

On the Determination of General Scientific Models with Application to Asset Pricing ¹

A. Ronald Gallant
Fuqua School of Business
Duke University
Durham NC 27708-0120 USA

Robert E. McCulloch
Graduate School of Business
University of Chicago
Chicago IL 60637-6040 USA

November 2003
This draft March 2008

Abstract

We consider a consumption based asset pricing model that uses habit persistence to overcome the known statistical inadequacies of the classical consumption based asset pricing model. We find that the habit model fits reasonably well and agrees with results reported in the literature if conditional heterogeneity is suppressed by a sharp prior but that it does not fit nor do results agree if conditional heterogeneity, well known to be present in financial market data, is allowed to manifest itself. We also find that it is the preference parameters of the model that are most affected by the presence or absence of conditional heterogeneity, especially the risk aversion parameter. The habit model exhibits four characteristics that are often present in models developed from scientific considerations: (1) a likelihood is not available; (2) prior information is available; (3) a portion of the prior information is expressed in terms of functionals of the model that cannot be converted into an analytic prior on model parameters; (4) the model can be simulated. The underpinning of our approach is that, in addition, (5) a parametric statistical model for the data, determined without reference to the scientific model, is known. In general one can expect to be able to determine a model that satisfies (5) because very richly parameterized statistical models are easily accommodated. We develop a computationally intensive, generally applicable, Bayesian strategy for estimation and inference for scientific models that meet this description together with methods for assessing model adequacy. An important adjunct to the method is that a map from the parameters of the scientific model to functionals of the scientific and statistical models becomes available. This map is a powerful tool for understanding the properties of the scientific model.

Keywords: Scientific models, simulation, Bayes, MCMC, estimation, inference, asset pricing.

1 Introduction

This article is motivated by an asset pricing application, namely the habit persistence asset pricing model, that has characteristics in common with all modern general equilibrium models. And this article contributes to that empirical literature. But many models derived from the principles of the physical, biological, social, or engineering sciences exhibit the same characteristics that make statistical analysis of the habit model challenging. These characteristics are (1) a likelihood is not available; (2) prior information is available; (3) a portion of the prior information is expressed in terms of functionals of the model that cannot be converted into an analytic prior on model parameters; (4) the model can be simulated. We outline a general approach for the analysis of such models and then apply it to the habit model.

In some instances other methods are available. For example, if the only cause of difficulties is a small number of latent variables, then a data augmentation approach will likely be applicable and be less computationally intensive than the methods proposed here. We are not concerned with such models. Our concern is with models such as our application where there seems to be little else available other than what we propose here. This becomes doubly true when data are sparse, as in our application, so that serious use of prior information becomes essential. Our proposals are especially helpful if, as in our application, some prior information may be expressed only in terms of functionals of the model. The methods proposed here also generate ancillary information that can help to interpret the scientific model in terms of its statistical properties and to cast model inadequacies into sharp relief.

Our approach depends on an assumption that (5) an adequate statistical model for the data is available. Because the statistical model is only estimated in large simulated data sets, richly parameterized models may be used so that (5) can usually be satisfied. Briefly, our proposal is as follows. Given (5), we can construct a map from the parameters of the scientific model to those of the statistical model such that a point in the parameter space of the scientific model and its image under the map both correspond to the same data generating process. Typically the parameters of the statistical model will live in a higher dimensional space than that of the scientific model. In Section 3 we obtain a Bayesian inference for the scientific model using the map to compute the likelihood. In Section 4, we work in terms of the statistical model and its parameters, but use the map as a key element in our prior

construction. This construction allows us to assess the adequacy of the scientific model. The methodology developed here allows the prior information to be expressed either directly on the parameters of the scientific model or on functionals of the scientific model that can be evaluated via simulation.

The idea that we intend to convey by the term scientific model is that the model is discipline-based and statistically intractable. The idea that we intend to convey by the term statistical model is that it has been obtained by data-analytic considerations with the specific intent of being statistically tractable.

The discovery of the mapping from the parameters of the scientific model to those of the statistical model, which is an intermediate step of the methods proposed here, is often itself of interest. For instance, the statistical model must, perforce, be expressed entirely in terms of observables whereas scientific models often contain unobservables. Having a mapping from the subset of the parameters that control the unobservable features of the scientific model to the parameters of a statistical model consisting entirely of observables can be helpful in understanding the observable consequences of changes in a model's unobservable internal structure. The utility of this approach can be extended by using the same methods to find the map from the parameters of the scientific model to functionals of both the scientific and statistical models.

A Bayesian approach suggests itself for problems that exhibit the five characteristics just listed because the methodology gracefully accepts prior information into the analysis. Nonetheless, comparison is of interest and in the application we compare to the results of Bansal, Gallant, and Tauchen (2007), referred to as BGT hereafter, who use a synthesis of frequentist methods proposed by Smith (1993), Gouriéroux, Monfort, and Renault (1993), and Gallant and Tauchen (1996). This approach is logically distinct from ours but does make use of an auxiliary statistical model as an adjunct to estimation and inference as we do here. As will be seen, sparse data forces a simplification on BGT estimates that inhibits discovery of our findings.

Although we know of nothing in print, several people, notably Anthony A. Smith, Jr., Yale University, and Alan E. Gelfand, Duke University, having seen this work presented, told us that they have had similar thoughts along the same lines of using a statistical model to synthesize a likelihood but either did not attempt or did not succeed in making them practicable. See also Del Negro and Schorfheide (2004), Dejong, Ingram, Whiteman (1996,

2000) for closely related ideas under linearity assumptions and the references therein. Our implementation relies on modern object oriented programming methods, modern data structures, and a discretization at a critical point in the computations. Bringing these elements to bear on the problem seems to be both novel to this work and essential to success. We believe that our proposals for model assessment are new.

2 Scientific and Statistical Models

We shall use the notational conventions of time series analysis because most models of the sort considered here are dynamic. This is in no way essential because the results apply equally well to other data with a few obvious changes to notation.

Let the transition density of the scientific model be denoted as $p(y_t|x_{t-1}, \theta)$, $\theta \in \Theta$, where $x_{t-1} = (y_{t-1}, \dots, y_{t-L})$ if Markovian and $x_{t-1} = (y_{t-1}, \dots, y_1)$ if not. We assume that there is no direct information about $p(\cdot|\cdot, \theta)$. All that we can do is simulate data from $p(\cdot|\cdot, \theta)$ for given θ . If the model produces ergodic output, then a single long simulation for each setting of θ suffices for our purposes. If not, then many independently simulated replicates of the data are used.

Since we do not have access to $p(\cdot|\cdot, \theta)$, there is no direct way to compute the likelihood. Our approach is to find a parametric family of distributions that is capable of representing the process $\{y_t\}$. Specifically,

ASSUMPTION 1 We assume that there is a transition density $f(y_t|x_{t-1}, \eta)$, $\eta \in H$, and that there is a one-to-one map $g : \theta \mapsto \eta$ such that

$$p(y_t|x_{t-1}, \theta) = f(y_t|x_{t-1}, g(\theta)) \quad \theta \in \Theta \tag{1}$$

and that the form of $f(\cdot|\cdot, \eta)$ is known.

When we need a likelihood based on the unknown $p(\cdot|\cdot, \theta)$, we substitute $f(\cdot|\cdot, g(\theta))$. The model $f(\cdot|\cdot, \eta)$ is a statistical description of the observed data that we call the statistical model. Often this model will be known from the literature. In other cases it must be determined as part of the analysis. As richly parameterized models are permitted, success in finding an acceptable statistical model can be anticipated. It is to be emphasized that we only use the statistical model to fit large simulations from the scientific model (Section 3) or

when augmented by a strong prior dictated by the scientific model (Section 4) so that the fact that the data may be too sparse to support it is not a consideration.

When Assumption 1 is satisfied the likelihood is exactly that implied by the scientific model. When Assumption 1 is violated the likelihood is different from that implied by the scientific model. To use Poirier’s (1988) terminology, when Assumption 1 holds one is looking at the world through the window implied by the scientific model. When Assumption 1 is violated one is looking at the world through a different window.

One might deliberately choose to violate Assumption 1. For example, if satisfaction of (1) leads to a statistical model $f(y|x, \theta)$ with characteristics markedly different from what is known about the distribution of the data, one might deliberately opt for a simplification that does not exhibit these characteristics. The issues that arise in this connection are discussed in general in Subsection 3.5 and specifically for our application in Subsection 5.4.

A key motivation for implementing a Bayesian approach to the problem is the importance of using prior information. This can be critical when data are sparse. When data are sparse, prior information can be used to fill in model features about which the data says little but the literature says much thereby enabling extraction of features about which the data are informative.

The scientific model is built using subject matter knowledge. Thus, we expect that real prior information is available. This prior information may be expressible either in terms of elements of θ or in terms of characteristics ψ of the process. For example, in our application an element of ψ is the unconditional moment of a variable for which no data that corresponds to its meaning within the model exists. In general, ψ can be regarded as a point in the range of a vector of functionals

$$\Psi : p(\cdot|\cdot, \theta) \mapsto \psi \tag{2}$$

that is computable from a simulation. Thus, ψ is a function of θ through the composition $\theta \mapsto p(\cdot|\cdot, \theta) \mapsto \psi$. We capture both of these types of information in our prior $\pi(\theta)$ through the construction

$$\pi(\theta) \propto h(\theta, \psi(\theta)). \tag{3}$$

Note that since we will be using the Metropolis-Hastings algorithm to compute the posterior, we only need a function proportional to the prior. We shall also discretize θ on a finite grid so that any positive h will be integrable.

Figure 1 about here

We may not want to impose the belief that the scientific model holds exactly. We can capture this idea by recasting the problem so that η of the statistical model is viewed as the parameter of interest and constructing a prior that expresses a preference for η that are close to the manifold

$$\mathcal{M} = \{\eta \in H : \eta = g(\theta), \theta \in \Theta\}, \quad (4)$$

where Θ is the parameter space of the scientific model and H is the parameter space of the statistical model. Our prior construction uses a single parameter that we call κ to control prior beliefs about how close η should be to the manifold (Section 4). The smaller κ is, the more prior weight is placed on η close to the manifold. We assess the scientific model by seeing if the marginal posterior distributions of interpretable features of the statistical model are sensitive to the choice of κ . If changing our prior so as to support η farther from the manifold results in location shifts of the posteriors that are appreciable from a practical point of view, then we conclude that the evidence in the likelihood is against the restriction corresponding to the scientific model.

We illustrate the ideas in Figure 1. The scientific model is $p(y|\theta) = n(y; \theta, \theta^2)$, and the statistical model is $f(y|\eta) = n(y; \eta_1, \eta_2)$, where $n(y; \mu, \sigma^2)$ denotes the normal density with mean μ and variance σ^2 . The mapping of the parameter θ of the scientific model to the parameters (η_1, η_2) of the statistical model is $g : \theta \mapsto (\theta, \theta^2)$. In each panel of Figure 1 the actors in our piece are displayed as follows: (i) the curve depicts the manifold (4), (ii) the dotted contours are those of the likelihood of the statistical model, (iii) the shading depicts the prior on η corresponding to a choice of κ , (iv) the dots are draws from the posterior for the parameter θ of the scientific model mapped into $\mathcal{M} \in H$ using the map $\eta = g(\theta)$ and then jittered (without the jittering, all the draws would be on the manifold), and (v) the solid contours represent the posterior of η given the prior (iii).

The likelihood in the two left panels was obtained by simulating 50 observations from the scientific model with $\theta = 2$, or equivalently, $(\eta_1 = 2, \eta_2 = 4)$. The likelihood in the two right panels was obtained by simulating 50 observations from the statistical model with $\eta = (2.8, 4)$. The prior off the manifold corresponds to a small value of κ in the top two panels, and a larger value in the bottom two. Our method for assessment of the scientific model is based on the observation that when the data support the scientific model, increasing κ may cause the posterior to become more spread out, but it will not dramatically shift location (panels (1,1) and (2,1)). Conversely, when the data does not support the scientific

model, increasing κ will result in a shift of the posterior (panels (1,2) and (2,2)).

