

Particle Learning for Fat-Tailed Distributions

Hedibert F. Lopes¹ and Nicholas G. Polson²

¹*INSPER Institute of Education and Research, São Paulo, Brazil*

²*University of Chicago Booth School of Business, Chicago, Illinois, USA*

It is well known that parameter estimates and forecasts are sensitive to assumptions about the tail behavior of the error distribution. In this article, we develop an approach to sequential inference that also simultaneously estimates the tail of the accompanying error distribution. Our simulation-based approach models errors with a t_ν -distribution and, as new data arrives, we sequentially compute the marginal posterior distribution of the tail thickness. Our method naturally incorporates fat-tailed error distributions and can be extended to other data features such as stochastic volatility. We show that the sequential Bayes factor provides an optimal test of fat-tails versus normality. We provide an empirical and theoretical analysis of the rate of learning of tail thickness under a default Jeffreys prior. We illustrate our sequential methodology on the British pound/U.S. dollar daily exchange rate data and on data from the 2008–2009 credit crisis using daily S&P500 returns. Our method naturally extends to multivariate and dynamic panel data.

Keywords Bayesian inference; Credit crisis; Dynamic panel data; Kullback–Leibler, MCMC.

JEL Classification C01; C11; C15; C16; C22; C58.

1. INTRODUCTION

Fat-tails are an important statistical property of time series prevalent in many fields, particularly economics and finance. Fat-tailed error distributions were initially introduced by Edgeworth (1888) and explored further by Jeffreys (1961) who once remarked that “... all data are t_4 .” They can be incorporated into dynamic models as latent variable scale mixtures of normals (Carlin et al., 1992). In this article, we develop a simulation-based sequential inference procedure for estimating the tail behavior of a time series using the t_ν -distribution. This family is attractive for this purpose due to its flexibility with normality ($\nu = \infty$) and Cauchy ($\nu = 1$) errors as special cases. Our method complements the existing literature by estimating the set of sequential posterior distributions $p(\nu|y')$

Q1

Address correspondence to Hedibert F. Lopes, INSPER Institute of Education and Research, Rua Quatá 300, Vila Olímpia, São Paulo, SP, 04546-042, Brazil; E-mail: hedibertFL@insper.edu.br

Color versions of one or more of the figures in the article can be found online at www.tandfonline.com/lecr.

44 for data $y^t = (y_1, \dots, y_t)$ and $t = 1, \dots, T$, as opposed to Markov chain Monte Carlo
 45 (MCMC) which estimates v given the full data history $p(v|y^T)$ (see Geweke, 1993; Eraker,
 46 Jacquier, and Polson (JPR), 1998; Jacquier et al., 2004; Fonseca et al., 2008). In other
 47 words, our methodology allows the researcher to estimate and update the tail-thickness
 48 of the error distribution as new data arrives.

49 The novel feature of our approach are the on-line estimates of the tail thickness of
 50 the error distribution using the marginal posterior distribution of the degrees of freedom
 51 parameter v . Being able to sequentially assess the degree of tail-heaviness is particularly
 52 important for dynamic portfolio and risk management strategies. For instance, $p(v|y^{t_0})$ and
 53 $p(\mu|y^{t_1})$ might resemble standard normal and t_4 distributions, respectively, for say t_0 much
 54 smaller than t_1 , which in turn would potentially affect decision making at both time points.

55 Our method is based on particle learning (PL, see Carvalho et al., 2010, and Lopes
 56 et al., 2010). We analyze two cases in detail: In the first observations y_t follow the
 57 independent and identically distributed (iid) standard t_v -distribution, i.e., $y_t \sim t_v(0, 1)$ (iid-
 58 t case), and in the second observations follow a non-identically distributed stochastic
 59 volatility model with fat-tails (SV- t case), i.e., $y_t|h_t \sim t_v(0, \exp\{h_t\})$ are conditionally
 60 independent given the T -dimensional latent vector of log-volatilities $h^T = (h_1, \dots, h_T)$, see
 61 JPR (2004) and Chib et al. (2002).

62 Our posterior distribution $p(v|y^t)$ on the tail thickness is sensitive to the choice of
 63 prior distribution, $p(v)$. We model the prior on the degrees of freedom v using a default
 64 Jeffreys prior (Fonseca et al., 2008). In this setting, we show that the Jeffreys prior
 65 has desirable properties. Primarily, it reduces bias for estimating the tail thickness in
 66 small sized data sets. Moreover, it is well known that more data helps to discriminate
 67 similar error distributions. Hence *a priori* we know that we will need a larger dataset to
 68 discriminate a t_{20} -distribution from a normal distribution than a t_4 -distribution from a
 69 normal. We develop a metric based on the asymptotic Kullback–Liebler rate of learning
 70 of tail thickness that can guide the amount of data required to discriminate two error
 71 distributions. Given the observed data, we then develop an empirical and theoretical
 72 analysis of the sequential Bayes factor which provides the optimal test of normality versus
 73 fat-tails in our sequential context.

