fact that the same attention choice both induces high-frequency mistakes and eliminates low-frequency mistakes. Furthermore, we see in the next section that the high-frequency mistakes can be quantitatively large and realistic.<sup>17</sup>

## 4.2 Numerical example

We examine the behavior of consumption under the numerical simulation when income has both transitory and persistent components. Figure 4 plots the log spectra of consumption growth under the various models.  $\bar{f}_*^w$  provides a closer fit to the utility optimal consumption spectrum at all frequencies. On the other hand, the statistical information policy produces a spectrum that is flatter – and closer to white noise – across most frequencies, but it has a very large peak at the lowest frequencies. The key question, then, will be which type of deviation – low- or high-frequency – is more relevant for utility.

To see how the fitting errors affect the behavior of consumption growth in the time domain, the right-hand panel of figure 4 plots the impulse response of the level of consumption to a unit shock to  $\varepsilon_t$  (i.e. a true innovation, not a filtered one) under the three consumption rules along with the cumulative impulse response of income (multiplied by  $(1 - R^{-1})$ ). As we would expect, the response of consumption under the full-information rule is flat: the permanent income hypothesis holds, and the response of consumption is approximately equal to the cumulative increase in income. The line for consumption under the optimal information policy shows that it inherits some of the short-run mean-reversion in income, rising and falling in the first few periods. It does not include the persistent component in income, though – consumption immediately jumps to approximately its long-run level, but the fluctuates around that level excessively. So the consumption policy is “right” in the long-run, but it is excessively sensitive to transitory variation in income in the short-run.

The behavior of a person using the model  $\bar{f}_*^w$  is again notably different from one using  $\bar{f}_F^w$ . The latter model does a better job of eliminating high-frequency fluctuations in consumption, but at the cost of inheriting the low-frequency behavior of income. The initial response of consumption under  $\bar{f}_F^w$  is too small, and consumption slowly drifts upward over the 80 periods of the IRF plotted here, eventually overshooting. So the  $\tau^F$  policy, counter to what is observed empirically, eliminates the sensitivity of consumption to transitory fluctuations in income, but causes consumption growth to deviate from white noise at long horizons. This result argues that empirically,  $\tau^*$  is a better description of consumption behavior than a setting where agents do not choose information optimally,  $\tau^F$ .

Those results may also be observed in more standard time series regressions for consumption growth. Table 1 below reports the coefficients from simulated regressions of consumption growth on

---

<sup>17</sup>It is also worth noting that the models in Sims (2003) and Luo (2008) can only be solved under quadratic utility, whereas we are able to accommodate CARA preferences here.

the predictable and unpredictable components of income growth under the two information policies and also under the full-information optimum.<sup>18</sup>

Information policy	Predictable income	Unpredictable income
$\tau^*$	0.21	0.95
$\tau^F$	0.06	0.74
Full-info. optimum	0	0.92

Table 1. Coefficients from regressions of consumption growth on income growth

The coefficient from the regression of consumption growth on the predictable part of income is of the same order of magnitude as the coefficient on the unpredictable part under  $\tau^*$ . The model can thus replicate the empirical result that consumption responds strongly to predictable income changes. That value is broadly consistent with the results of Parker (1999) and Souleles (1999), who both find that consumption rises by approximately 0.5 percent following a 1-percent anticipated increase in income.

Under the statistical benchmark,  $\tau^F$ , on the other hand, that relationship is much weaker, with the response to predictable income being, at 0.06, smaller by a factor of nearly 4. It is precisely the fact that agents optimally (under  $\tau^*$ ) fail to learn about high-frequency features of the model that causes them to overreact to predictable parts of income. Furthermore, note that the response of consumption to true income shocks is far closer to the full-information optimum under  $\tau^*$  than under  $\tau^F$ . This again demonstrates that in many ways,  $\tau^*$  helps agents get long-run responses right.

An alternative way to examine the behavior of consumption in the time domain is to study its autocorrelations. The left-hand panel of figure 5 plots the autocorrelations of consumption growth under  $\tau^*$  and  $\tau^F$ . Obviously under the full-information optimum, the autocorrelations are zero. At short lags, the autocorrelations are higher under  $\tau^*$ . Subsequently, though, the autocorrelations are substantially lower – by nearly a factor of 10. The right-hand panel plots the first autocorrelation of consumption growth over different spans. For a horizon denoted by  $n$  on the x-axis, we plot  $\text{corr}\left(\sum_{j=0}^{n-1} \Delta C_{t+j}, \sum_{j=0}^{n-1} \Delta C_{t-n+j}\right)$ . So the figure represents how consumption growth is correlated over neighboring intervals of length  $n$ . Consistent with the left-hand panel, for short intervals the correlations are higher under  $\tau^*$  than  $\tau^F$ . As we claimed above, though, the figure shows that consumption growth over long periods is substantially less autocorrelated under  $\tau^*$  than  $\tau^F$ .

To summarize, this example confirms the analytic results above that the optimal information policy does a good job of generating consumption growth that is close to white noise in the long-run, but that it causes consumption to be excessively sensitive to variation in income in the short-run. It also shows that the model can generate the empirical result that consumption responds to predictable variation in income.

---

<sup>18</sup>Here we use the version of the model in which income is difference-stationary. As discussed above, the results go through identically in that case. The difference is simply that then consumption and income have volatilities that are of the same order of magnitude, as observed in the data.

### 4.3 Empirical evidence

Since the optimal information policy implies that people learn the most about low-frequency features of the income process, it says that deviations of consumption growth from white noise should be observed primarily at high frequencies. Specifically, if the agent’s model of income dynamics,  $f^w(\omega; x, \tau^*)$ , is flat at high frequencies, then any variation in the shape of the true spectrum passes directly into consumption. The shape of the spectrum of  $f_{\Delta C}^w(\omega; x, \tau^*)$  will typically be similar to that of  $f(\omega)$  at high frequencies as the model predicts that people use simple (flat) models there.

Another way to build intuition for that prediction of the model is to note that high-frequency shocks also have relatively small effects on the net present value of income compared to more persistent shocks (which is why the function  $Z$  is relatively small at high frequencies). So the model essentially predicts that people spend excessively out of relatively small high-frequency increases in income compared to the larger low-frequency shocks.

Those predictions of the model are consistent with recent empirical evidence. Parker (1999) and Souleles (1999) provide classic evidence on the response of consumption to predictable changes in income due to the tax code (the cap on social security taxes and tax refunds, respectively). The shocks studied in those papers essentially shift income over time, exactly as in our numerical example. The results above show that consumption in the model does in fact respond to such variation in income, and that it tracks predictable income variation strongly.

Kaplan and Violante (2014; see references therein) review extensive evidence on the effectiveness of fiscal stimulus payments, finding that people tend to spend approximately 25 percent of these transitory payments in the quarter that they are received, even though the standard frictionless model would imply that they should spend a fraction near the level of the real interest rate (i.e. less than 1 percent per quarter). Moreover, these responses occur even among people with high incomes, who are less likely to be liquidity constrained (see also Kueng (2016)).

Kaplan and Violante explain the empirical evidence by arguing that when people hold illiquid assets, their consumption is excessively sensitive to transitory shocks because the benefit of smoothing is smaller than the cost of adjusting the stock of illiquid assets (e.g. housing). The intuition behind our results is similar to theirs (and also that of Cochrane (1989)) in that our results are also driven by the relatively small welfare benefit of smoothing transitory shocks. We differ in emphasizing the cost of learning about high-frequency dynamics, as opposed to assuming that saving is costly. Kaplan and Violante (2016) note that their model is consistent with the finding of Hsieh (2003) that consumption seems to respond relatively more to small than to large income shocks. That intuition is consistent with our argument that it is most natural for people to learn about shocks that have large effects on human wealth.

While the key source of variation for Kaplan and Violante (2014) is the size of shocks to income, for us it is their duration. Consumption mistakes should appear in response to short-duration shocks in our setting, and the empirical research finding violations of the permanent income hypothesis typically studies transitory income shocks.

Cochrane and Sbordone (1988) examine the joint relationship between aggregate consumption and output at long horizons and find that consumption helps forecast future output growth, but output does not help forecast consumption (nor do lags of consumption itself), implying that consumption growth is approximately white noise at long horizons. In other words, our model is consistent with the view that consumption growth may deviate from white noise and respond excessively to income in the short-term, but at longer horizons it is well described as white noise.

That implication requires aggregation, though, which is a nontrivial step. Since the consumption function in our model is linear, it will have desirable aggregation properties, but the exact details will depend on how income is driven by aggregate and idiosyncratic shocks at each frequency. Aggregate empirical results are thus not an ideal test of the model. The most direct test would be to measure the extent to which individual consumption growth is close to white noise over long horizons.

An alternative way to test the model, instead of examining consumption, would be to directly ask people what they are willing to pay for information. If they are at the optimum  $\tau^*$ , then information is equally valuable at all frequencies. On the other hand, under the standard models of ambiguity aversion without endogenous information acquisition, people would value low-frequency information most highly and be willing to pay the most for it.

#### 4.4 Relationship with the full-information optimal consumption rule

Our information-constrained agent uses a consumption rule that is suboptimal to the extent that  $b^w(L)$  differs from  $b(L)$ .  $b^w$  is not chosen to directly generate a path for consumption that necessarily maximizes realized utility; rather, ambiguity aversion causes it to be chosen to maximize utility under a pessimistic probability measure. We now show, though, that the agent's worst-case optimization problem is closely related to an optimization that approximates the correct consumption rule.