Of course, in high dimensional problems, we cannot simply look at the contours of the η posterior. Instead we must examine low dimensional marginals of interest. The densities depicted with the solid curves in the (1,2) and (2,2) panels of Figure 6 show the effect of increasing κ on a marginal of interest in our asset-pricing application. We view the shift in the distribution as evidence against the model.

At this point, the main conceptual ideas that we shall propose have been set forth. The devil is in the details, to which we now proceed. The reader who would rather see our substantive results first can skip to Section 5.

3 Bayesian Estimation of Scientific Models

We have two cases to consider. In the first case, we assume the scientific model is true and seek inference for θ . Here the statistical model is just a tool for computing the likelihood. In the second case we work in the context of the statistical model and η is the parameter. In this case, the scientific model is a source of prior information. We will consider the first case in this section and the second in Section 4. Note that we are never working on the product space $\Theta \times H$ as is the case in work closely related to ours (Del Negro and Schorfheide, 2004).

3.1 Computing the Map

Given θ , we need to uncover the map $g : \theta \mapsto \eta$ that satisfies (1) from a simulation $\{\hat{y}_t, \hat{x}_{t-1}\}_{t=1}^N$ of $p(y_t|x_{t-1}, \theta)$. The basic idea is twofold: one notes that, as defined, the map minimizes the Kullback-Liebler divergence between the models $p(y_t|x_{t-1}, \theta)$ and $f(y_t|x_{t-1}, \eta)$ and that one can use simulation to integrate with respect to the scientific model $p(y_t|x_{t-1}, \theta)$. Intuitively this means that we can choose N , the simulation size, so large that the simulated data gives us all the information we need about the nature of the process at the given θ and then find the corresponding η by maximizing the likelihood of the simulated data under the statistical model $f(\cdot|\cdot, \eta)$. We find the η which gives us the same kind of data as θ . More formally, we are finding the η that puts the Kullback-Liebler divergence $d(f, p) = \iint [\log p(y|x, \theta) - \log f(y|x, \eta)] p(y|x, \theta) dy p(x|\theta) dx$ to zero by minimizing $d(f, p)$ with respect to η and are noting that $\iint \log p(y|x, \theta) p(y|x, \theta) dy p(x|\theta) dx$ does not have to be computed to solve this minimization problem. We approximate the integral that does have to be

computed in the usual way: $\int f \log f(y|x, \eta) p(y|x, \theta) dy p(x|\theta) dx \approx \frac{1}{N} \sum_{t=1}^N \log f(\hat{y}_t | \hat{x}_{t-1}, \eta)$. (Or by $\frac{1}{R} \sum_{r=1}^R \frac{1}{n} \sum_{t=1}^n \log f(\hat{y}_{t,r} | \hat{x}_{t-1,r}, \eta)$ if not ergodic. We assume ergodicity hereafter; if not, the requisite modifications are obvious.) Thus, upon dropping the division by N , the map is computed as

$$g : \theta \mapsto \operatorname{argmax}_{\eta} \sum_{t=1}^N \log f(\hat{y}_t | \hat{x}_{t-1}, \eta).$$

Our algorithm will incorporate a simple approach for computing this mle.

3.2 A Metropolis Algorithm for θ

We will use the Metropolis algorithm to compute the posterior distribution of θ . This algorithm will have to accommodate various forms of prior information and a computational scheme for obtaining the mle defining the map g .

The Metropolis algorithm is an iterative scheme generating a sequence of θ values according to a Markov chain whose stationary distribution is the posterior. As in all Bayesian analysis, we must specify our prior and likelihood. For our Metropolis chain we must also specify a Markov chain in θ used to propose new values.

Let $\mathcal{L}(\theta)$ denote the likelihood assuming that (1) holds. To compute it we may use

$$\mathcal{L}(\theta) = \prod_{t=1}^n f(y_t | x_{t-1}, g(\theta)),$$

where (y_t, x_{t-1}) denotes the observed data and n the sample size. Let $\pi(\theta)$ denote the prior distribution on θ . As discussed in Section 2, in order to compute this prior $\pi(\theta)$ we may need the value ψ taken on by the functionals Ψ given by (2). Let q denote our Metropolis proposal. For a given θ , $q(\theta, \theta^*)$ defines a distribution of potential new values θ^* .

Given a current θ^o and the corresponding $\eta^o = g(\theta^o)$, we obtain the next pair (θ', η') as follows:

1. Draw θ^* according to $q(\theta^o, \theta^*)$.
2. Draw $\{\hat{y}_t, \hat{x}_{t-1}\}_{t=1}^N$ according to $p(y_t | x_{t-1}, \theta^*)$.
3. Compute $\eta^* = g(\theta^*)$ and ψ^* from the simulation $\{\hat{y}_t, \hat{x}_{t-1}\}_{t=1}^N$.
4. Let $\alpha = \min \left(1, \frac{\mathcal{L}(\theta^*) \pi(\theta^*) q(\theta^*, \theta^o)}{\mathcal{L}(\theta^o) \pi(\theta^o) q(\theta^o, \theta^*)} \right)$.
5. With probability α , $(\theta', \eta') = (\theta^*, \eta^*)$, otherwise $(\theta', \eta') = (\theta^o, \eta^o)$.

Steps 1, 4, and 5 are just the standard Metropolis algorithm. Steps 2 and 3 are essential features of our approach. If the proposed θ in Step 1 violates a support condition that can be checked without running Step 2, one skips Step 2 because α in Step 4 will be zero.

In order to complete the specification of our algorithm we need to choose a q . We shall also propose a particular approach to the computation of η in Step 3. These are the next topics.

3.3 Choice of θ Proposal

To specify our algorithm we must choose a proposal transition density q for θ . To compute the likelihood at a proposed θ , the scientific model must be simulated. For a sophisticated scientific model, this simulation may involve significant computation. Moreover, there could well be a call to a nonlinear optimizer or nonlinear equation solver that needs starting values involved in this simulation. This motivates us to consider proposing small changes in θ so that computational results from the old θ may be used in doing the computations for the proposed θ . In particular, if θ is not changed too much, results from the previous computation can be used as starting values for the new one. The cost of this strategy is in dependence in the Markov Chain. If we limit ourselves to small changes, it may take us a while to navigate from one place in the parameter space to another.

We start by discretizing θ because, as seen later, discretization permits significant improvements in computational efficiency. For the i^{th} component of θ we choose $a_i < b_i$, and s_i . We then let θ_i take on the values $a_i + js_i$ where j ranges from 1 to g_i which is equal to the integer part of $(b_i - a_i)/s_i$. Thus, θ_i takes values between a_i and b_i on a grid of mesh s_i .

To propose a new θ we first randomly choose a component to change, with each component having the same chance of being chosen. If the i^{th} component is chosen, there is some j such that the current $\theta_i = a_i + js_i$. We choose a set of distributions $q_i(j, k)$ on $\{1, 2, \dots, g_i\}$ where i is the θ component, j is the current grid position of that component, and k denotes the random new grid position to be drawn. We draw $k \sim q_i(j, \cdot)$ and let θ^* be obtained from θ by changing the i^{th} component from $a_i + js_i$ to $a_i + ks_i$.

To specify the $q_i(j, k)$ we choose a σ_i for the i^{th} component of θ and let

$$q_i(j, k) \propto \begin{cases} \exp(-\frac{1}{2\sigma_i^2}(k - j)^2) & k \neq j \\ 0 & \text{else} \end{cases}$$

The choice of σ_i determines the number of s_i that we tend to move. We assign 0 probability to proposing that we stay put as there is no point in proposing that we go to where we are.

To run the θ Metropolis chain, we have to choose a starting value for each θ_i . The choice of a_i and b_i is not critical; a_i and b_i can be set so that the intervals (a_i, b_i) cover the support of the posterior by a wide margin without noticeably degrading performance. The choice of s_i is crucial. We will move away from the starting value in integer multiples of s_i . The combination of the choice of s_i and σ_i determines the size of the changes that q proposes. The choice of s_i determines the accuracy of our inference. When we choose s_i we are saying that, as a practical matter, we only need to know θ_i in terms of s_i units. Two θ 's that differ in component i by less than s_i are virtually the same as a practical matter. Since computation is expensive, we should not waste resources by determining θ on a finer scale than we actually care about.

3.4 Computing the MLE of η with the Simulated Data

Step 3 of the Metropolis algorithm presented in Section 3.2 is the computation of the mle of η under the statistical model given the large simulated data set.

Since the $f(\cdot|\cdot, \eta)$ family is generally chosen to be flexible and high dimensional, this likelihood can be complicated. However, the simulated data set is large and we have a good starting value. In the notation of Section 3.2, η^o should be a good starting value in the search for η^* . This assumes that θ^* is not too different from θ as discussed in Section 3.3. In order to keep our analytical requirements to a minimum, we would like our method only to require the computation of the objective $\mathcal{L}(\eta) = \prod_{t=1}^N f(\hat{y}_t|\hat{x}_{t-1}, \eta)$.

Given these considerations, to find the mle we run a Markov chain for η using the simulated data. Since our goal is to find the mle and the sample size is large, we use a flat prior on η when running this chain. With the large sample size, the Markov chain will quickly move from the η to values close to η^* . We use a normal random walk Metropolis within Gibbs approach. That is, we first subdivide the η vector into subvectors. In the manner of a Gibbs sampler, we cycle through the subvectors one at a time. For subvector η_i , we use the normal proposal $q(\eta_i, \eta_i^*) \sim n(\eta_i, \Sigma_i)$ in a standard random walk Metropolis algorithm. Effectively, this is a simulated annealing optimization algorithm where the simulation size N is the temperature parameter because N is what controls the peakedness of the likelihood. A side benefit is that the chain for η also provides the scaling for the model assessment

strategy proposed in Section 4.

The computation of η requires start values, as does any nonlinear optimization. If θ is only moved slightly between iterates of in an MCMC chain to compute the posterior for θ , then the last computed value of η will be a good start for the next. This consideration becomes doubly important when the scientific model $p(y|x, \theta)$ contains an embedded nonlinear computation requiring start values as does the habit model (Section 5.2). These requirements argue for a random walk proposal density that makes small moves in connection with the MCMC chain used to get the posterior distribution of θ .

We choose a fixed number of steps to run this chain, and keep the visited η which has the highest likelihood under the simulated data. In our experience, it is relatively straightforward to choose (i) a simulation sample size which is large enough to ensure that the map is adequately recovered by the mle and (ii) a number of steps to iterate the Markov chain in η that will ensure the chain has finished moving away from the starting value of η . We might also note that putting θ on a grid is a considerable help here because it reduces the accuracy to which η^* needs to be computed.

This is a computationally costly part of our overall procedure. Since the simulation sample size N is large, each computation of the likelihood for the η chain can take a long time. Nonetheless, we have found that, because of the large N , this part of the procedure is remarkably stable, even though the statistical model may actually be difficult to estimate on data samples of the size n that we actually observe.

The main reason for placing θ on a grid is that a significant reduction in computational time can be achieved. With θ on a grid, it takes only a modest amount of memory to store all previously computed values of η , $\mathcal{L}(\eta)$, $\pi(\theta)$, etc. in a binary tree indexed by θ . When θ is revisited, both the fact that it is a revisit and the information required for Step 4 of the Metropolis algorithm for θ (Subsection 3.2) can be quickly obtained by traversing the tree. The two costliest Steps 2 and 3 are thereby eliminated. We have found the C++ associative map to be an exceptionally convenient implementation of such a tree. By storing previous results in a tree and looking them up, the θ chain runs faster as it becomes longer.