**Remark 1** *A second-order expansion of the Kullback–Leibler (KL) divergence between the full-information rational expectations consumption process and that used by an agent with model  $f^w$  around the point  $f^w = f$  is*

$$KL(f_{\Delta C}; f_{\Delta C}^w) \approx \frac{1}{4\pi} \int_0^{2\pi} \left( (Z(\omega) - 1)^2 + 2 \left( R^{-1} \frac{\alpha}{2} (1 - R^{-1}) \right)^2 Z(\omega)^2 \right) (f^w(\omega) - f(\omega))^2 d\omega. \quad (37)$$

The KL divergence is a likelihood-based measure of the deviation between the two random processes (one interpretation is that it measures how likely one would be to reject the hypothesis that consumption is driven by one process after observing data generated by the other). Squared errors in the model  $f^w$  are weighted by a quadratic function of  $Z(\omega)$ . As long as  $R$  is close enough to 1, this weighting function is strictly maximized at  $\omega = 0$ , meaning that reducing the distance between  $f^w$  and  $f$  at low frequencies reduces the KL divergence the most. The optimization



problem that our agent solves involves minimizing squared errors in  $f^w(\omega)$  weighted by  $\tau^*(\omega)$ , and proposition 2 shows that  $\tau^*(\omega) \propto Z(\omega)$ . The estimations of the agent and of someone minimizing KL divergence both involve using the weights given by  $Z$  to put more emphasis on the precision of the estimate at low frequencies.

#### 4.5 Extension: frequency-dependent information costs

In the baseline model, assumption 4 implies that agents have equal ability to learn about all frequencies. That assumption is most natural in the limited attention interpretation of the model, and it can also be supported when agents can always income sufficiently long to measure any frequency. A natural question, though, is how our results are changed when the cost of acquiring information varies across frequencies.

In this subsection, we consider the following alternative to the constraint in assumption 4:

$$\sum_{j=1}^n \phi(j) \tau(\omega_j) d\omega \leq \bar{\tau} \quad (38)$$

for some cost function  $\phi$ . It does not appear possible to obtain a closed-form solution for optimal attention,  $\tau^*$ , under general  $\phi$ . However, in the special case where there is no smoothing across frequencies –  $\lambda = 0$  – there is an analytic solution:

$$\text{for } \lambda = 0, \tau^*(\omega_j) = Z(\omega_j) \phi(\omega_j)^{-1/2} \theta^{-1/2} \psi^{1/2}. \quad (39)$$

(39) is a simple generalization of the result in the baseline case; the only difference is that now  $\tau^*(\omega_j)$  is decreasing in the cost of obtaining information at frequency  $\omega_j$ . If low frequencies are more expensive to learn about than higher frequencies, then  $\tau^*$  will have a less extreme tilt toward low frequencies than in the baseline case.

Recall our motivation for the learning framework in which agents get information about the dynamics of income by examining income histories of other people. In order to have information about a particular frequency  $\omega$ , an income history must have at least  $2\pi/\omega$  periods (that is the first periodogram ordinate; intuitively, one does not have any direct information about fluctuations that last longer than the data sample). So if only some fraction  $F(\omega)$  of people have been alive for at least  $2\pi/\omega$  periods, then on average an agent must look at  $1/F(\omega)$  histories in order to find one that can inform them about frequency  $\omega$ .

More specifically, suppose people die with a probability  $\delta \in (0, 1)$  in every period. Then as long as the birth rate is constant, the fraction of people who have been alive for at least  $k$  periods is  $\delta^{k-1}$ , implying that  $F(\omega) = \delta^{(2\pi/\omega)-1}$ . A reasonable functional form for  $\phi$  is therefore

$$\phi(\omega) = \delta^{1-(2\pi/\omega)}. \quad (40)$$

As  $\omega \rightarrow 0$ ,  $\phi(\omega) \rightarrow \infty$ , which means that in general this cost function will cause agents to learn less

about low frequencies than in the baseline. However,  $\phi'(0) = -\infty$ , while  $Z'(0) = 0$ . So attention should be increasing with frequency local to zero.

We calibrate  $\delta = 0.975$ , corresponding to an annual death probability of 2 percent, which we motivate as equivalent to people having a 50-year working life on average. The top panels of figure 5 then plot the optimal information policies  $\tau^* \propto Z$  and  $\tau^\phi \propto Z\phi^{-1/2}$  (normalized to have equal integrals). Both lines again peak at low frequencies, but whereas  $\tau^*$  peaks at frequency zero,  $\tau^\phi$  peaks at a slightly interior frequency. That peak comes at a frequency corresponding to cycles lasting approximately 160 years, though. So while the function  $\phi$  causes agents to learn less about the very lowest frequencies, they still very much focus their attention on long-term cycles.

To see how that change affects our calibration, the bottom panels of figure 5 plot the worst-case spectra under various  $\tau$  policies now also including  $\tau^\phi(\omega)$ .<sup>19</sup> That policy leads to results between the benchmark  $\tau^*$  and the constant  $\tau$  policy. At the very lowest frequencies, the  $\tau^\phi$  model does not match the true spectrum as well as  $\tau^*$ , but it still does much better than  $\tau^F$ . At the middle frequency peak and at higher frequencies, on the other hand, the policy  $\tau^\phi$  does a better job of matching the log spectrum than  $\tau^*$  but still worse than  $\tau^F$ .

This section therefore shows, as one might expect, that when low frequencies are more costly to learn about, the main results are weakened somewhat. We continue to find that agents allocate the most attention to low frequencies, just not to the *very* lowest – the peak is at an interior frequency, but one corresponding to cycles lasting a century or more. The impact on the worst-case model is to put it somewhere between that induced in the baseline optimum and that induced by the policy that puts equal weight on all frequencies.

## 5 Conclusion

This paper studies how people can direct their attention to different features of a model. We consider a nonparametric class of income processes and show precisely how agents optimally allocate attention to the behavior of income at different frequencies. The utility maximizing policy is to pay the most attention to the behavior of income at very low frequencies, and use a relatively simple and inaccurate model at high frequencies.

While there is extensive past work on learning, the innovation of this paper is to provide an exactly solvable framework for studying how learning can be applied to different aspects of a model of the world, as opposed to learning about state variables. The theory can be used to describe what people pay attention to, what aspects of the world they try to model accurately and what they use coarser approximations for, and the set of mistakes that people should be expected to make.

We show that optimal learning implies people are most likely to make mistakes at high frequencies, as those are the aspects of the income process least important for utility. Consistent with

---

<sup>19</sup>For non-zero  $\lambda$  with  $\phi$  varying across frequencies,  $\tau^\phi(\omega) \propto Z(\omega)\phi(\omega)^{-1/2}$  is not technically the optimal policy – it must be solved for numerically. We focus on the analytic case for the sake of simplicity. Furthermore, the calibration in figure 5 is set up so that the total precision under  $\tau^*$  is the same as that under  $\tau^\phi$  – they differ only in how that precision is allocated across frequencies.

empirical evidence, the model implies that consumption tends to track transitory fluctuations in income in the short-run, but at lower frequencies consumption growth is close to white noise (which it would be under the full-information optimal policy). In other words, the consumption mistakes that the empirical literature has documented are consistent with optimal learning.

## References

- Abel, Andrew B., Janice C. Eberly, and Stavros Panageas**, “Optimal Inattention to the Stock Market,” *The American Economic Review*, 2007, *97* (2), 244–249.
- , —, and —, “Optimal Inattention to the Stock Market with Information Costs and Transactions Costs,” *Econometrica*, 2013, *81* (4), 1455–1481.
- Akaike, Hirotugu**, “Smoothness Priors and the Distributed Lag Estimator.,” Technical Report, DTIC Document 1979.
- Bansal, Ravi and Ivan Shaliastovich**, “Confidence Risk and Asset Prices,” *The American Economic Review*, 2010, *100* (2), 537–541.
- Barron, John M. and Jinlan Ni**, “Endogenous Asymmetric Information and International Equity Home Bias: The Effects of Portfolio Size and Information Costs,” *Journal of International Money and Finance*, 2008, *27* (4), 617–635.
- Bidder, Rhys and Ian Dew-Becker**, “Long-Run Risk is the Worst-Case Scenario,” *The American Economic Review*, September 2016, *106* (9), 2494–2527.
- Bordalo, Pedro, Nicola Gennaioli, and Andrei Shleifer**, “Diagnostic Expectations and Credit Cycles,” *Working Paper*, 2016.
- Brillinger, David R.**, *Time Series: Data Analysis and Theory*, McGraw Hill, 1981.
- Cochrane, John H.**, “The Sensitivity of Tests of the Intertemporal Allocation of Consumption to Near-Rational Alternatives,” *The American Economic Review*, 1989, *79* (3), 319–337.
- and **Argia Sbordone**, “Multivariate Estimates of the Permanent Components of GNP and Stock Prices,” *Journal of Economic Dynamics and Control*, 1988, *12*(2–3), 255–296.
- Collin-Dufresne, Pierre, Michael Johannes, and Lars A Lochstoer**, “Parameter Learning in General Equilibrium: The Asset Pricing Implications,” *The American Economic Review*, 2016, *106* (3), 664–698.
- Dew-Becker, Ian**, “How Risky is Consumption in the Long-run? Benchmark Estimates From a Robust Estimator,” *Review of Financial Studies*, 2017, *30* (2), 631–666.

- and **Stefano Giglio**, “Asset pricing in the frequency domain: theory and empirics,” *Review of Financial Studies*, 2016, *29* (8), 2029–2068.
- Eichenbaum, Martin**, “Comment on “Natural Expectations, Macroeconomic Dynamics, and Asset Pricing”,” in “NBER Macroeconomics Annual 2011, Volume 26,” University of Chicago Press, 2011, pp. 49–60.
- Friedman, Milton**, *A Theory of the Consumption Function*, Princeton University Press, 1957.
- Fuster, Andreas, Benjamin Hebert, and David Laibson**, “Natural Expectations, Macroeconomic Dynamics, and Asset Pricing,” *NBER Macroeconomics Annual*, 2011, *26* (1), 1–48.
- Gersch, Will and Genshiro Kitagawa**, “Smoothness Priors Transfer Function Estimation,” *Automatica*, 1989, *25* (4), 603–608.
- Hall, Robert E.**, “Stochastic Implications of the Life Cycle-Permanent Income Hypothesis: Theory and Evidence,” *Journal of Political Economy*, 1988, *86* (6), 971–987.
- Hansen, Lars Peter and Thomas J. Sargent**, “Formulating and Estimating Dynamic Linear Rational Expectations Models,” *Journal of Economic Dynamics and Control*, 1980, *2*, 7–46.
- and — , “A note on Wiener-Kolmogorov Prediction Formulas for Rational Expectations Models,” *Economics Letters*, 1981, *8* (3), 255–260.
- and — , “Fragile Beliefs and the Price of Uncertainty,” *Quantitative Economics*, 2010, *1*(1), 129–162.
- and — , “Sets of Models and Prices of Uncertainty,” 2015. Working paper.
- , — , and **Thomas D. Tallarini**, “Robust Permanent Income and Pricing,” *Review of Economic Studies*, 1999, *66* (4), 873–907.
- Hsieh, Chang-Tai**, “Do Consumers React to Anticipated Income Changes? Evidence from the Alaska Permanent Fund,” *The American Economic Review*, 2003, *93* (1), 397–405.
- Jappelli, Tullio and Luigi Pistaferri**, “The Consumption Response to Income Changes,” *Annual Review of Economics*, 2010, *2*, 479–506.
- Ju, Nengjiu and Jianjun Miao**, “Ambiguity, Learning, and Asset Returns,” *Econometrica*, 2012, *80*(2), 559–591.
- Kacperczyk, Marcin, Stijn van Nieuwerburgh, and Laura Veldkamp**, “A Rational Theory of Mutual Funds’ Attention Allocation,” *Econometrica*, 2016, *84* (2), 571–626.
- Kaplan, Greg and Giovanni L Violante**, “A Model of the Consumption Response to Fiscal Stimulus Payments,” *Econometrica*, 2014, *82* (4), 1199–1239.

- Kitagawa, Genshiro and Will Gersch**, “A Smoothness Priors–State Space Modeling of Time Series with Trend and Seasonality,” *Journal of the American Statistical Association*, 1984, 79 (386), 378–389.
- and — , “A smoothness priors long AR model method for spectral estimation,” *IEEE transactions on automatic control*, 1985, 30 (1), 57–65.
- and — , *Smoothness Priors Analysis of Time Series*, Springer Science & Business Media, 1996.
- Kueng, Lorenz**, “Explaining Consumption Excess Sensitivity with Near-Rationality: Evidence from Large Predetermined Payments,” 2016. Working paper.
- Luo, Yulei**, “Consumption Dynamics under Information Processing Constraints,” *Review of Economic Dynamics*, 2008, 11 (2), 366–385.
- and **Eric R. Young**, “Risk-Sensitive Consumption and Savings Under Rational Inattention,” *American Economic Journal: Macroeconomics*, 2010, 2 (4), 281–325.
- Parker, Jonathan A.**, “The Reaction of Household Consumption to Predictable Changes in Social Security Taxes,” *The American Economic Review*, 1999, 89 (4), 959–973.
- Peng, Lin and Wei Xiong**, “Investor Attention, Overconfidence and Category Learning,” *Journal of Financial Economics*, 2006, 80 (3), 563–602.
- Shiller, Robert J.**, “A Distributed Lag Estimator Derived from Smoothness Priors,” *Econometrica*, 1973, pp. 775–788.
- Sims, Christopher A.**, “Implications of Rational Inattention,” *Journal of Monetary Economics*, 2003, 50 (3), 665–690.
- Souleles, Nicholas S.**, “The Response of Household Consumption to Income Tax Refunds,” *The American Economic Review*, 1999, 89 (4), 947–958.
- van Nieuwerburgh, Stijn and Laura Veldkamp**, “Information acquisition and under-diversification,” *The Review of Economic Studies*, 2010, 77 (2), 779–805.
- Veldkamp, Laura L.**, “Information Markets and the Comovement of Asset Prices,” *Review of Economic Studies*, 2006, 73 (3), 823–845.
- , *Information Choice in Macroeconomics and Finance*, Princeton University Press, 2011.
- Wahba, Grace**, “Automatic Smoothing of the Log Periodogram,” *Journal of the American Statistical Association*, 1980, 75 (369), 122–132.

**Wang, Neng**, “Precautionary Saving and Partially Observed Income,” *Journal of Monetary Economics*, 2004, 51 (8), 1645–1681.

—, “Optimal consumption and asset allocation with unknown income growth,” *Journal of Monetary Economics*, 2009, 56 (4), 524–534.

Figure 1: Weighting function  $Z(\omega)$  and its multiplicative inverse

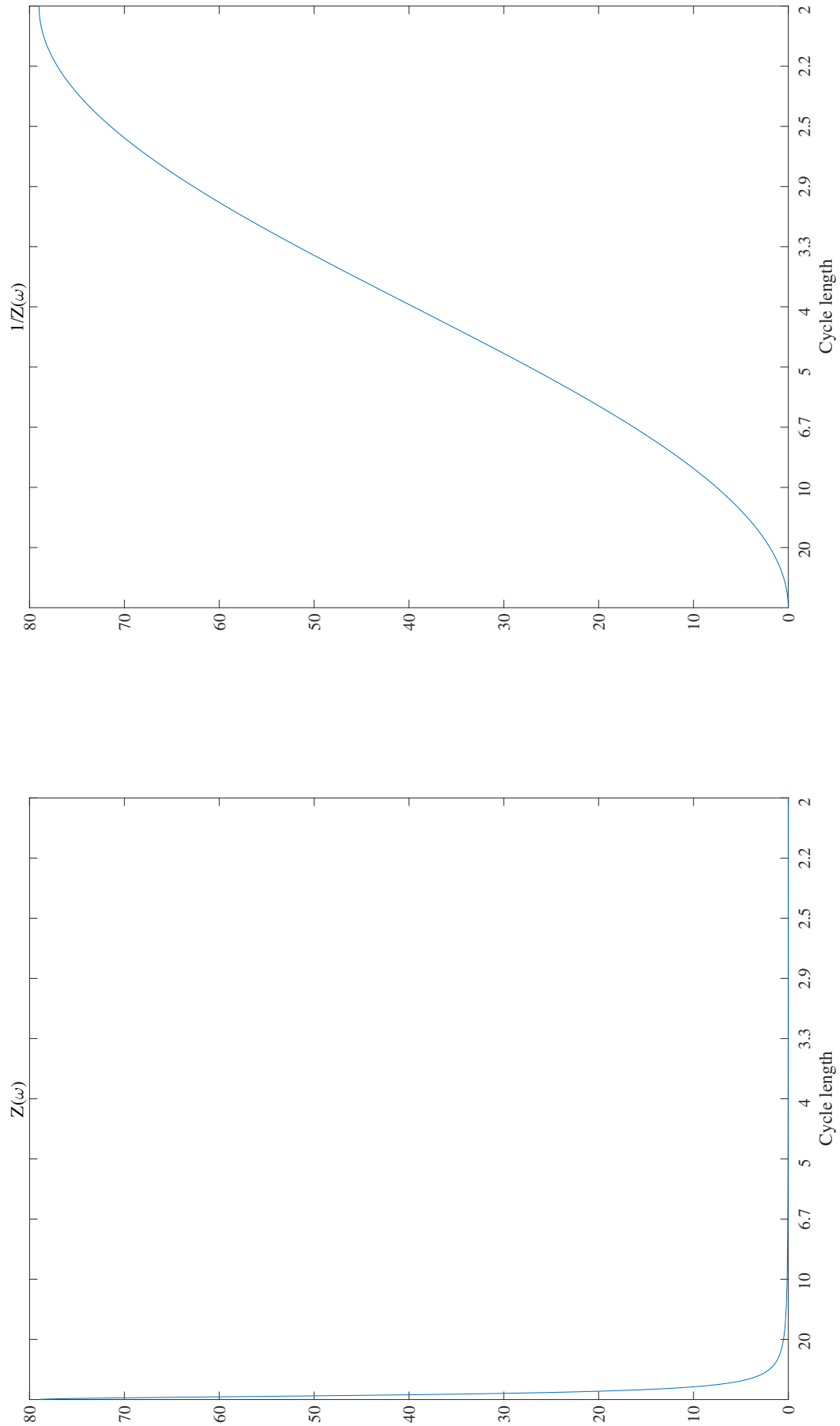
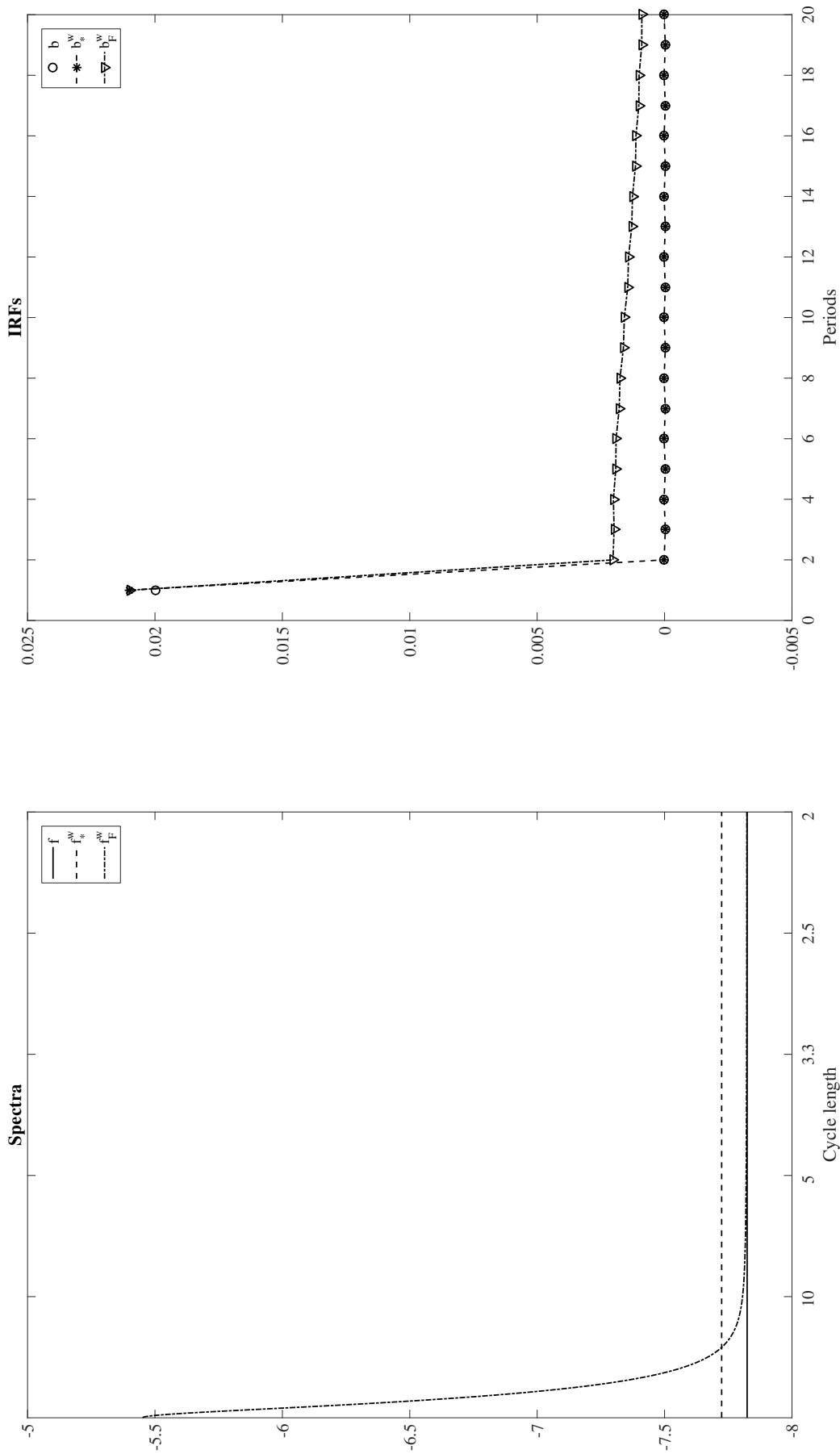