One might consider alternative parameterization of the statistical model in order to improve performance of the MCMC chain for η . This will have no effect on the chain for θ using the algorithm of Subsection 3.2 because the maximum of the likelihood is not affected by reparameterization.

Figure 2 displays the results of ten runs of an η chain (from Section 5). Every run is clearly visible in the figure as a segment of 200 iterations. In the notation of Section 3.2, each segment displays the results of a Markov chain in η that is started at η^o , uses the data simulated from the scientific model at θ^* , and uses the likelihood of the statistical model coupled with a flat prior on η . On the vertical axis, the log-likelihood from the statistical model is plotted. We can see the likelihood quickly increase as the η value moves toward the mle. The segments level off at different likelihoods because they represent the likelihoods of different simulated data sets. Because of the large size of each simulated data set, $N = 50,000$ in this instance, the posterior is very tight around the mle and the chain quickly moves to a new level.

Figure 2 about here

3.5 Identification and Map Recovery

The scientific model and the statistical model must work in concert with one another. As Poirier (1988) argues, we think persuasively, what is most relevant is that the target audience views the likelihood as a reasonable approximation to the characteristics of the phenomenon under study. As he puts it, the target audience must accept the likelihood as the window through which they are willing to view the world. As he points out, this creates a tension between parsimony and generality. He suggests resolving this tension by a sensitivity analysis. We are in this situation with respect to our application in Subsection 5.4 and we address the tension by following Poirier’s proscription in Subsection 5.7.

The scientific and the statistical models must satisfy three requirements if they are to work concert with one another. We consider each in turn.

The first requirement is that the map $g(\theta) = \eta$ should not be one-to-many. This is equivalent to stating that the statistical model should be identified by simulations from the scientific model. Identification of the statistical model by simulations from the scientific model entails some obvious conditions such as that the support of the statistical model should include the support of the scientific model. A violation of this requirement can occur if the scientific model has fewer random shocks than the dimension of y thereby causing the support of the scientific model to be a lower dimensional submanifold of the support of the statistical model. This can happen inadvertently if a proposed θ implies a singular variance matrix somewhere within the scientific model. However, this is usually an easily checked

support condition.

Often the statistical model is effectively linear in its parameters, as it is in our application, so that a lack of identification can only occur by the scientific model generating observables that lie in a lower dimensional space than presumed by the statistical model. One can protect against this by checking support conditions on θ as just discussed. If the statistical model is not linear then identification problems can arise. An example is the statistical model $y = \eta_1 \exp(\eta_2 x) + \eta_3 \exp(\eta_4 x) + e$ with simulated data from a scientific model that is actually $y = \theta_1 \exp(\theta_2 x) + e$ in disguise. In this case, either $\eta_3 = 0$ and η_4 is not identified or $\eta_2 = \eta_4$ and only the sum $\eta_1 + \eta_3$ is identified. We have found through experience in a wide variety of contexts that one can detect this situation by checking the η chain described in Subsection 3.4 for a unit root. One way to do this is to find the eigen vector ℓ of the variance matrix of the chain $\{\eta_t\}_{t=1}^N$ with largest eigen value and examine the sequence of inner products $\{\ell' \eta_t\}_{t=1}^N$ for a unit root.

Lack of identification of the statistical model does not matter in the computation of the θ chain described in Subsection 3.2 as long as the likelihood is actually maximized at the computed η . The implied map $g(\theta) = \eta$ will be one-to-many but often considerations similar to estimability in less than full rank linear models come into play so that the features of the statistical model that are of interest in an application have the same value regardless of which maximizing value of η is chosen. These considerations are discussed in Gallant (1987). The methods proposed in Section 4 would have to be modified if η is not identified.

An omnibus check for identification failure and nearly anything else that can go wrong with the algorithm of Subsection 3.4, such as poor start values or not running the chain long enough to compute η^* to sufficient accuracy, is to run a regression of all computed η^* on low degree polynomials in the corresponding values of θ . If the R^2 are high one has some assurance that the statistical model is identified and that accuracy is adequate. For the application described in Section 5, the 25%, 50%, and 75% quantiles of the R^2 for regressions of each component of η^* on a cubic in θ are 0.91, 0.97, and 0.98.

The second requirement is that (1) should hold. One may check this requirement for specific values of θ by (a) fitting the statistical model to a scientific model simulation, (b) simulating from the fitted statistical model, and (c) checking to see if the empirical distributions of the two simulations match. An illustration of this exercise is Figure 3. This can serve as a partial empirical check if (1) cannot be established analytically.

The third requirement is that the statistical model $f(y|x, \eta)$ should be a reasonable model for the data. One would usually prefer not to allow the statistical model to exhibit characteristics that are known to be counterfactual or, to use Poirier’s phrasing, are known to cause a large fraction of one’s audience not to be willing to accept the view through such a window. If a statistical model that satisfies (1) is at odds with the data, then the second and third requirements are in conflict. If one opts for the third requirement in violation of the second, then one owes both the scientist who proposed the model and other readers a sensitivity analysis. This must be within reason, however, because, as Poirier remarks, both editors and readers have limits to their patience when it comes to sensitivity analyses. Violating the second requirement may make the map $g(\theta) = \eta$ many-to-one but this is not a problem with a proper prior on θ .

4 Inference Off the Manifold: Model Assessment

In this section $\eta \in H$ becomes the parameter of interest. The scientific model $p(y|x, \theta)$ may be viewed as a sharp prior that restricts η to lie on the manifold $\mathcal{M} \subset H$. What one would like to do is see how results change as this prior is relaxed. However, once we have moved off the manifold \mathcal{M} we can no longer view results from the perspective of the scientific model $p(y|x, \theta)$ and must view them from the perspective of the statistical model $f(y|x, \eta)$ because the scientific model loses meaning off the manifold. Therefore, seeing how results change must be taken to mean seeing how the marginal posterior distribution of a parameter or functional of the statistical model changes. Denote the vector of functionals of the statistical model of interest by

$$\Upsilon : f(\cdot|\cdot, \eta) \mapsto v. \tag{5}$$

For convenience, if an element of η is of interest, we make it an element of v .

We assume that we have a discrete set of points on the manifold $\{\eta_j \in \mathcal{M} : j = 1, \dots, G\}$. The analysis of Section 3 generates a discrete set of points $\{\theta_j \in \Theta : j = 1, \dots, G\}$ at which the map g has been evaluated; putting $\eta_j = g(\theta_j)$ provides such a set of points. Relaxation of the prior will be formulated in terms of a weighted distance of η from the manifold \mathcal{M} . We can cheaply compute the distance from η to the manifold as

$$d(\eta, \mathcal{M}) = \min_{j=1, \dots, G} (\eta - \eta_j)' A_j (\eta - \eta_j), \tag{6}$$

where A_j are scaling matrices. Let \hat{j} denote the index j at which the minimum occurs, let \hat{J} denote the map $\hat{J} : \eta \mapsto \hat{j}$, and let $\hat{h}(\eta) = \eta_{\hat{j}(\eta)}$.

The prior we propose is the product of preferences along the manifold, preferences about how close the statistical model is to the manifold, and general preferences about η

$$\pi_\kappa(\eta) \propto w_1[\hat{h}(\eta)] \exp\left(-\frac{d(\eta, \mathcal{M})}{2\kappa}\right) w_3(\eta), \quad (7)$$

where $w_1(\eta)$ and $w_3(\eta)$ are suitably chosen positive functions and we assume the middle $w_2(\eta, \kappa)$ term assures integrability. The three terms in the product (7) correspond to our three kinds of preferences regarding η . The prior becomes more diffuse and the scientific model less influential as the scale factor κ increases. Note that none of the individual terms is thought of as being a prior in its own right, each is just a part of the overall construction. To check that (7) results in a reasonable prior, draws of η may be simulated from the prior and prior marginals of interest checked (see Figure 6).

The functions $w_1(\eta)$ and $w_3(\eta)$ in (7) may be related to the preferences along \mathcal{M} that we had for θ of the scientific model, but there is no logical necessity that this be the case. If we desire to use roughly the same kind of preferences along \mathcal{M} for η as we used for θ , we can use $\eta_j = g(\theta_j)$ to generate our points on the manifold, put $w_1[\hat{h}(\eta)] \propto \pi(\theta_{\hat{j}(\eta)})$, where $\pi(\theta)$ is given by (3), and put $w_3(\eta) = 1$. Note that $\pi(\theta_{\hat{j}(\eta)})$ is a composite function that depends only on η that is easily retrieved from our stored map. With these choices, if η_1 and η_2 both have the same distance from the manifold, then $\pi_\kappa(\eta_1)/\pi_\kappa(\eta_2) = \pi(\theta_{\hat{j}(\eta_1)})/\pi(\theta_{\hat{j}(\eta_2)})$. Recall that one of our motivating considerations was the problem of sparse data and that this problem is overcome by the introduction of prior information. A small κ may imply sufficient prior information for inference. With large κ , additional prior information may be needed in the form of a choice of $w_3(\eta)$.

The computation of the posterior distribution of η using the statistical model $f(y|x, \eta)$ and prior $\pi_\kappa(\eta)$ can be accomplished by a routine application of the Metropolis algorithm because $\pi_\kappa(\eta)$ is easily computable and an analytic expression for $f(y|x, \eta)$ is available.

Our proposal is that the scientific model be assessed by plotting a suitable measure of the location and scale of the posterior distribution of v against κ or, better, sequential density plots as in Figure 6. What one expects to see, for a well fitting scientific model, is that the location measure does not move by a scientifically meaningful amount as κ increases, which indicates that the model fits, and that the scale measure increases, which indicates that the

scientific model has empirical content. What we see in (Figure 6) is that as κ increases the scale of two functionals of interest increases, indicating empirical content, but that location also shifts, providing evidence against the model.

Choosing the scaling matrices A_j can be difficult. We suggest two automatic choices. In the application we put $A_1 = \dots = A_G = \Sigma_\eta^{-1}$ in (6) where Σ_η is computed as follows: Initialize to zero. Whenever the Metropolis-Hastings chain for computing the mle of η described in Subsection 3.4 must be run, update

$$\Sigma_\eta \leftarrow \Sigma_\eta + (\eta_1 - \eta_2)(\eta_1 - \eta_2)' \quad (8)$$

where η_1 is a point on the chain immediately after transients have died out and η_2 is the last point on the chain. This method of scaling the distance measure is reasonable because it puts η on the scale of the posterior: Distance is being measured in units of standard deviation.

The distance function (6) can be made invariant to a linear reparameterization of the statistical model and first order invariant to a nonlinear transformation by letting A_j be the inverse of the covariance of the draws η_t^j from the chain used to obtain $\eta_j = g(\theta_j)$ (Section 3.4). Set $A_j = \Sigma_{\eta_j}^{-1}$, where

$$\Sigma_{\eta_j} = \sum_{k=1}^K (\eta_{t_1+k\Delta}^j - \bar{\eta}^j)(\eta_{t_1+k\Delta}^j - \bar{\eta}^j)',$$

$K > \dim(\eta)$, $\Delta = (t_2 - t_1)/K$, t_1 is the index of the point where transients die out, t_2 is the index of the last point of the chain, and $\bar{\eta}^j = (1/K) \sum_{k=1}^K \eta_{t_1+k\Delta}^j$. The scaling proposed in the paragraph above will achieve approximate invariance if the Σ_{η_j} are relatively homogeneous.

One should note that using posterior draws from the scientific model to compute the image \mathcal{M} of the map $g(\cdot)$ does make use of the data to determine the prior $\pi_\kappa(\eta)$. We do not regard this as a problem because the θ draws will contain enough extreme values to make sure that the extent of \mathcal{M} is large enough. If one is particularly worried about this, one could run the θ chain with a smaller amount of data to make sure that \mathcal{M} is over explored.