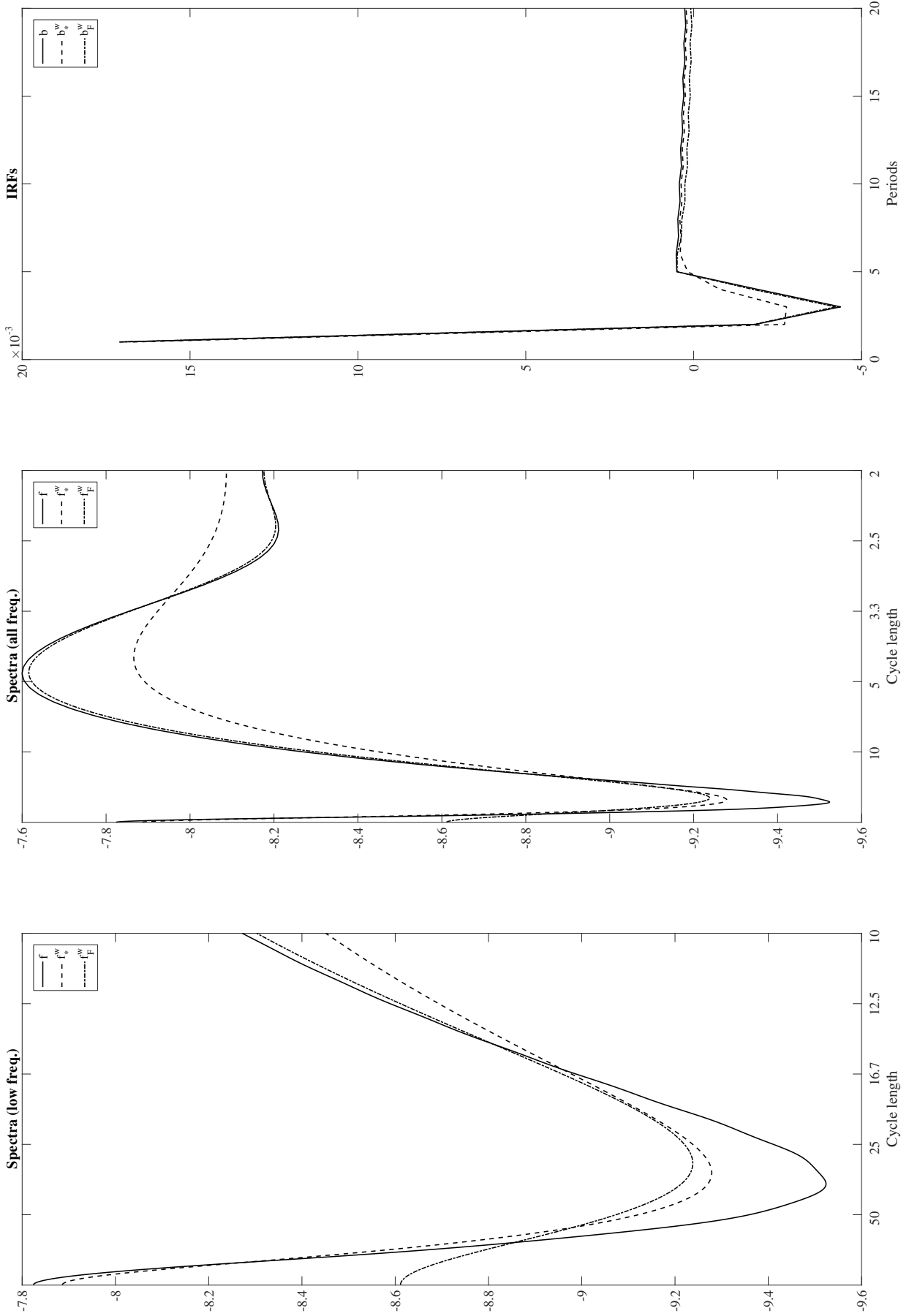
Figure 2: Average estimated log spectra and IRFs for white-noise income



Notes: In the left-hand panel,  $f$  is the true log spectrum of income (flat), the line for  $f_*^w$  is the average worst-case log spectrum under the optimal information policy, and the line for  $f_F^w$  is the average worse-case log spectrum under the statistical benchmark that yields equally precise signals at all frequencies. The right-hand panel plots the impulse response functions (Wold moving average representations) for income corresponding to the three log spectra.



Figure 3: Average estimated log spectra and IRFs with transitory and persistent components in income



Notes: The middle and left-hand panel correspond to the left-hand panel in figure 2, except for a different value for the true spectrum,  $f$ . The right-hand panel here corresponds to the right-hand panel in figure 2, but for this alternative example with an income process that has both persistent and transitory components.

Figure 4: Behavior of consumption with permanent and transitory income fluctuations

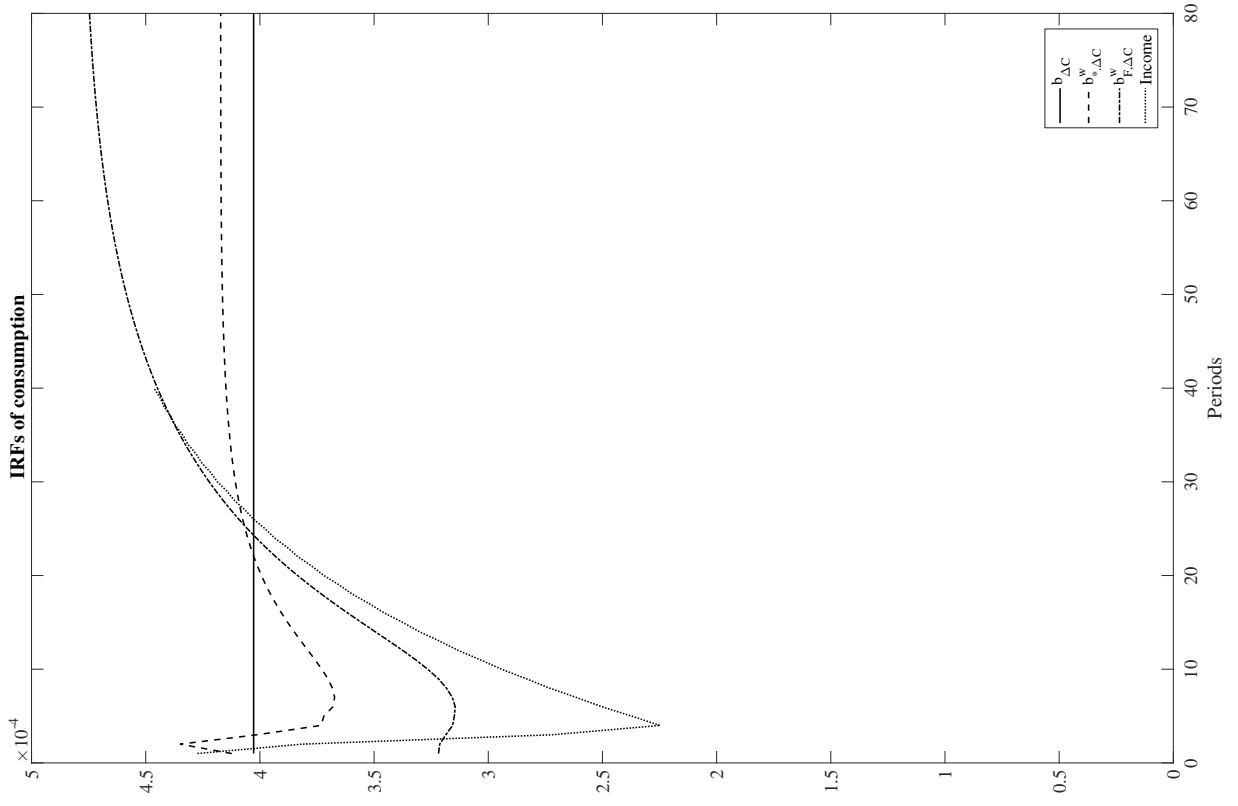
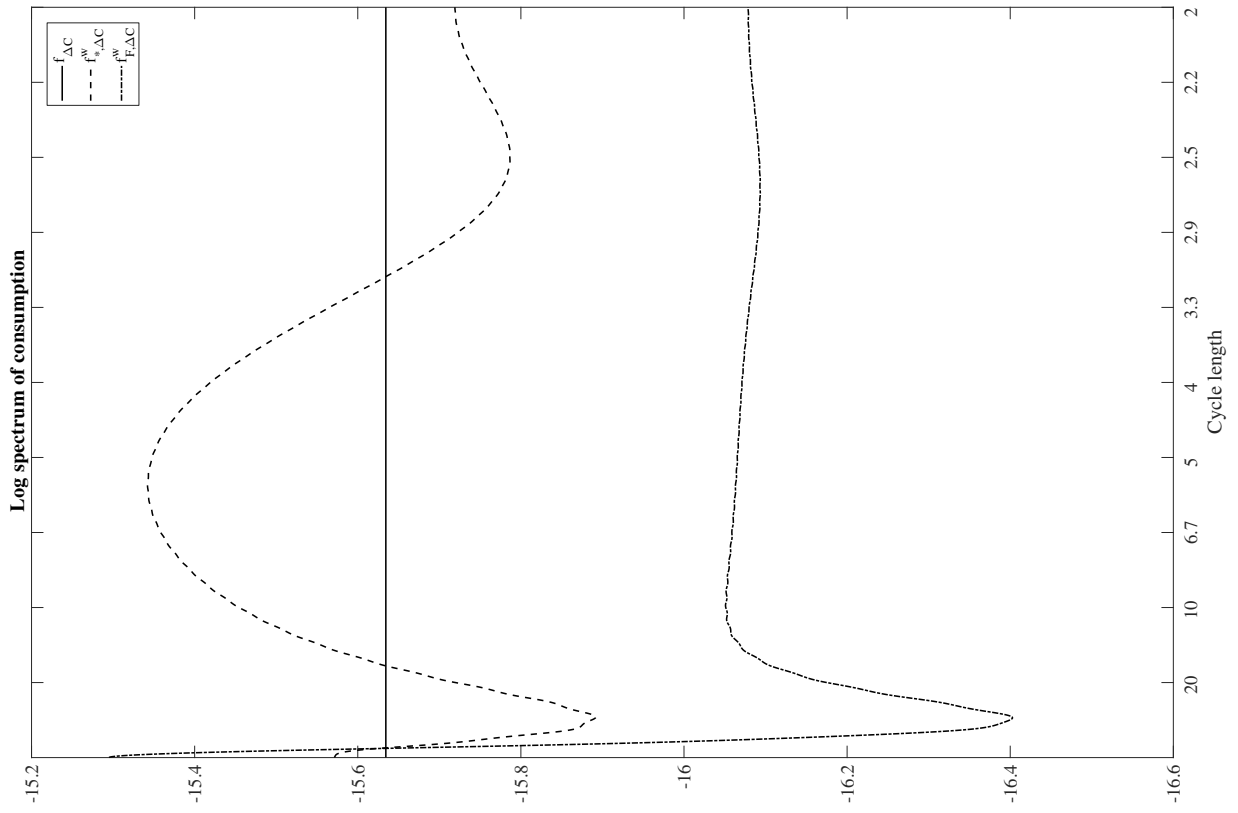
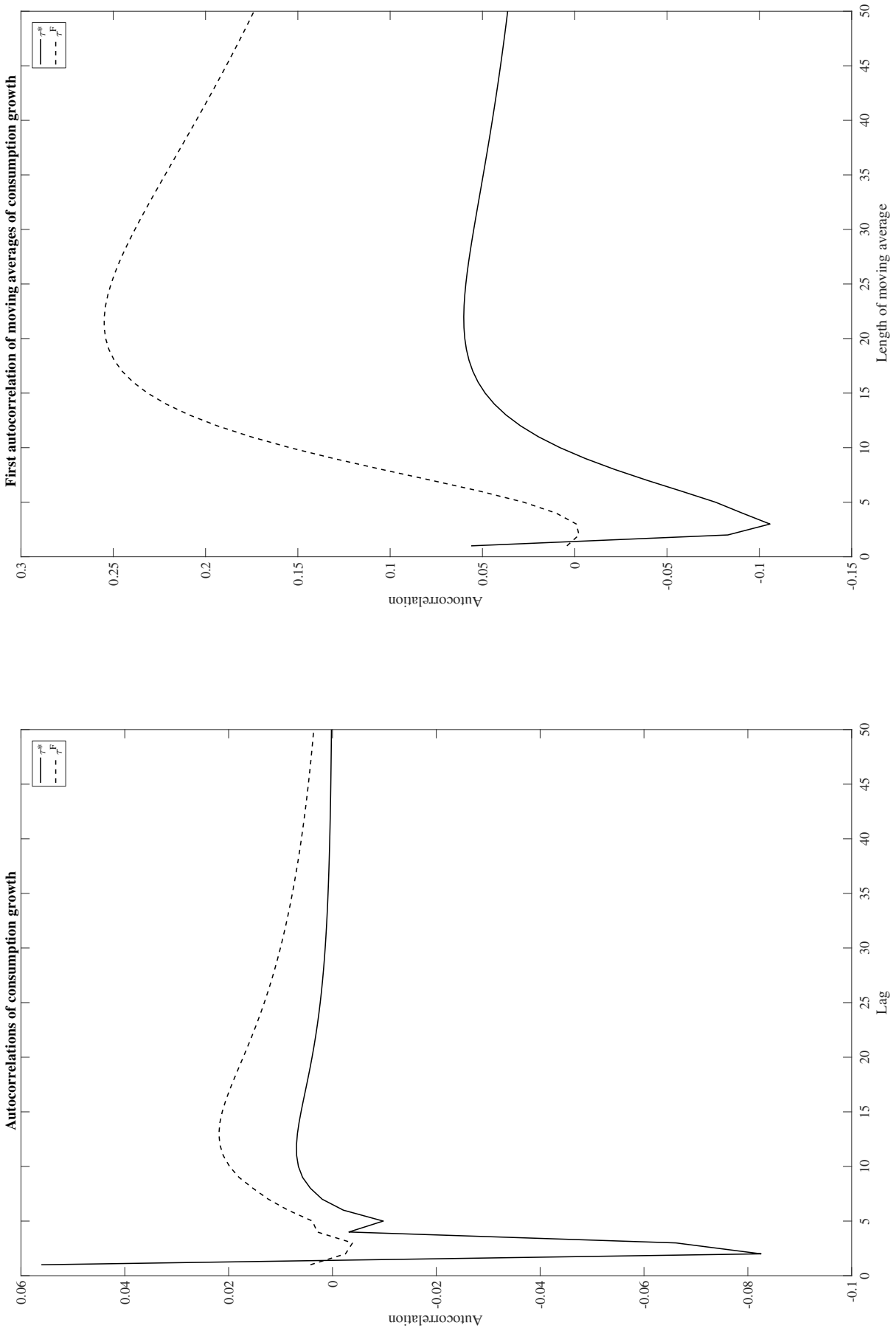
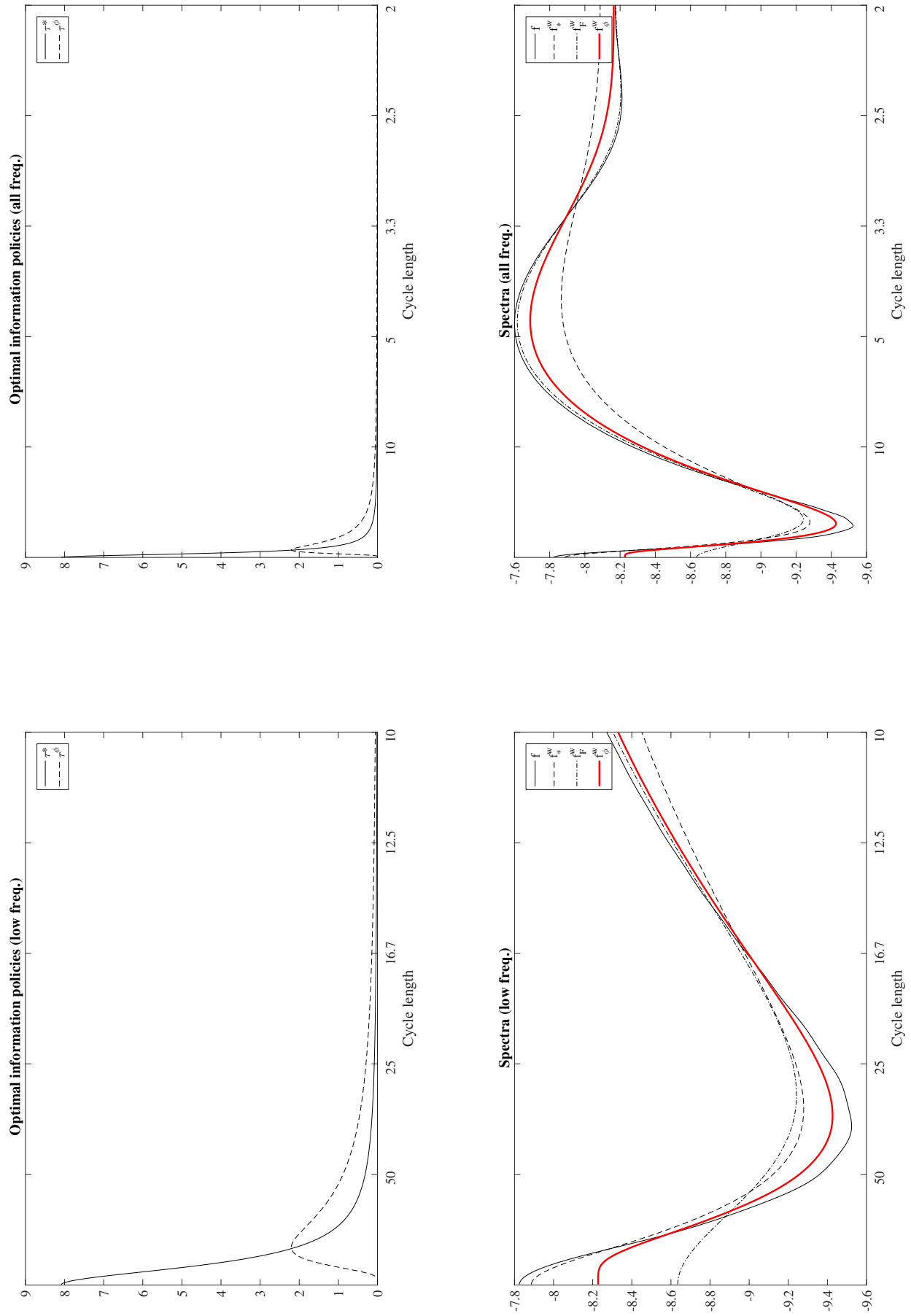


Figure 5: Persistence of consumption growth with transitory and persistent components in income



Notes: The left-hand panel plots the autocorrelations of consumption growth,  $\text{corr}(\Delta C_t, \Delta C_{t-1})$ . The right-hand panel plots the autocorrelation of moving averages,  $\text{corr}\left(\sum_{j=0}^{n-1} \Delta C_{t+j}, \sum_{j=0}^{n-1} \Delta C_{t-n+j}\right)$ , where  $n$  varies along the x-axis.