5 Habit Persistence Asset Pricing Model

In this section we shall apply the proposed methods to the habit persistence asset pricing model of Campbell and Cochrane (1999), referred to as CC hereafter. Although it is widely viewed as a behavioral model and it is the result of an admitted attempt to reverse engineer

away the statistical inadequacies of the classical consumption based asset pricing model (Lucas, 1978), the habit model actually can be justified from plausible micro-foundations (Guisar, 2005). The habit model exhibits all the characteristics discussed in Section 1: (1) the likelihood is not available; (2) prior information on model parameters is available; (3) prior information in the form of restrictions on model functionals is available; (4) the model can be simulated; (5) a generally accepted statistical model for its data is available.

In the remainder of this section, we describe the data, introduce the habit model, and apply the methods that we have proposed. We conclude that the habit model is not supported by the data using indirect evidence, Figure 6, and a direct statistical test.

We also find, Table 1, that statistical methods that rely on a countefactual assumption that the conditional variance of the data is homogeneous generate reasonable parameter estimates but that when conditional heterogeneity is taken into account the parameter estimates become implausible. From the implied map we discover in Figure 5 that a sharply delineated subset of the habit model’s parameters control the observable conditional heterogeneity. This subset contains the parameters that describe preferences of which the risk aversion parameter is particularly important. BGT argue that the habit model is not supported by the data using indirect evidence that differs from ours but their direct frequentist test accepts the habit model. Our results suggest that BGT’s failure to formally reject is a consequence of using statistical methods that impose homogeneity of conditional variance under both the null and the alternative. However, without the use of prior information, one is nearly forced to impose homogeneity and other undesirable restrictions on the alternative due to the sparseness of the data. Commentary in BGT speaks to this issue.

5.1 The Data

The data are annual, 1929–2001, non-durables and services consumption C_t^a and price P_{dt}^a of a value weighted portfolio comprised of all stocks listed on the New York and American exchanges; both series are adjusted for inflation and population growth. Data sources and collection protocol are described in BGT. The corresponding annual consumption growth series and stock returns series are $\Delta c_t^a = \log(C_t^a) - \log(C_{t-1}^a)$ and $r_{dt}^a = \log(P_{dt}^a) - \log(P_{d,t-1}^a)$. In our previous vector notation, we have $y_t = (\Delta c_t^a, r_{dt}^a)'$, $t = 1, \dots, n = 72$. Annual data is used rather than monthly or quarterly data because seasonality issues are avoided and the historical record is longer.

We next discuss how to simulate the habit model. The model is simulated at the monthly frequency and then aggregated to the annual frequency. This is consistent with CC and BGT and is common in this literature.

5.2 Model Description

The intuitive notions behind any consumption based asset pricing model are that agents receive wage income and dividend income from which they purchase consumption. Agents seek to reallocate their consumption over time by trading shares of stock that pay a random dividend and bonds that pay interest with certainty. The driving processes of such a model are the wage process and dividends or, equivalently, consumption and dividends because wages can be recovered by subtracting dividends from consumption due to the fact that someone must own the stock so the dividends must be received while for bonds someone pays interest and another receives so there are no net bond receipts. Agents are endowed with a utility function that depends on the entire consumption process. The parameters of this function determine their preferences. The first order conditions of their utility maximization problem determine the prices at which they are willing to trade securities. We shall describe the driving processes, the utility function, and the first order conditions of the habit model, in that order.

The driving processes of the habit model are consumption and dividend growth

$$\begin{aligned} c_t - c_{t-1} &= g + v_t \\ d_t - d_{t-1} &= g + w_t \end{aligned} \tag{9}$$

We follow the standard convention that lower case denotes the logarithm of an upper case variable; e.g. $c_t = \log(C_t)$, $d_t = \log(D_t)$. The errors (v_t, w_t) in (9) are normal with mean zero and variance $\text{Var}(v_t, w_t) = RR'$, where R is upper triangular with nonzero elements r_{11} , r_{12} , and r_{22} . At times it is more convenient to express the variance matrix in terms of $\sigma^2 = \text{Var}(v_t)$, $\sigma_w^2 = \text{Var}(w_t)$, and $\rho = \text{corr}(v_t, w_t)$.

Upon exponentiation to get C_t and summing over adjacent blocks of twelve, the consumption process does correspond conceptually to the data series C_t^a described above. On the other hand, the dividend process does not correspond conceptually to observable data primarily because what can be observed is strongly influenced by tax policy causing, e.g., corporations to shift dividend payments into or out of stock repurchases. Therefore, d_t is

to be regarded as a latent variable. It remains to consider how the returns process r_{dt}^a is simulated.

The habit model asserts that all agents in the economy are endowed with the same utility function $\mathcal{E}_0 \sum_{t=0}^{\infty} \delta^t [(C_t - X_t)^{1-\gamma} - 1]/(1 - \gamma)$ where X_t is habit, δ the time discount factor, and γ is the risk aversion parameter. Habit is determined by $X_t = C_t - S_t C_t$, where S_t , called the surplus ratio, is the state variable of the model and is designed to behave as if it followed a discretely sampled conditionally heteroskedastic diffusion. Specifically, and recalling the upper and lower case convention,

$$s_t - \bar{s} = \phi(s_{t-1} - \bar{s}) + \lambda(s_{t-1})v_t, \quad (10)$$

where $\lambda(s) = \frac{1}{\bar{s}}\sqrt{[1 - 2(s - \bar{s})] - 1}$ if $s_t \leq s_{\max}$ and zero else; \bar{s} and s_{\max} are computed from model parameters $\theta = (g, r_{11}, r_{12}, r_{22}, \phi, \delta, \gamma)$ as $\bar{S} = \sqrt{[(r_{11}^2 + r_{12}^2)(\gamma/(1 - \phi))]}$ and $s_{\max} = \bar{s} + (1 - \bar{S}^2)/2$.

Agents are presumed to be so numerous that each can solve their own utility maximization problem without regard to the actions of the others. Under this assumption, the price dividend ratio satisfies $P_{dt}/D_t = V_{\theta}(S_t)$, where $V_{\theta}(\cdot)$ solves the integral equation

$$0 = \mathcal{E}_t \left\{ V_{\theta}(S_t) - \delta \left(\frac{S_{t+1}C_{t+1}}{S_t C_t} \right)^{-\gamma} \left(\frac{D_{t+1}}{D_t} \right) [1 + V_{\theta}(S_{t+1})] \right\} \quad (11)$$

and \mathcal{E}_t denotes conditional expectation given S_t . These are the first order conditions of each agent's optimization problem. The solution method we use is described in BGT. Two aspects of the solution method are relevant to us: the method uses a long simulation of (C_t, D_t, S_t) to approximate an integral by an average and uses a nonlinear equation solver that needs starting values to compute $V(\cdot)$. An important consequence is that if two adjacent draws θ_i and θ_{i+1} in an MCMC chain are close together then the answer left over from solving for draw i makes a good starting value for solving for draw $i + 1$.

What one does, then, for θ given, is first simulate consumption and dividends using (9) and compute the surplus ratio from them using (10) to get the simulation $\{\hat{C}_t, \hat{D}_t, \hat{S}_t\}_{t=1}^{12N}$, where $12N$ is to have a simulation of length N after aggregation. Second use the simulation and the BGT method to solve (11) for $V_{\theta}(\cdot)$. Third, compute stock price as $\hat{P}_{dt} = \hat{D}_t V_{\theta}(\hat{S}_t)$. Fourth, compute geometric stock returns using $\hat{r}_{dt} = \log(\hat{P}_{dt}) - \log(\hat{P}_{d,t-1})$. A computation similar to the second through third steps provides the geometric risk free interest rate \hat{r}_{ft} . Fifth, compute functionals of interest from the simulation. An important functional is

$\mathcal{E}(r_{ft}^a)$ computed by annualizing $\mathcal{E}(r_{ft}) \doteq (\sum_{t=1}^{12N} r_{ft})/(12N)$. Last, aggregate $\{\hat{C}_t, \hat{r}_{dt}\}_{t=1}^{12N}$ by summing consecutive blocks of twelve to get the annual series $\{\hat{C}_t^a, \hat{r}_{dt}^a\}_{t=1}^N$ and put $\hat{y}_t = (\log \hat{C}_t^a - \log \hat{C}_{t-1}^a, \hat{r}_{dt}^a)$.

5.3 Prior Information

Our discretization of θ (see Section 3.3) automatically restricts our prior support to the set of θ such that $a_i \leq \theta_i \leq b_i$ for each i . From the marginal posteriors (Figure 4) we can see that our choices of a_i and b_i are not influential since none of the distributions appear to be truncated.

A more complex support condition derives from the fact that a solution $V_\theta(\cdot)$ to (11) does not exist for all θ . We view this as an indication that such a θ is unreasonable and assign it zero prior probability. To protect against spurious failure to solve (11) we rely on small proposals for θ and a long list of recent solutions to (11) from which we choose the most promising start. Thus, if we fail to solve (11) we simply set the prior $\pi(\theta)$ equal to zero. In our Metropolis algorithm (Section 3.2), this determination is made at step 2. In the case of failure we set α to zero and proceed directly to step 5.

The rest of our prior information will be expressed through $\psi(\theta) = \mathcal{E}(r_{ft}^a)$ (a functional of $p(\cdot|\cdot, \theta)$, computed from the simulation as described above), the correlation between monthly consumption and dividend growth $\rho(\theta)$ (a simple function of θ), and ϕ (an element of θ). In the notation of (3), Section 2, we must define a function $h(\theta, \psi(\theta))$ that, at least approximately, captures this information.

In thinking about the risk free rate we first note that it is not directly observable because its computation from observed interest rates requires estimation of anticipated inflation (Mishkin, 1981). Because the evidence indicates that the risk free rate is low with a variance that is much smaller than the variance of inflation (Campbell, 2002), an attempt to generate a risk free rate series over our sample period would produce a series that would consist almost entirely of measurement error. We base our prior information on the values that Campbell determined from several long historical time series. We put most of our prior mass on θ such that $\psi(\theta) = \mathcal{E}(r_{ft}^a)$ is in the interval $.89\% \pm 1\%$.

A consequence of treating dividends as latent is that $\rho(\theta)$ is poorly determined by the observables. Thus, prior information is important. The parameter ϕ is an autoregressive parameter and we find that it is necessary to prevent the MCMC chain from putting ϕ too

close to one. We put prior mass on θ such that $\rho(\theta)$ is in the interval $.2 \pm .1$ and ϕ is in the (tight) interval $.9884 \pm .01$. The locations of these intervals are based on CC. Note that in terms of θ , $\rho(\theta) = \frac{r_{12}}{\sqrt{r_{12}^2 + r_{22}^2}}$.

We let

$$\pi(\theta) \propto e^{-\frac{1}{2}\left(\frac{\psi(\theta)-.89}{0.5}\right)^2} e^{-\frac{1}{2}\left(\frac{\rho(\theta)-.2}{0.5}\right)^2} e^{-\frac{1}{2}\left(\frac{\phi-0.9984}{0.005}\right)^2} I(\theta)$$

where $I(\theta)$ is an indicator function capturing the support conditions discussed above.

Although we use the normal kernel in our construction, nothing is normally distributed. For example, the reader is reminded that we discretize θ so that its support lies on a finite number of atoms. We use the normal kernel as a simple way of making sure that the resulting prior $\pi(\theta)$ is relatively small for θ corresponding to ψ , ρ , and ϕ outside the intervals corresponding to our prior beliefs. Certainly, none of the normal kernels in the expression above represents a prior in its own right. Each is merely a part of the overall construction just as in the case of the indicator function.

The inclusion of the middle term involving $\rho(\theta)$ creates prior dependence between r_{12} and r_{22} . Although we have no analytical knowledge of the function $\psi(\theta)$, CC's study of numerical solutions suggest that it is a complex non-linear function involving all components of θ . Thus, inclusion of the first term captures rich and complex prior information, involving every element of θ , that would be very hard to impose in any other way.

5.4 The Statistical Model and Prior

The statistical model is a bivariate GARCH system with normal innovations:

$$\begin{pmatrix} c_t^a - c_{t-1}^a \\ r_{dt}^a \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} + \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix} \begin{pmatrix} c_{t-1}^a - c_{t-2}^a - b_1 \\ r_{d,t-1}^a - b_2 \end{pmatrix} + \begin{pmatrix} r_{11,t-1} & r_{12,t-1} \\ 0 & r_{22,t-1} \end{pmatrix} \begin{pmatrix} z_{1t} \\ z_{2t} \end{pmatrix} \quad (12)$$

$$\begin{pmatrix} r_{11,t-1} \\ r_{12,t-1} \\ r_{22,t-1} \end{pmatrix} = \begin{pmatrix} \rho_1 \\ \rho_2 \\ \rho_3 \end{pmatrix} + \begin{pmatrix} \rho_{11} & \rho_{12} \\ \rho_{21} & \rho_{22} \\ \rho_{31} & \rho_{32} \end{pmatrix} \begin{pmatrix} |z_{1,t-2}| \\ |z_{2,t-2}| \end{pmatrix} + \tau \begin{pmatrix} r_{11,t-2} - \rho_1 \\ r_{12,t-2} - \rho_2 \\ r_{22,t-2} - \rho_3 \end{pmatrix} \quad (13)$$

where (z_{1t}, z_{2t}) are iid normal. The statistical model $f(y_t|x_{t-1}, \eta)$ has sixteen parameters $\eta = (b_1, b_2, b_{11}, b_{21}, b_{12}, b_{22}, \rho_1, \rho_2, \rho_3, \rho_{11}, \rho_{21}, \rho_{31}, \rho_{12}, \rho_{22}, \rho_{32}, \tau)$.

This is a standard type of model for these data (Tauchen and Hussey, 1991). One must remember that these are bivariate, annual data on consumption growth and stock

returns. While the BIC criterion indicates that a fat-tailed conditional density is required for univariate, daily stock returns (Chernov, Gallant, Ghysels, and Tauchen, 2003), it does not indicate that fat tails are required for bivariate, annual consumption growth and stock returns. The BGT analysis of the habit model assumes normality as implicitly does the Campbell and Cochrane (1999) calibration as indicated by their choice of which moments to match. These considerations and a desire to be able to compare with the results of others lead us to adopt this specification for the statistical model.

Figure 3 about here

On the other hand, the habit model behaves more like a model for high frequency data than a model for annual data. If one applies the methods for finding an auxiliary model proposed by Chernov, Gallant, Ghysels, and Tauchen (2003) to a simulation from the habit model at the parameters shown in column Bayes-GARCH in Table 1, one arrives at a model with a two-lag nonlinear mean function and fat-tailed GARCH innovations. In our view, a two-lag nonlinear mean function is too countefactual for annual data to be taken seriously. Fat-tailed innovations may be a different matter.

Figure 3 indicates quite clearly that the habit model does put fat tails into the annual returns series. Figure 3 is an application of the diagnostic technique proposed in Subsection 3.5 where the statistical model is fitted to simulations from the scientific model and then itself simulated. Densities estimated from these two simulations are then compared. Figure 3 indicates that the tails of the returns innovation in the statistical model must be like the t -distribution on 3.5 degrees of freedom to match the simulation of the habit model.

As discussed in Subsection 3, this puts the second requirement that (1) be satisfied in conflict with the third requirement that the the statistical model should be a reasonable model for the data. We address the conflict as proposed in Subsection 3 by conducting a sensitivity analysis.

What follows suggests that our statistical model is rich enough to extract information in the data that relates to the statistical adequacy of the habit model. When on the manifold, we are not limited by the number of observations. The relevant sample size on the manifold is the simulation size N , which is entirely under our control.

Off the manifold, when we iterate the η chain directly to explore model adequacy as described in Section 4, the relevant sample size is the number of observations n . For small κ , the influence of the prior $\pi_\kappa(\eta)$ is strong enough that the MCMC chain for η behaves well

with $w_3(\eta) \propto 1$. But when κ is large, the parameter τ that appears in (13) can move to explosive values and stick. We let $w_3(\eta) \propto e^{-\frac{1}{2}\left(\frac{\tau-0.6}{0.08}\right)^2}$. As in Section 5.3, this does not mean that the marginal prior of τ is normal. The other terms in (7) will involve τ . The normal kernel is a simple way for us to downweight η with τ far from .6. As discussed in Section 4, we use $w_1[\hat{h}(\eta)] \propto \pi(\theta_{\hat{J}(\eta)})$.

The statistical model with recursion (13) is denoted as GARCH hereafter. To compare our method to BGT, we shall also apply our methods using the statistical model with homogeneous variance denoted VAR hereafter.

5.5 Model Estimation: Scientific Prior Imposed

We ran the θ chain for 800,000 iterations discarding every 8 leaving 100,000. The chain was started near the mode of the posterior density to avoid problems with transients. The mode was determined from several initial runs of size 100,000. Visual inspection of time series plots of the draws (not shown) indicated that this strategy for eliminating transients was successful.

Table 1 about here

Table 1 reports our results for the VAR and GARCH statistical models together with BGT estimates. Estimates for the model as described are shown in the top half of the table. These estimates are for the monthly frequency. For convenience in interpretation, annualized estimates are displayed in the bottom half. Further down are estimates for $\mathcal{E}(r_{f,t}^a), \sqrt{\text{Var}}(r_{f,t}^a), \mathcal{E}(r_{d,t}^a), \sqrt{\text{Var}}(r_{d,t}^a)$.

In the columns of Table 1 that are labeled mode, the values shown are the mode of the multivariate posterior for $\theta = (g, r_{11}, r_{21}, r_{22}, \phi, \delta, \gamma)$, not the mode of each marginal. The values of the multivariate mode appear in the MCMC chain for θ and therefore do correspond to an actual simulation of the habit model whereas the vector of marginal modes or means might not and might not even satisfy our prior support conditions. The values for $\mathcal{E}(r_{f,t}^a), \sqrt{\text{Var}}(r_{f,t}^a), \mathcal{E}(r_{d,t}^a), \sqrt{\text{Var}}(r_{d,t}^a)$ in the columns labeled mode are computed from the mode of θ . Those in the columns labeled mean are the means of the marginal posterior distribution of $\mathcal{E}(r_{f,t}^a), \sqrt{\text{Var}}(r_{f,t}^a), \mathcal{E}(r_{d,t}^a), \sqrt{\text{Var}}(r_{d,t}^a)$. All standard deviations in the columns headed Bayes are computed from marginal posteriors. The BGT estimates optimize a criterion function and therefore do correspond to a simulation; the values shown for $\mathcal{E}(r_{f,t}^a), \sqrt{\text{Var}}(r_{f,t}^a), \mathcal{E}(r_{d,t}^a), \sqrt{\text{Var}}(r_{d,t}^a)$ are computed from that simulation.

As with any discrete time dynamic model, annualized parameter values for the habit model do not have the property that simulations obtained by running the model at the annual frequency using the annualized parameter values will have the same distribution as simulations obtained by aggregating a monthly simulation. For this reason, the values at the monthly frequency in the upper half of of Table 1 should be regarded as the definitive estimates of model parameters. The annualized values in the lower half of the table are to be regarded as only an aid in their interpretation.

Figure 4 about here

The Bayes-GARCH columns of Table 1 report computations that use the GARCH statistical model. Their most striking feature is the large deviation of the Bayes-GARCH estimate of mean stock returns $\mathcal{E}(r_{d,t}^a) \doteq 11\%$ from the value computed from the data of 6.02% and from the BGT estimate of 6.54%. The reason for this is that the BGT estimator completely ignores the conditional heterogeneity in the data. The auxiliary model of the BGT estimator, whose impact on the BGT estimate is analogous to the impact of the statistical model on the Bayes procedure we propose, is a VAR with homogeneous conditional variance. If we too use a VAR with homogeneous conditional variance to implement our Bayes estimator, the results in the column labeled Bayes-VAR are obtained, bringing the Bayes estimate of $\mathcal{E}(r_{d,t}^a)$ into agreement with the data and with the BGT estimates. The Bayes-VAR estimates are also in good agreement with the values reported in CC that were determined by matching to a set of moments that are not rich enough to identify models with conditional heterogeneity.

The reasons that BGT had to suppress GARCH in their auxiliary model are twofold. Observable data rather than prior information are used to achieve identification in BGT thereby creating the need for a four dimensional series. The data ($n = 72$) are too sparse to support a GARCH specification in four dimensions. We, on the other hand, have no such difficulties because we are fitting to large simulations, not to data.

It is of interest to view the effect that using the VAR statistical model has on marginal posteriors. These are shown in Figure 4. Elimination of the requirement that the habit model confront the conditional heterogeneity in the data results in large shifts in the marginal posterior of the utility parameters (ϕ, γ, δ) . This accounts for the dramatic left shift in the marginal posterior for $\mathcal{E}(r_{d,t})$ and dramatic variance reduction in the posterior for $\sqrt{\text{Var}(r_d)}$ that is seen in Table 1. This, following the logic of Section 4, can be regarded as strong evidence against the VAR specification of the statistical model.

Figure 5 about here

It is also of interest to note that it is the preference parameters (ϕ, δ, γ) of the scientific model that control the GARCH parameters $(\rho_{11}, \rho_{21}, \rho_{31}, \rho_{12}, \rho_{22}, \rho_{32}, \tau)$ of the statistical model with γ having by far the most influence. This can be discovered by inspecting the map $g : \theta \mapsto \eta$. Since this map is high-dimensional we focus our interest on the conditional correlation between annual stock returns and consumption growth for the year 2002. This correlation is important substantively because a high positive correlation implies that risk averse investors will require high expected returns to induce them to invest. The conditional correlation is a functional of the statistical model $f(\cdot|\cdot, \eta)$ determined by the map $\eta = g(\theta)$ and computed by applying the GARCH recursion (13) to the data. The last variance matrix given by the recursion is the conditional variance for the year 2002. In Figure 5 we plot the conditional correlation against the parameters $\theta = (g, r_{11}, r_{21}, r_{22}, \phi, \delta, \gamma)$ of the scientific model. The dots show the conditional correlation when the statistical model is GARCH and the circles when VAR. The difference between these two curves is the GARCH effect. In each panel, a θ_i is varied and the remaining θ held fixed at their posterior means. The solid vertical line is the posterior mean of θ_i when the statistical model is GARCH and the dashed line when VAR. The point where the solid line crosses the dots gives the conditional correlation when $\eta = g(\theta)$ is evaluated at the posterior mean under GARCH; similarly VAR for the dashed line and circles.

5.6 Model Assessment: Scientific Prior Relaxed

We now apply the methodology described in Section 4 to the habit model.

The functionals of interest are the mapping Υ_1 of $f(\cdot|\cdot, \eta)$ to the mean return on the stock portfolio over the period 2002–2102 and the mapping Υ_2 to the conditional correlation between the return on the stock portfolio and consumption growth for the year 2002. The conditioning event for both functionals is the 73 years of observed data. Υ_1 is computed from a realization obtained by simulating the GARCH model over the period 2002–2102; Υ_2 is obtained from the variance matrix for the year 2002 computed as described at the end of Subsection 5.5. Both depend on the data and η ; Υ_1 also depends on an initial seed that was the same for each η .

Figure 6 about here

Using the methods proposed in Section 4, we computed MCMC chains for three values of κ (1,20,100). The stationary distribution of each chain is the posterior for η under prior κ . We then evaluated the two functionals (Υ_1, Υ_2) at each η in the chain and plotted their densities in Figure 6. κ increases as we go down the rows. The left hand panels show prior (dashed lines) and posterior densities (solid lines) for the mean return and the right hand panels display the same for the correlation. For small κ , we have a tight prior and the prior and posterior are similar.