Figure 6: Effects of information cost varying across frequencies



Notes: The bottom panels of this figure replicate the left and middle panels in figure 3, but using the optimal information policy when information costs vary across frequencies, denoted  $\tau^\phi$ .  $f^\phi$  denotes the average worst-case log spectrum under that information policy.

## A Proof of lemma 2

From Dew-Becker (2016), the optimal consumption rule is

$$C_t = (R-1)W_{t-1} + -\frac{\alpha}{2}R^{-1}(1-R^{-1})\hat{b}(R^{-1})^2 + z(L)\hat{\varepsilon}_t - \alpha^{-1}\frac{\log \beta R}{R-1} \quad (1)$$

for a lag polynomial  $z(L)$ . Dew-Becker (2016) also shows that

$$E_t \left[ -\alpha^{-1} \sum_{j=0}^{\infty} \beta^j \exp(-\alpha C_{t+j}) \right] = \frac{-\alpha^{-1}}{1-R} \exp(-\alpha C_t) \quad (2)$$

(note that the probability measure for the expectation operator here is arbitrary) which implies

$$\begin{aligned} -\alpha^{-1} \log E_t \left[ (1-\beta) \sum_{j=0}^{\infty} \beta^j \exp(-\alpha C_{t+j}) \right] &= -\alpha^{-1} \log \frac{(1-\beta)}{1-R} + (R-1)W_{t-1} \\ &\quad -\frac{\alpha}{2}R^{-1}(1-R^{-1})\hat{b}(R^{-1})^2 + z(L)\hat{\varepsilon}_t - \alpha^{-1}\frac{\log \beta R}{R-1}. \end{aligned} \quad (3)$$

The result in the text then immediately follows.

## B Proof of lemma 3

$\hat{b}(L)$  is the Wold representation associated with the spectrum  $\exp \hat{f}(\omega)$  and is obtained through the canonical factorization of the spectrum (see Priestley (1981)). Define Fourier coefficients of  $\hat{f}$  as  $\hat{c}_k$ ,

$$\hat{c}_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} \cos(\omega k) \hat{f}(\omega) d\omega. \quad (4)$$

Given  $\hat{c}_k$ , the coefficients  $\hat{b}_k$  are constructed as

$$\hat{B}(\omega) = \exp \left( \hat{c}_0/2 + \sum_{j=1}^{\infty} \hat{c}_j e^{-i\omega j} \right) \quad (5)$$

$$\hat{b}_k = \int \hat{B}(\omega) \exp(i\omega k) d\omega \quad (6)$$

where from here on  $\int$  denotes  $\frac{1}{2\pi} \int_{-\pi}^{\pi}$ .

We can then obtain  $\hat{b}(R^{-1})$ ,

$$\hat{b}(R^{-1}) = \sum_{m=0}^{\infty} R^{-m} \int \exp \left( \hat{c}_0/2 + \sum_{k=1}^{\infty} \hat{c}_k e^{-i\omega k} \right) e^{i\omega m} d\omega. \quad (7)$$

Consider the derivative of  $\hat{b}(R^{-1})$  with respect to  $\hat{c}_k$  ( $k > 0$ ),

$$\frac{d\hat{b}(R^{-1})}{d\hat{c}_k} = \sum_{m=0}^{\infty} R^{-m} \int \exp \left( \hat{c}_0/2 + \sum_{k=1}^{\infty} \hat{c}_k e^{-i\omega k} \right) e^{i\omega(m-k)} d\omega \quad (8)$$

$$= \sum_{m=0}^{\infty} R^{-m} \hat{b}_{m-k} \quad (9)$$

$$= R^{-k} \hat{b}(R^{-1}), \quad (10)$$

where we use the fact that  $\hat{b}_{m-k} = 0$  for  $k > m$ . For  $k = 0$ ,

$$\frac{d\hat{b}(R^{-1})}{d\hat{c}_0} = \sum_{m=0}^{\infty} R^{-m} \int \exp\left(\hat{c}_0/2 + \sum_{k=1}^{\infty} \hat{c}_k e^{-i\omega k}\right) \frac{1}{2} d\omega \quad (11)$$

$$= \frac{1}{2} \hat{b}(R^{-1}). \quad (12)$$

We then have for  $k > 0$

$$\frac{d\hat{b}(R^{-1})^2}{d\hat{c}_k} = 2\hat{b}(R^{-1}) \times R^{-k} \hat{b}(R^{-1}) \quad (13)$$

$$\frac{d \log \hat{b}(R^{-1})^2}{d\hat{c}_k} = 2R^{-k} \quad (14)$$

and

$$\frac{d \log \hat{b}(R^{-1})^2}{d\hat{c}_0} = 1 \quad (15)$$

(where we square  $\hat{b}(R^{-1})$  so that we are taking the log of a positive number). In other words, then,  $\log \hat{b}(R^{-1})^2$  is linear in the  $\hat{c}_k$  and depends only on them and possibly a constant. We can therefore write

$$\log \hat{b}(R^{-1})^2 = cons + 2 \left( \hat{c}_0/2 + \sum_{k=1}^{\infty} R^{-k} \hat{c}_k \right) \quad (16)$$

$$= cons + \sum_{k=-\infty}^{\infty} R^{-|k|} \hat{c}_k \quad (17)$$

for some unknown constant *cons*. Using the definition of  $\hat{c}_k$ , we have

$$\sum_{k=-\infty}^{\infty} R^{-|k|} \hat{c}_k = \sum_{k=-\infty}^{\infty} R^{-|k|} \int \cos(\omega k) f(\omega) d\omega \quad (18)$$

$$= \int \sum_{k=-\infty}^{\infty} \cos(\omega k) R^{-|k|} f(\omega) d\omega \quad (19)$$

$$= \int Z(\omega) f(\omega) d\omega, \quad (20)$$

where

$$Z(\omega) \equiv \sum_{k=-\infty}^{\infty} \cos(\omega k) R^{-|k|} = 1 + 2 \sum_{k=1}^{\infty} R^{-k} \cos(\omega k). \quad (21)$$

Finally, note that if  $f(\omega) = 0$ , then the process is white noise with unit variance, so  $b(R^{-1})^2 = 1$ . That immediately implies that the constant term is zero, yielding the desired result,

$$\log \hat{b}(R^{-1})^2 = \int Z(\omega) \hat{f}(\omega) d\omega. \quad (22)$$

## C Finding the worst-case spectrum

Nature chooses  $\{\hat{f}(\omega_j)\}$  to solve

$$\begin{aligned} \{f^w(\omega_j)\} = \arg \max_{\{\hat{f}(\omega_j)\}} & \sum_{j=1}^n Z(\omega_j) \hat{f}(\omega_j) d\omega \\ & - \frac{\psi^{-1}}{2} \sum_{j=1}^n \left(x(\omega_j) - \hat{f}(\omega_j)\right)^2 \tau(\omega_j) d\omega - \frac{\psi^{-1}}{2} \lambda \sum_{j=2}^n \left(\frac{\hat{f}(\omega_j) - \hat{f}(\omega_{j-1})}{d\omega}\right)^2 d\omega. \end{aligned} \quad (23)$$

The first-order conditions for interior points ( $1 < j < n$ ) are

$$0 = Z(\omega_j) + \psi^{-1} (x(\omega_j) - f^w(\omega_j)) \tau(\omega_j) + \frac{\psi^{-1} \lambda}{d\omega} \left( \left( \frac{f^w(\omega_{j+1}) - f^w(\omega_j)}{d\omega} \right) - \left( \frac{f^w(\omega_j) - f^w(\omega_{j-1})}{d\omega} \right) \right).$$

At the boundaries they are

$$0 = Z(\omega_1) + \psi^{-1} (x(\omega_1) - f^w(\omega_1)) \tau(\omega_1) + \psi^{-1} \lambda \frac{f^w(\omega_2) - f^w(\omega_1)}{d\omega^2} \quad (24)$$

$$0 = Z(\omega_n) + \psi^{-1} (x(\omega_n) - f^w(\omega_n)) \tau(\omega_n) - \psi^{-1} \lambda \frac{f^w(\omega_n) - f^w(\omega_{n-1})}{d\omega^2}. \quad (25)$$

We define here vectors containing the various objects at the frequencies  $\omega_j$  using variables with no subscript. For example,  $\tau \equiv [\tau(\omega_1), \tau(\omega_2), \dots, \tau(\omega_n)]'$ . We can then write the first-order conditions as

$$0 = Z + \psi^{-1} \text{diag}(\tau) (x - f^w) + \psi^{-1} \lambda D f^w, \quad (26)$$

where  $\text{diag}(\tau)$  is a matrix with  $\tau$  on the diagonal and zero elsewhere and  $D$  is a differencing matrix:

$$D \equiv \begin{bmatrix} -1 & 1 & 0 & 0 & \dots & 0 \\ 1 & -2 & 1 & 0 & & \\ 0 & 1 & -2 & 1 & & \vdots \\ \vdots & & & \ddots & & 0 \\ & & 0 & 1 & -2 & 1 \\ 0 & \dots & 0 & 0 & 1 & -1 \end{bmatrix} d\omega^{-2}. \quad (27)$$

The second-order condition is that

$$-\text{diag}(\tau) + \lambda D \quad (28)$$

is negative definite, i.e. that all of its eigenvalues are negative.