For larger κ , the posterior shifts. In the three right hand plots, the posterior distribution of the correlation shifts as κ increases while the prior remains reasonable. As κ increases we place less weight on η close to the scientific model. This is simple, interpretable evidence that the scientific model has trouble tracking a feature of the data. This feature is important substantively because a high positive correlation implies that risk averse investors will require high expected returns to induce them to invest.

The mean stock return, unlike a correlation, has a substantive meaning without having an independent statistical meaning. Nonetheless, coping with retirement plan options and the laws governing bequests has made the notion of a mean return over a long planning horizon meaningful to most of us. For the mean return, we also see that as κ increases, the distribution shifts. The posterior mean of the mean return over the hundred year horizon is 0.1 in the top left panel and 0.08 in the bottom left so that while the shift may not appear as dramatic as for the correlation functional it is substantively large. Also, the standard deviation increases from 0.0127 to 0.0172. In the context of the asset pricing literature, this is substantive evidence against the scientific model.

Figure 6 focuses on the marginal priors and posteriors of two particular functionals of interest. Of course, as κ changes, the entire sixteen dimensional prior and posterior of η is changing. To get a sense of how the posterior is moving, we consider κ equal 1, 20, or 100 as a choice of model and compute the posterior probability for the three different models with each having prior probability 1/3. That is, the pair $(f(\cdot|\cdot, \eta), \pi_\kappa(\eta))$ is considered to be a model and the posterior probability of each κ choice is proportional to $\int f(y|x, \eta) \pi_\kappa(\eta) d\eta$. The integrals were computed using method f_5 from Gamerman and Lopes (2006) section 7.2.1. Method f_5 requires draws from both the prior and the posterior but does not require the normalizing constant for the prior. We find that the posterior probability that κ equals 100 is 0.999. As we loosen up the prior and allow the posterior to move away from the

scientific manifold, it is able to find much higher likelihoods. If we include the scientific model ($\kappa = 0$) in our list of models and use prior probabilities 1/4, the posterior probability that κ equals 100 is 0.996.

5.7 Sensitivity Analysis

To check the sensitivity of our results to the normality assumption in (12), we recomputed the columns labeled Bayes-GARCH in Table 1 using a statistical model where the normal distribution for z_{2t} in (12) is replaced by the t on 3.5 degrees of freedom. The results are shown in Table 2. Comparing Tables 1 and 2 we see that the posterior distribution of θ is not affected in any substantive way. If anything, the evidence against the habit model is slightly stronger because the already overly large risk premium of $11.14\% - 1.21\% = 9.93\%$ has increased to $11.07\% - 0.73\% = 10.34\%$ and the mean g of consumption growth has become even more counterfactual.

Table 2 about here

6 Conclusion

We considered a consumption based asset pricing model that uses habit persistence to overcome the known statistical inadequacies of the classical consumption based asset pricing model. We found that the habit model fits reasonably well and agrees with results reported in the literature if conditional heterogeneity is suppressed but that it does not fit nor do results agree if conditional heterogeneity is allowed to manifest itself. We also found that it is the preference parameters of the model that are most affected by the presence or absence of conditional heterogeneity, especially the risk aversion parameter.

To obtain these results we proposed and implemented a general purpose Bayesian methodology for the analysis of complex models from the sciences. It relies on the ability to simulate from the scientific model, upon the availability of substantive prior information, and upon the willingness to use that prior information. Analysis is carried out by means of a richly parameterized statistical model $f(\cdot|\cdot, \eta)$ that is viewed as being the correct description of the distribution of the data. To assess the scientific model, we view it as imposing a severe prior $\pi_\kappa(\eta)$ on the statistical model. The correctness of the scientific model is assessed by relaxing $\pi_\kappa(\eta)$ and assessing the posteriors of scientifically meaningful functionals of $f(\cdot|\cdot, \eta)$.

If location does not change more than scientifically meaningful magnitudes as $\pi_\kappa(\eta)$ is relaxed, then the model is supported. If scale increases as $\pi_\kappa(\eta)$ is relaxed, then the model has empirical content.

References

- Bansal, Ravi, A. Ronald Gallant, and George Tauchen (2007), “Rational Pessimism, Rational Exuberance, and Asset Pricing Models,” *Review of Economic Studies* 5, 523–590.
- Campbell, John Y. (2002), “Consumption Based Asset Pricing,” Manuscript, Department of Economics, Harvard, Cambridge MA.
- Campbell, John Y., and John Cochrane (1999) “By Force of Habit: A Consumption-based Explanation of Aggregate Stock Market Behavior,” *Journal of Political Economy* 107, 205–251.
- Dejong, David N., Beth F. Ingram, Charles H. Whiteman (1996), “A Bayesian Approach to Calibration,” *Journal of Business and Economic Statistics* 14, 1–9.
- Dejong, David N., Beth F. Ingram, Charles H. Whiteman (2000), “Keynesian Impulses versus Solow Residuals: Identifying Sources of Business Cycle Fluctuations,” *Journal of Applied Econometrics* 15, 311–329.
- Del Negro, Marco, and Frank Schorfheide (2004), “Priors from General Equilibrium Models for VARS,” *International Economic Review* 45, 643–673.
- Gallant, A. Ronald (1987), “Identification and Consistency in Semiparametric Regression,” in Bewley, Truman F., ed. (1987), *Advances in Econometrics Fifth World Congress, Volume 1*, Cambridge University Press, New York, 145–170.
- Chernov, Mikhail, A. Ronald Gallant, Eric Ghysels, and George Tauchen (2003), “Alternative Models for Stock Price Dynamics,” *Journal of Econometrics* 116, 225–257 .
- Gallant, A. R. and G. Tauchen (1996), “Which Moments to Match?” *Econometric Theory* 12, 657–681.
- Gamerman, D., and H. F. Lopes (2006), *Markov Chain Monte Carlo: Stochastic Simulation for Bayesian Inference (2nd Edition)*, Chapman & Hall, Boca Raton, FL.

- Gourieroux, C., A. Monfort and E. Renault (1993) "Indirect Inference," *Journal of Applied Econometrics*, 8, S85–S118.
- Guvenen, Fatih (2005), "A Parsimonious Macroeconomic Model for Asset Pricing: Habit Formation or Cross-sectional Heterogeneity?" Manuscript, Department of Economics, University of Texas at Austin.
- Lucas, R. E., Jr. (1978), "Asset Prices in and Exchange Economy," *Econometrica* 46, 1429–1445.
- Mishkin, Frederick S. (1981) "The Real Rate of Interest: An Empirical Investigation," *Carnegie-Rochester Conference Series on Public Policy, The Cost and Consequences of Inflation* 15, 151–200.
- Poirier, Dale J., (1988) "Frequentist and Subjectivist Perspectives on the Problems of Model Building in Economics," *Journal of Economic Perspectives* 2, 121–144.
- Smith, A. A. (1993), "Estimating Nonlinear Time Series Models Using Simulated Vector Autoregressions," *The Journal of Applied Econometrics* 8, S63–S84.
- Tauchen, George, and Robert Hussey (1991), "Quadrature-Based Methods for Obtaining Approximate Solutions to Nonlinear Asset Pricing Models," *Econometrica*, 59, 371–396.
- Theil, Henri, and Arthur S. Goldberger (1961), "On Pure and Mixed Estimation in Economics," *International Economic Review* 2, 65–78.

Table 1. Parameter Estimates, Normal Errors

Parameter	BGT Estimates		Bayes-VAR			Bayes-GARCH		
	Estimate	Std. Err.	Mode	Mean	Std. Dev.	Mode	Mean	Std. Dev.
Monthly								
g	0.002116	0.000250	0.001639	0.001739	0.000258	0.001803	0.001780	0.000684
ψ_{11}	0.006151	0.000896						
ψ_{22}	0.036503	0.007716						
ρ_s	0.971900	0.015449						
μ_{dc}	-3.3587	0.0380						
r_{11}			0.006753	0.007326	0.000627	0.007254	0.007417	0.001903
r_{12}			0.001350	0.001451	0.000403	0.001350	0.001336	0.001068
r_{22}			0.003125	0.008109	0.006205	0.003125	0.018852	0.034435
ϕ	0.9853	0.0026	0.9861	0.9857	0.0024	0.9804	0.9818	0.0095
δ	0.9939	0.0005	0.9955	0.9937	0.0030	0.9898	0.9907	0.0070
γ	0.8386	0.2462	0.5726	0.9463	0.6179	1.0744	1.1747	1.7638
Annualized								
g	2.539	0.0866	1.9672	2.087	0.0895	2.164	2.136	0.2369
σ	2.1308		2.3857	2.5870	0.2202	2.5589	2.6106	0.2513
ρ	0.1650		0.1960	0.1943	0.0508	0.1830	0.1773	0.0507
σ_w	12.9118		1.0825	2.8090	2.1496	1.0825	6.5306	4.4984
ϕ	0.8372	0.0090	0.8450	0.8412	0.0084	0.7890	0.8023	0.0328
δ	0.9292	0.0018	0.9477	0.9269	0.0102	0.8845	0.8934	0.0244
γ	0.8386	0.2462	0.5726	0.9463	0.6179	1.0744	1.1747	1.7638
$\mathcal{E}(r_{dt}^a)$	6.54		6.00	7.58	0.6930	11.14	10.45	0.5487
$SDev(r_{dt}^a)$	16.9		20.49	21.48	1.5343	24.22	26.14	2.4735
$\mathcal{E}(r_{ft}^a)$	1.07		1.20	1.16	0.1389	1.21	0.99	0.1451
	$\chi^2(5) = 7.11 (0.21)$		reps = 800,000 by 8			reps = 800,000 by 8		

Notes: The values shown as monthly are the actual parameters of the habit model when simulated at a monthly frequency. To aid interpretation the location, scale, and discount parameters $g, \sigma, \sigma_w, \delta$ have been re-expressed as annual rates in percentage terms. The autoregressive parameter ϕ is adjusted from a monthly to annual frequency; ρ, γ do not require adjustment. BGT estimates, from Bansal, Gallant, and Tauchen (2004), use data on the price dividend ratio and the consumption dividend ratio in addition to consumption growth and stock returns and impose $\mathcal{E}(r_{f,t}^a) = 0.89\%$ and cointegration among consumption, dividends, and price. Variance parameters relate as

$$\text{Var} \begin{pmatrix} c_t - c_{t-1} \\ d_t - d_{t-1} \end{pmatrix} = \begin{pmatrix} \sigma^2 & \rho \sigma \sigma_w \\ \text{sym} & \sigma_w^2 \end{pmatrix} = \begin{pmatrix} r_{11}^2 + r_{12}^2 & r_{12} r_{22} \\ \text{sym} & r_{22}^2 \end{pmatrix} = \begin{pmatrix} \psi_{11}^2 & \psi_{11}^2 \\ \text{sym} & \psi_{11}^2 + 2\psi_{22}^2(1 - \rho_s)^{-1} \end{pmatrix}$$

where BGT parameters $\mu_{dc}, \rho_s,$ and ψ_{22} are the location, autoregressive, and scale parameters of the cointegration relation between c_t and d_t . In the data, the mean of r_d^a is 6.02% and the standard deviation is 19.29%; for consumption growth the values are 1.95% and 2.24%. The mode of $\theta = (g, \sigma, \rho, \sigma_w, \delta, \gamma)$ is the mode of the multivariate posterior. All standard deviations are from marginal posteriors.