The solution to nature's optimization problem is then obtained by directly solving (26):

$$f^w = (\text{diag}(\tau) - \lambda D)^{-1} (\psi Z + \text{diag}(\tau) x) \quad (29)$$

$$= (I - \lambda \text{diag}(\tau^{-1}) D)^{-1} (\psi \text{diag}(\tau^{-1}) Z + x), \quad (30)$$

where  $\tau^{-1}$  here is an elementwise inverse of the vector  $\tau$ . Since this is a linear problem, the solution is unique as long as the matrix inverse exists.

## D Proposition 1

Consider a total derivative of (26) with respect to  $\tau'$  at the point  $x = \bar{f}$ :

$$0 = \psi^{-1} \text{diag}(\bar{f} - f^w) - \psi^{-1} \text{diag}(\tau) \frac{df^w}{d\tau'} + \psi^{-1} \lambda D \frac{df^w}{d\tau'}. \quad (31)$$

We can then solve for  $\frac{df^w}{d\tau'}$ :

$$\frac{df^w}{d\tau'} = (\lambda D - \text{diag}(\tau'))^{-1} \text{diag}(f^w - \bar{f}) \quad (32)$$

Now the objective is to minimize

$$\{\tau^*(\omega_j)\} = \arg \min_{\{\tau(\omega_j)\}} \log b^w (R^{-1})^2 + \theta \sum_j \tau(\omega_j) d\omega \quad (33)$$

$$= \arg \min_{\{\tau(\omega_j)\}} Z' f^w d\omega + \theta \sum_j \tau(\omega_j) d\omega. \quad (34)$$

The first-order condition for that problem is

$$0 = Z' \frac{df^w}{d\tau'} + \theta \mathbf{1}_{1 \times n}, \quad (35)$$

where  $\mathbf{1}_{1 \times n}$  is a  $1 \times n$  vector of ones. Inserting the formula for  $\frac{df^w}{d\tau'}$  yields

$$0 = Z' (\lambda D - \text{diag}(\tau^{*'}))^{-1} \text{diag}(f^w - \bar{f}) + \theta \mathbf{1}_{1 \times n} \quad (36)$$

$$Z' = -\theta \mathbf{1}_{1 \times n} \text{diag}(f^w - \bar{f})^{-1} (\lambda D - \text{diag}(\tau^{*'})). \quad (37)$$

Now we conjecture that  $f^w - \bar{f}$  is equal to a constant  $c$  multiplied by a column of ones. We then have

$$Z' = -\theta c^{-1} \mathbf{1}_{1 \times n} (\lambda D - \text{diag}(\tau^{*'})) \quad (38)$$

$$= \theta c^{-1} \tau^{*'}, \quad (39)$$

where the second line uses the fact that  $\mathbf{1}_{1 \times n} D = \mathbf{0}_{1 \times n}$  since the columns of  $D$  sum to zero.

In order to confirm that result, we must now show that when  $Z = \theta c^{-1} \tau^*$ ,  $f^w - \bar{f} = c \mathbf{1}_{n \times 1}$ . Inserting  $Z = \theta c^{-1} \tau^*$  into (26) yields

$$0 = \theta c^{-1} \tau^* + \psi^{-1} \text{diag}(\tau^*) (\bar{f} - f^w) + \psi^{-1} \lambda D f^w. \quad (40)$$

In order for it to be the case that  $f^w - \bar{f} = c \mathbf{1}_{n \times 1}$ , we must have

$$0 = \theta c^{-1} \tau^* - \psi^{-1} \text{diag}(\tau^*) \mathbf{1}_{n \times 1} c + \psi^{-1} \lambda D \mathbf{1}_{n \times 1} (\bar{f} + c) \quad (41)$$

$$= \theta c^{-1} \tau^* - \psi^{-1} \tau^* c. \quad (42)$$

where the second line uses the fact that  $D \mathbf{1}_{n \times 1} = \mathbf{0}_{n \times 1}$ . This is solved by

$$\sqrt{\theta \psi} = c \quad (43)$$

$$Z = \theta c^{-1} \tau^* \quad (44)$$

$$\tau^* = (\theta/\psi)^{-1/2} Z. \quad (45)$$

We can then plug the value of  $\tau^*$  into the equation for  $f^w$ :

$$f^w = (I - \lambda \text{diag}(\tau^{*-1}) D)^{-1} \left( \psi \text{diag} \left( (\theta/\psi)^{1/2} Z^{-1} \right) Z + x \right) \quad (46)$$

$$= (I - \lambda \text{diag}(\tau^{*-1}) D)^{-1} \left( \mathbf{1}_{n \times 1} \theta^{1/2} \psi^{1/2} + x \right), \quad (47)$$

where, as with  $\tau^{-1}$ ,  $Z^{-1}$  is an elementwise inverse of the vector  $Z$ . It follows that

$$E[f^w] = (I - \lambda \text{diag}(\tau^{*-1}) D)^{-1} \left( \mathbf{1}_{n \times 1} \theta^{1/2} \psi^{1/2} + \bar{f} \right) \quad (48)$$

$$= \mathbf{1}_{n \times 1} \theta^{1/2} \psi^{1/2} + \bar{f}, \quad (49)$$



where the last line follows from the fact that the rows of  $(I - \lambda \text{diag}(\tau^{*-1}) D)^{-1}$  sum to 1. To see why, note that

$$(I - \lambda \text{diag}(\tau^{*-1}) D)^{-1} = I + \lambda \text{diag}(\tau^{*-1}) D + (\lambda \text{diag}(\tau^{*-1}) D)^2 + \dots \quad (50)$$

The rows of  $\lambda \text{diag}(\tau^{*-1}) D$  sum to zero, meaning that  $1_{n \times 1}$  is an eigenvector with eigenvalue zero. When a matrix is raised to a power, its eigenvectors are unchanged and its eigenvalues are raised to the same power, meaning that  $1_{n \times 1}$  remains an eigenvector with 0 the associated eigenvalue, and the rows sum to zero. Since the rows of  $I$  sum to 1, the rows of  $(I - \lambda \text{diag}(\tau^{*-1}) D)^{-1}$  then do also.

## D.1 Bias of $f^w(\omega; x, \tau)$

From above, the solution for the vector  $f^w$  is

$$f^w(x, \tau) = (I - \lambda \text{diag}(\tau^{-1}) D)^{-1} (\psi \text{diag}(\tau^{-1}) Z + x) \quad (51)$$

$$f^w(x, \tau) = \left( I + \sum_{j=1}^{\infty} (\lambda \text{diag}(\tau^{-1}) D)^j \right) (\psi \text{diag}(\tau^{-1}) Z + x) \quad (52)$$

$$f^w(x, \tau) - \psi \text{diag}(\tau^{-1}) Z - x = \left( \sum_{j=1}^{\infty} (\lambda \text{diag}(\tau^{-1}) D)^j \right) (\psi \text{diag}(\tau^{-1}) Z + x). \quad (53)$$

Now scale  $\tau^{-1}$  by  $c$  and divide both sides by  $c$

$$c^{-1} f^w(x, \tau/c) - \psi \text{diag}(\tau^{-1}) Z - c^{-1} x = \left( \begin{array}{c} \lambda \text{diag}(\tau^{-1}) D \\ + \sum_{j=2}^{\infty} c^{j-1} (\lambda \text{diag}(\tau^{-1}) D)^j \end{array} \right) (c \psi \text{diag}(\tau^{-1}) Z + x). \quad (54)$$

Since both sides are linear in  $x$ , we can take the expectation and then the limit as  $c \rightarrow 0$  to yield

$$\lim_{c \rightarrow 0} \frac{E[f^w(x, \tau/c)] - f}{c} = \psi \text{diag}(\tau^{-1}) Z + \lambda \text{diag}(\tau^{-1}) D f. \quad (55)$$

In the limit as  $n \rightarrow \infty$ ,  $Df$  becomes  $f''$ .

## E Consumption and income forecasts

### E.1 The behavior of consumption

From Dew-Becker (2016), consumption follows

$$C_t = (R - 1) W_{t-1} + Z_t - (R - 1)^{-1} \alpha^{-1} \log(\beta R) \quad (56)$$

$$W_t = W_{t-1} + Y_t - Z_t + (R - 1)^{-1} \alpha^{-1} \log(\beta R), \quad (57)$$

where

$$Z_t = (1 - R^{-1}) Y_t - \frac{1}{\alpha} R^{-1} \log E_t \exp(-\alpha Z_{t+1}). \quad (58)$$

We then have

$$C_{t+1} = (R - 1) W_t + Z_{t+1} - (R - 1)^{-1} \alpha^{-1} \log(\beta R) \quad (59)$$

$$\Delta C_{t+1} = (R - 1) W_t + Z_{t+1} - (R - 1) W_{t-1} - Z_t \quad (60)$$

$$\Delta C_{t+1} = (R - 1) (R W_{t-1} + Y_t - C_t) + Z_{t+1} - (R - 1) W_{t-1} - Z_t \quad (61)$$

$$\Delta C_{t+1} = (R - 1) Y_t + Z_{t+1} - R Z_t + \alpha^{-1} \log(\beta R). \quad (62)$$

Now define  $H$  as follows:

$$Z_t = (1 - R^{-1}) Y_t - \frac{1}{\alpha} R^{-1} \log E_t \exp(-\alpha Z_{t+1}) \quad (63)$$

$$H_t \equiv Z_t - (1 - R^{-1}) Y_t \quad (64)$$

$$H_t = -\frac{1}{\alpha} R^{-1} \log E_t \exp(-\alpha (H_{t+1} + (1 - R^{-1}) Y_{t+1})). \quad (65)$$

This definition yields

$$\Delta C_{t+1} = R((1 - R^{-1}) Y_t - Z_t) + Z_{t+1} + \alpha^{-1} \log(\beta R) \quad (66)$$

$$= H_{t+1} - R H_t + (1 - R^{-1}) Y_{t+1} + \alpha^{-1} \log(\beta R), \quad (67)$$

with the recursion

$$\bar{h} + h(L) \varepsilon_t = -\frac{1}{\alpha} R^{-1} \log E_t [\exp(-\alpha (\bar{h} + h(L) \varepsilon_{t+1} + (1 - R^{-1}) b(L) \varepsilon_{t+1})) | f^w] \quad (68)$$

$$= R^{-1} \left( \bar{h} + \sum_{j=1}^{\infty} (h_j + (1 - R^{-1}) b_j^w) \varepsilon_{t+1-j} \right) - R^{-1} \frac{\alpha}{2} (h_0 + (1 - R^{-1}) b_0^w)^2 \quad (69)$$

and solution

$$h_j = R^{-1} h_{j+1} + R^{-1} (1 - R^{-1}) b_{j+1}^w \quad (70)$$

$$\bar{h} = -\frac{R^{-1}}{1 - R^{-1}} \frac{\alpha}{2} (h_0 + (1 - R^{-1}) b_0^w)^2 \quad (71)$$

$$= -R^{-1} \frac{\alpha}{2} (1 - R^{-1}) b^w (R^{-1})^2 \quad (72)$$

$$\Delta C_{t+1} = (1 - R^{-1}) Y_{t+1} + H_{t+1} - R H_t + \alpha^{-1} \log \beta R. \quad (73)$$

Now we can insert the formulas for the various objects:

$$\Delta C_t = (1 - R^{-1}) b(L) \varepsilon_{t+1} + (1 - R) \bar{h} + h_0 \varepsilon_{t+1}^w + \sum_{j=0}^{\infty} (h_{j+1} - R h_j) \varepsilon_{t-j}^w + \alpha^{-1} \log \beta R \quad (74)$$

$$= (1 - R^{-1}) b(L) \varepsilon_{t+1} + (1 - R) \bar{h} + (1 - R^{-1}) (b^w (R^{-1}) - b_0^w) \varepsilon_{t+1}^w \quad (75)$$

$$- \sum_{j=0}^{\infty} (1 - R^{-1}) b_{j+1}^w \varepsilon_{t-j}^w + \alpha^{-1} \log \beta R \quad (76)$$

$$= (1 - R^{-1}) b(L) \varepsilon_{t+1} + (1 - R) \bar{h} + (1 - R^{-1}) (b^w (R^{-1}) - b_0^w) \frac{b(L)}{b^w(L)} \varepsilon_{t+1} \quad (77)$$

$$- \sum_{j=0}^{\infty} (1 - R^{-1}) b_{j+1}^w \frac{b(L)}{b^w(L)} \varepsilon_{t-j} + \alpha^{-1} \log \beta R \quad (78)$$

$$= (1 - R^{-1}) b(L) \varepsilon_{t+1} + (1 - R) \bar{h} + (1 - R^{-1}) b^w (R^{-1}) \frac{b(L)}{b^w(L)} \varepsilon_{t+1} \quad (79)$$

$$- (1 - R^{-1}) b^w(L) \frac{b(L)}{b^w(L)} \varepsilon_{t+1} + \alpha^{-1} \log \beta R \quad (80)$$

$$= (1 - R^{-1}) b^w (R^{-1}) \frac{b(L)}{b^w(L)} \varepsilon_{t+1} + (1 - R) \bar{h} + \alpha^{-1} \log \beta R. \quad (81)$$

So consumption growth is equal to a constant plus  $(1 - R^{-1}) b^w (R^{-1}) \frac{b(L)}{b^w(L)} \varepsilon_{t+1}$ . The dynamic behavior of consumption growth is therefore determined by  $b^w (R^{-1}) \frac{b(L)}{b^w(L)}$ . The spectral density of consumption growth is

$$f_{\Delta C}^w(\omega) = b^w (R^{-1})^2 \frac{f(\omega)}{f^w(\omega)}. \quad (82)$$

## F KL divergence for consumption process

We consider the relative entropy of consumption growth under the worst-case model compared to the true model. If we have two models of consumption growth defined by their spectra and means,  $\{f_{\Delta C}(\omega), \mu_{\Delta C}\}$ , then the Kullback–Leibler divergence is

$$\int \frac{\exp f_{\Delta C}^w(\omega)}{\exp f_{\Delta C}(\omega)} - \log \frac{\exp f_{\Delta C}^w(\omega)}{\exp f_{\Delta C}(\omega)} d\omega + \frac{(\mu_{\Delta C}^w - \mu_{\Delta C})^2}{\exp f_{\Delta C}(0)}. \quad (83)$$

In our case, the ratio of the spectra is

$$\frac{\exp f_{\Delta C}^w(\omega)}{\exp f_{\Delta C}(\omega)} = \frac{b^w (R^{-1})^2 \frac{\exp f(\omega)}{\exp f^w(\omega)}}{b (R^{-1})^2} \quad (84)$$

$$= \exp \left( \int Z(\kappa) f^w(\kappa) d\kappa \right) \exp(-f^w(\omega)) \frac{\exp f(\omega)}{b (R^{-1})^2}, \quad (85)$$

and the difference in the means is

$$\mu_{\Delta C}^w - \mu_{\Delta C} = -R^{-1} \frac{\alpha}{2} (1 - R^{-1}) \left( b^w (R^{-1})^2 - b (R^{-1})^2 \right). \quad (86)$$

So the KL divergence is, ignoring additive constants,

$$KL = \exp \left( \int Z(\kappa) f^w(\kappa) d\kappa \right) \int \exp(-f^w(\omega)) \frac{\exp f(\omega)}{b (R^{-1})^2} d\omega \quad (87)$$

$$+ \int f^w(\omega) d\omega - \int Z(\kappa) f^w(\kappa) d\kappa \quad (88)$$

$$+ \frac{(R^{-1} \frac{\alpha}{2} (1 - R^{-1}))^2 \left( b^w (R^{-1})^2 - b (R^{-1})^2 \right)^2}{\exp f_{\Delta C}(0)}. \quad (89)$$

The derivative with respect to  $f^w(m)$  is

$$\frac{dKL}{df^w(m)} = - \exp \left( \int Z(\kappa) f^w(\kappa) d\kappa \right) \exp(-f^w(m)) \frac{\exp f(m)}{b (R^{-1})^2} \quad (90)$$

$$+ Z(m) \exp \left( \int Z(\kappa) f^w(\kappa) d\kappa \right) \int \exp(-f^w(\omega)) \frac{f(\omega)}{b (R^{-1})^2} d\omega + 1 - Z(m) \quad (91)$$

$$+ 2 \frac{(R^{-1} \frac{\alpha}{2} (1 - R^{-1}))^2 \left( b^w (R^{-1})^2 - b (R^{-1})^2 \right)}{\exp f_{\Delta C}(0)} \exp \left( \int Z(\kappa) f^w(\kappa) d\kappa \right) Z(m). \quad (92)$$

Evaluating at  $f^w = \bar{f}$ , we obtain  $\frac{dKL}{df^w(m)}|_{f^w=\bar{f}} = 0$ , as we would expect. The second derivative is

$$\frac{d^2 KL}{d(f^w(m))^2} = \exp \left( \int Z(\kappa) f^w(\kappa) d\kappa \right) \exp(-f^w(m)) \frac{\exp f(m)}{b (R^{-1})^2} \quad (93)$$

$$- \exp \left( \int Z(\kappa) f^w(\kappa) d\kappa \right) \exp(-f^w(m)) \frac{\exp f(m)}{b (R^{-1})^2} Z(m) \quad (94)$$

$$+ Z(m)^2 \exp \left( \int Z(\kappa) f^w(\kappa) d\kappa \right) \int \exp(-f^w(\omega)) \frac{\exp f(\omega)}{b (R^{-1})^2} d\omega \quad (95)$$

$$- Z(m) \exp \left( \int Z(\kappa) f^w(\kappa) d\kappa \right) \exp(-f^w(m)) \frac{\exp f(m)}{b (R^{-1})^2} \quad (96)$$

$$+ 2 \frac{(R^{-1} \frac{\alpha}{2} (1 - R^{-1}))^2 \left( b^w (R^{-1})^2 - b (R^{-1})^2 \right)}{\exp f_{\Delta C}(0)} \exp \left( \int Z(\kappa) f^w(\kappa) d\kappa \right) Z(m)^2 \quad (97)$$

$$+ 2 \frac{(R^{-1} \frac{\alpha}{2} (1 - R^{-1}))^2}{\exp f_{\Delta C}(0)} \left( \exp \left( \int Z(\kappa) f^w(\kappa) d\kappa \right) Z(m) \right)^2. \quad (98)$$

Evaluating now at  $f^w = f$ ,

$$\frac{d^2 KL}{dl f^w(m)^2} \Big|_{f^w=f} = (Z(m) - 1)^2 + 2 \left( R^{-1} \frac{\alpha}{2} (1 - R^{-1}) \right)^2 \frac{b(R^{-1})^2}{\exp f_{\Delta C}(0)} Z(m)^2 \quad (99)$$

$$= (Z(m) - 1)^2 + 2 \left( R^{-1} \frac{\alpha}{2} (1 - R^{-1}) \right)^2 Z(m)^2. \quad (100)$$

So the weights across frequencies in the KL divergence are approximately a function of  $Z(\omega)^2$ .

## References

**Dew-Becker, Ian**, “The pricing of economic risks under time-separable and recursive preferences,” 2016. Working paper.