Table 2. Parameter Estimates, t-Errors

Parameter	Bayes-GARCH-t		
	Mode	Mean	Std. Dev.
Monthly			
g	0.001968	0.001883	0.000247
r_{11}	0.006753	0.007263	0.000659
r_{12}	0.001350	0.001318	0.000404
r_{22}	0.003125	0.021659	0.013042
ϕ	0.9804	0.9810	0.0034
δ	0.9904	0.9902	0.0025
γ	1.0744	1.1988	0.6744
Annualized			
g	2.361	2.260	0.2961
σ	2.3857	2.5569	0.2320
ρ	0.1960	0.1786	0.0514
σ_w	1.0825	7.5028	4.5167
ϕ	0.7889	0.7944	0.0331
δ	0.8907	0.8887	0.0266
γ	1.0745	1.1988	0.6745
$\mathcal{E}(r_{dt}^a)$	11.07	10.99	0.4828
$SDev(r_{dt}^a)$	25.58	27.58	2.7776
$\mathcal{E}(r_{ft}^a)$	0.73	0.88	0.1477

reps = 800,000 by 8

Notes: As for Table 1 except the error distribution of the auxilliary model is Student-t on 3.5 degrees of freedom.

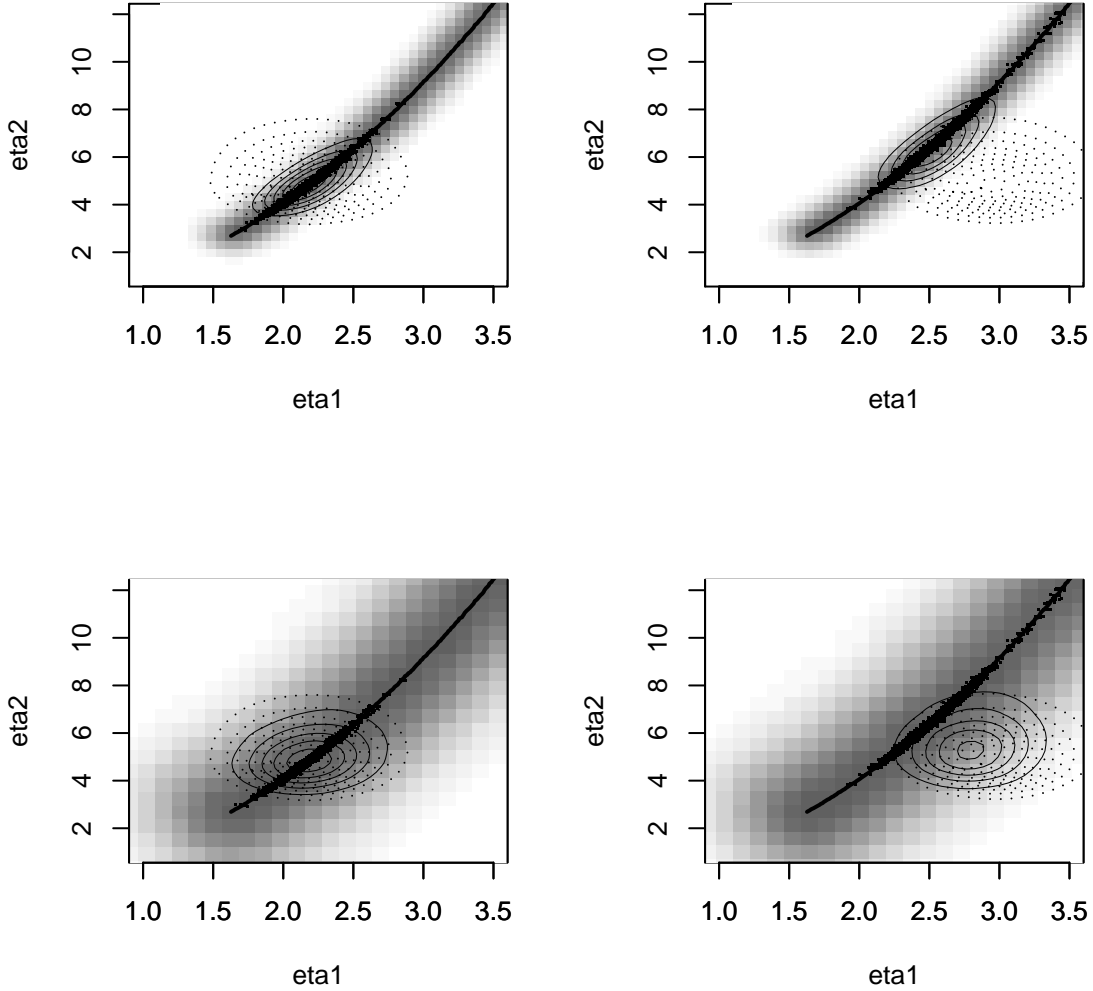


Figure 1. Priors and posteriors for the statistical model, tinker toy example.

The dotted lines are contours of the likelihood of the statistical model $f(y|x, \eta)$ of the tinker toy example. The line is the prior on η determined by the implied map $\eta = g(\theta)$ from the parameters θ of scientific model $p(y|x, \theta)$ to the parameters η of the statistical model. In the left panels the scientific model is true, in the right it is false. The thickness of the line is proportional to the posterior of η . The prior $\pi(\eta)$ can be relaxed as indicated by the shading. The lower panels are more relaxed than the upper. The solid contours show the posterior under the relaxed prior. Relaxation causes the contours to enlarge in all cases. When the scientific model is false, the posterior shifts in search of the likelihood.

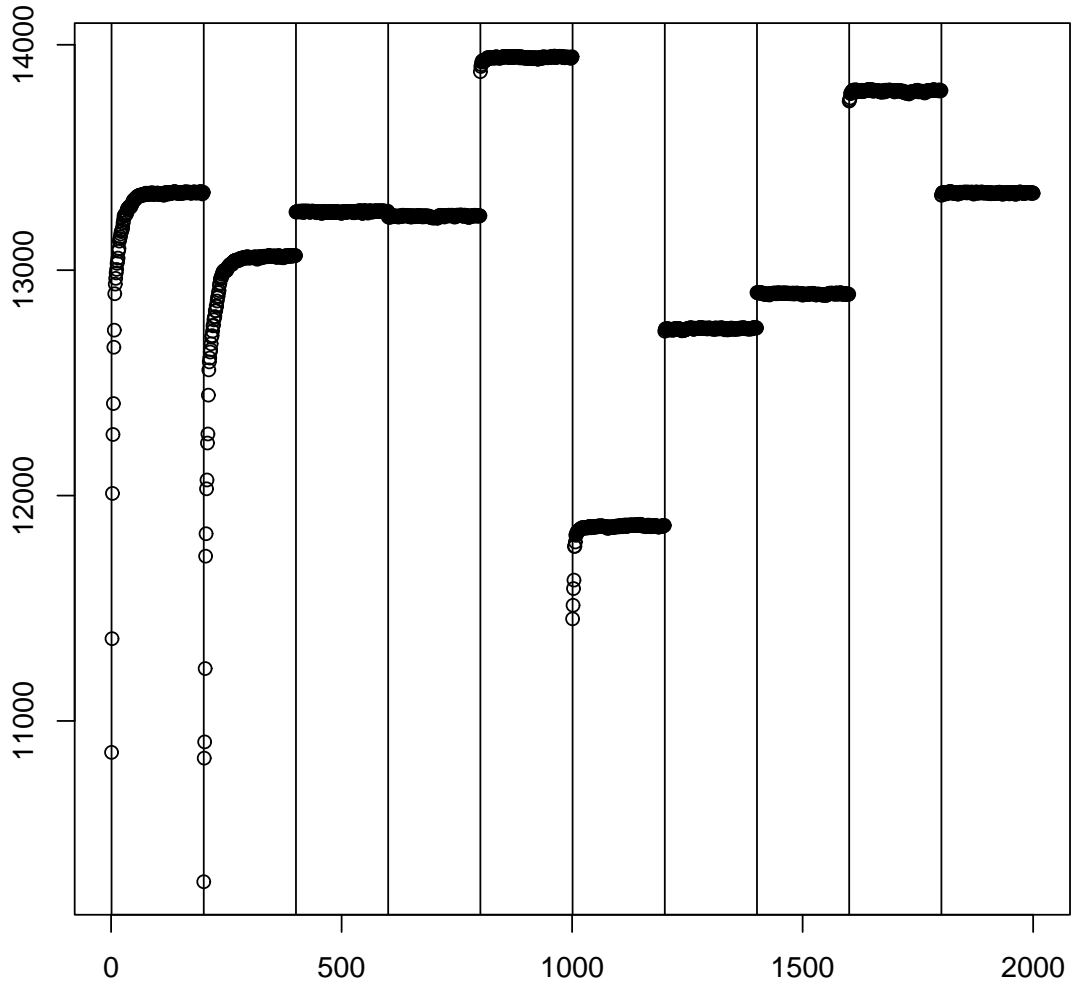


Figure 2. η Chain for the habit model. Ten successive runs of the η chain. For the habit model, the dimension of η is 16 and the dimension of θ is 7. Each run is 200 iterations. The log-likelihood of the simulated data set is plotted on the vertical axis. Vertical bars mark where θ changes. Jumps are because $\{\hat{y}_t\}_{t=1}^N$ changes at each vertical bar.

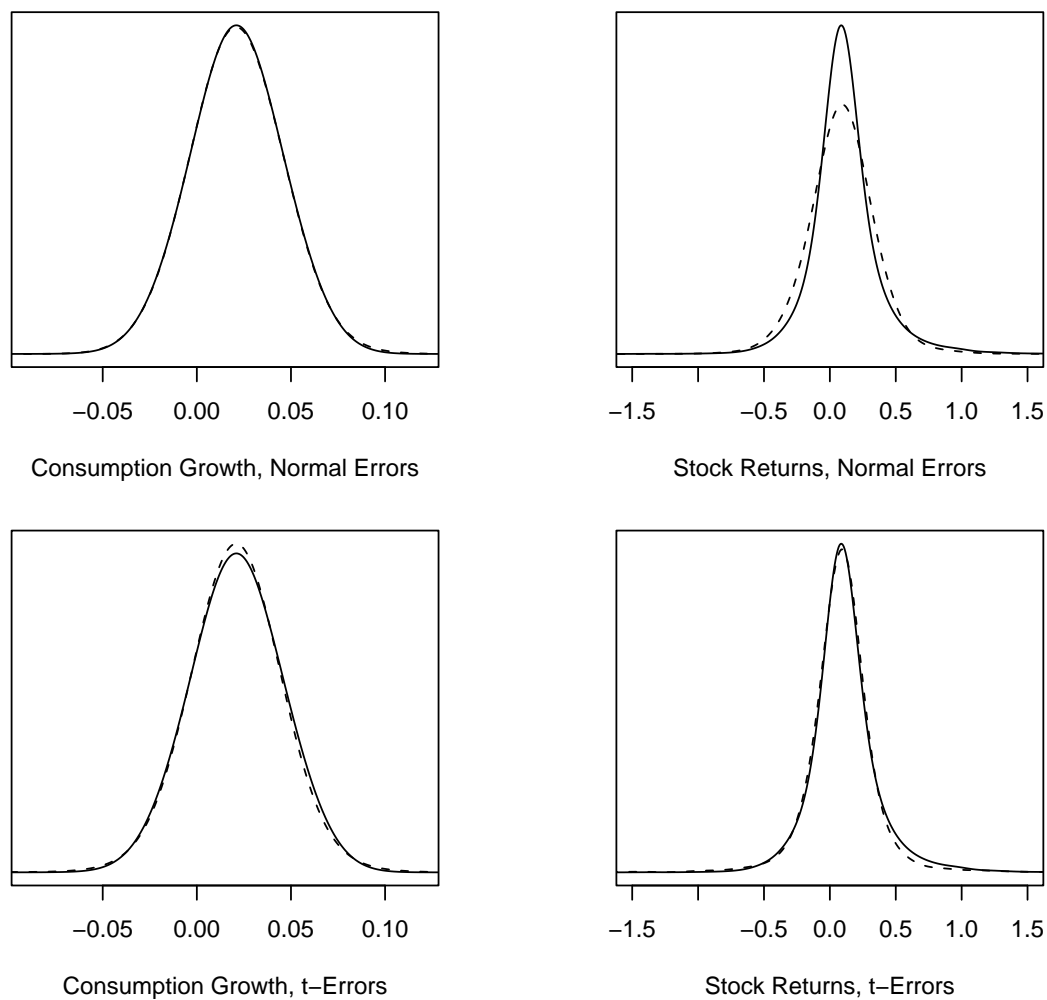


Figure 3. Statistical Model Error Distribution. The marginal densities of the scientific model at the parameter values in column Bayes-Garch Mode of Table 1 are plotted as the solid lines; they are kernel estimates from a simulation. In the upper panels, the marginal densities of the statistical model determined by this simulation are plotted as dashed lines. In the lower panels, the marginal densities of the statistical model with z_{2t} in (12) distributed as a t on 3.5 degrees of freedom instead of a normal are plotted as dashed lines.

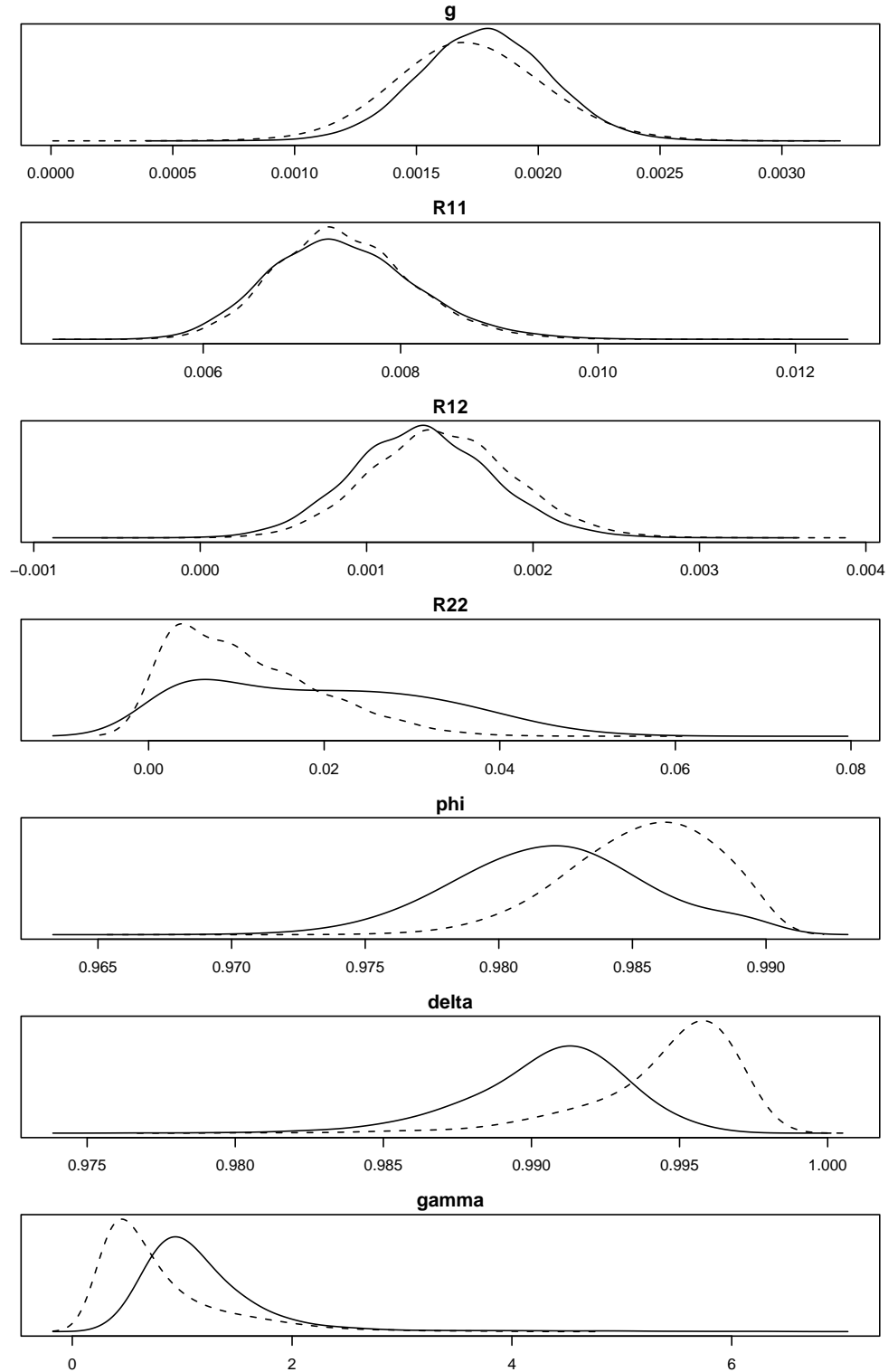


Figure 4. Density of the MCMC chain for θ . Shown is a kernel density estimate from iterates 1 to 800,000 by 8 of the MCMC chain for $\theta = (g, r_{11}, r_{12}, r_{22}, \psi, \delta, \gamma)$ at the monthly frequency. The dashed line is for the Bayes-VAR chain and the solid line is for the Bayes-GARCH chain.

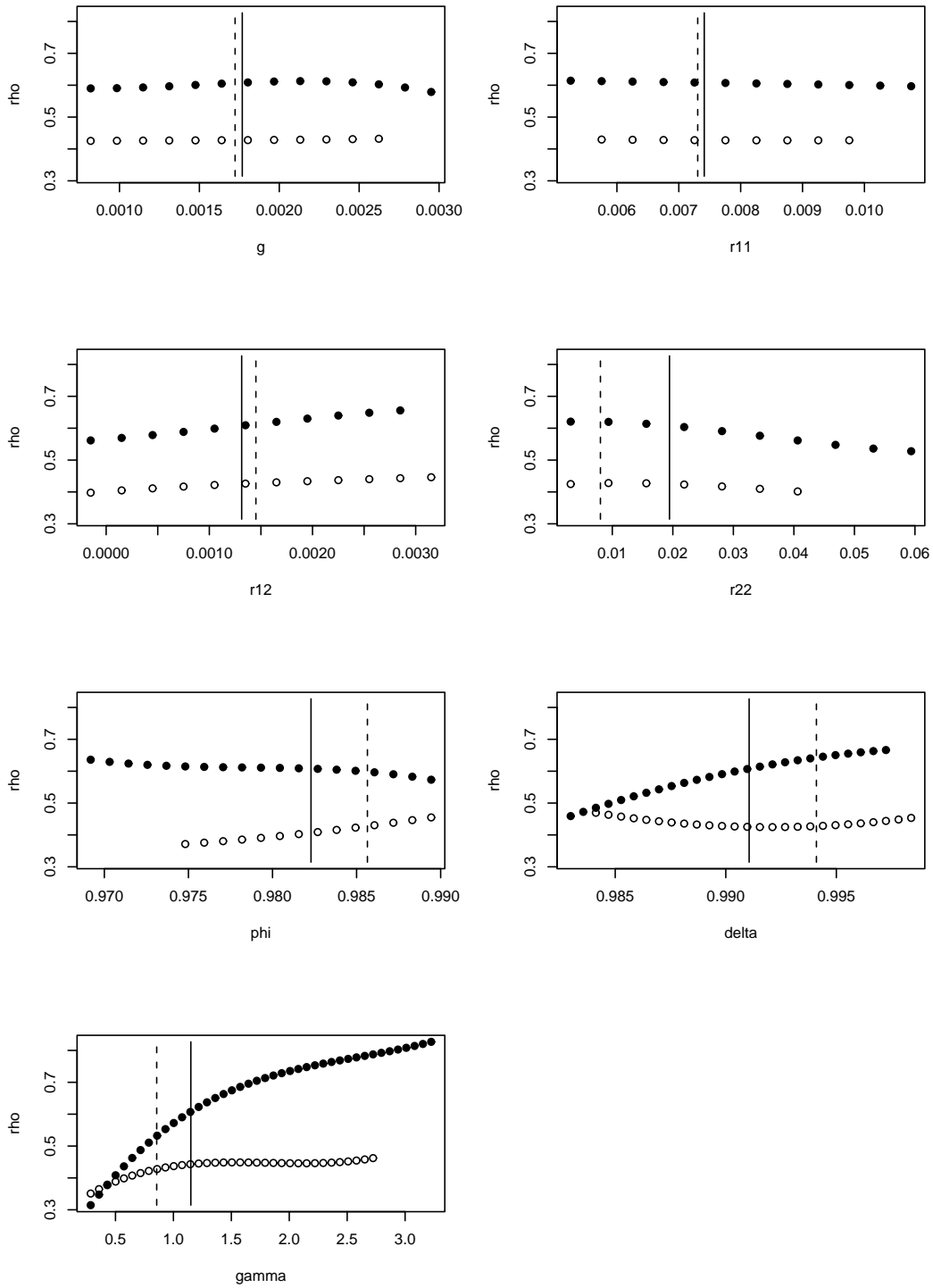


Figure 5. Conditional Correlation. The conditional correlation between annual consumption growth and stock returns for 2002 plotted against scientific model parameters $\theta = (g, r_{11}, r_{21}, r_{22}, \phi, \delta, \gamma)$ as determined from the map $g : \theta \mapsto \eta$. Dots are for the GARCH statistical model; circles for the VAR. The solid vertical line is the GARCH posterior mean and the dashed VAR.

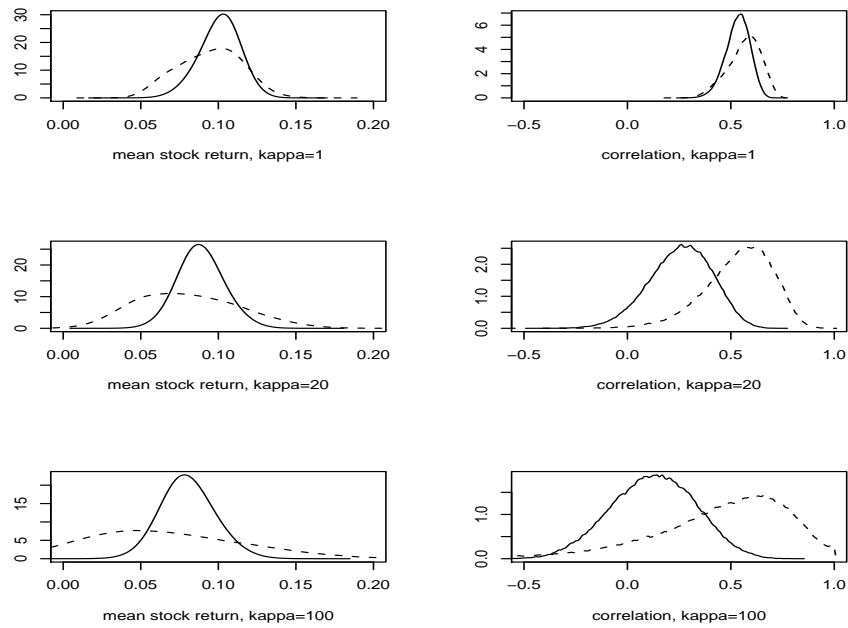


Figure 6. Priors and posterior for two functionals of the statistical model, habit model example. The statistical model is a bivariate GARCH model on annual consumption growth and stock returns. The scientific model is a habit persistence asset pricing model. The data covers 1929–2001. The left panels are the mean stock return for the period 2002–2102 implied by the GARCH model. The right panels are the conditional correlation between stock returns and consumption growth implied by the GARCH model for the year 2002. The prior is more relaxed in the lower panels than it is in the upper panels. The solid line is the posterior and the dashed line is the prior.