74 Recent estimation approaches for fat-tails use approximate latent Gaussian models
 75 (McCausland, 2012). We use the traditional data augmentation with a vector of latent
 76 scale variables λ_t to avoid evaluating the likelihood (a T -dimensional integral). We
 77 develop a particle learning algorithm for sampling from the sequential set of joint
 78 posterior distributions $p(\lambda_t, v|y^t)$, for the iid- t case, and from $p(\lambda_t, h_t, v|y^t)$, for the SV- t
 79 case, for $t = 1, \dots, T$. The marginal posterior distribution $p(v|y^t)$ provides estimates of
 80 the tail-thickness of the error distribution. The purpose for developing new estimation
 81 methods is apparent from a remark of Smith (2000) who warns that the likelihood for
 82 non-Gaussian models can have several local maxima, be very skewed, or have modes
 83 on the boundary of the parameter space, making estimating tail behavior a complex
 84 statistical problem.
 85
 86

87 The rest of the article is outlined as follows. Section 2 describes how to sequentially
 88 learn the tail of the t_ν -distribution under iid- t and SV- t models. Section 3 discusses our
 89 particle learning implementation. We focus on using a default Jeffreys prior, showing that
 90 this has a number of desirable properties when learning the fat-tailed error distribution
 91 with finite samples. Section 4 provides an analysis of the sequential Bayes factor
 92 for testing normality versus fat-tails. Section 5 provides our empirical analysis and
 93 comparisons including an analysis of the British pound and U.S. dollar daily exchange
 94 rate and daily S&P500 returns from the credit crisis. Jacquier et al. (2004) apply MCMC
 95 methods to the SV- t model to daily exchange rate on the British pound versus the U.S.
 96 dollar, and we provide a sequential analysis for comparative purposes. Finally, Section 6
 97 concludes.
 98
 99

100 **2. T_ν -DISTRIBUTED ERRORS**

101 Consider data $y^t = (y_1, \dots, y_t)$ arising from a fat-tailed t_ν -distribution. The data are
 102 observed on-line, and we wish our estimation procedure to take this into account. Given
 103 a prior distribution $p(v)$, the aim is to compute a set of sequential marginal posterior
 104 distributions $p(v|y^t)$ which are given by Bayes rule
 105

$$106 \quad p(v|y^t) = \frac{p(y^t|v)p(v)}{\int p(y^t|v)p(v)dv}.$$

107
 108 The marginal likelihood is given by $p(y^t|v)$. In an iid setting, this likelihood is simply
 109 $p(y^t|v) = \prod_{i=1}^t p(y_i|v)$, a product of marginals. In the SV- t setting, it is more complicated
 110 and requires integrating out the unobserved t -dimensional vector of log-volatilities $h^t =$
 111 (h_1, \dots, h_t) , namely
 112

$$113 \quad p(y^t|v) = \int \prod_{i=1}^t p(y_i|h_i, v)p(h^t)dh^t,$$

114 where $p(y_i|h_i, v) \sim t_\nu(0, \exp\{h_i\})$. One advantage of particle methods is that this
 115 computation will naturally occur within the procedure. Our task is to provide sequential
 116 inference for the degrees of freedom or tail thickness parameter, ν , via the set of marginal
 117 posterior distributions $p(v|y^t)$, for $t = 1, \dots, T$. To do this, we will first use a standard
 118 data augmentation and then provide a sequential Monte Carlo algorithm to sample from
 119 $p(\lambda_t, \nu|y^t)$ which we now describe for the iid- t and SV- t models.
 120
 121
 122
 123
 124

125 **2.1. The iid- t Model**

126 Consider iid observations y_t , for $t = 1, \dots, T$, from a fat-tailed location-scale model
 127

$$128 \quad y_t = \mu + \sigma\eta_t \quad \text{where } \eta_t \stackrel{\text{iid}}{\sim} t_\nu(0, 1).$$

129

Data augmentation uses a scale mixture of normals representation by writing η_t in the following two steps: i) $\eta_t = \sqrt{\lambda_t}\epsilon_t$ and ii) $\lambda_t \stackrel{\text{iid}}{\sim} IG(v/2, v/2)$, where IG denotes the inverse gamma distribution. The marginal data distribution, integrating out λ_t , is then the fat-tailed t_v -distribution $p(y_t|v, \mu, \sigma^2) \sim t_v(\mu, \sigma^2)$, where σ^2 can be interpreted as a scale parameter. This leads to a hierarchical specification of the model

$$y_t = \mu + \sigma\sqrt{\lambda_t}\epsilon_t \quad \text{where } (\lambda_t|v) \stackrel{\text{iid}}{\sim} IG(v/2, v/2) \text{ and } \epsilon_t \stackrel{\text{iid}}{\sim} N(0, 1).$$

These specifications lead to a likelihood function $p(y|\mu, \sigma^2, v)$ of the form

$$p(y|\mu, \sigma^2, v) = \prod_{t=1}^T \frac{\Gamma(\frac{v+1}{2})}{\sqrt{v}\Gamma(\frac{v}{2})} \left[1 + \frac{1}{v} \left(\frac{y_t - \mu}{\sigma} \right)^2 \right]^{-\frac{v+1}{2}}$$

with marginal distribution $p(y_t|v) = \int p(y_t|v, \mu, \sigma^2)p(\mu, \sigma^2)d\mu d\sigma^2$. Fonseca et al. (2008) make the important observation that the marginal likelihood for v becomes unbounded as $v \rightarrow \infty$ and the maximum likelihood estimator is not well defined. This leads us to further develop an approach based on prior regularization, namely that the degree of freedom parameter v is random with a prior distribution $p(v)$ which we further discuss in Section 2.3.

Inference on the parameters (μ, σ^2) is not the focus of our study, and for simplicity we assume that either they are known quantities or taken from a standard diffuse prior, $p(\mu) \propto 1$, and inverse-gamma prior $\sigma^2 \sim IG(n_0/2, n_0\sigma_0^2/2)$ given hyper-parameters n_0 and σ_0^2 . These parameters control, respectively, the shape and the location of the distribution.

2.2. The SV- t Model

A common model of time-varying volatility is the stochastic volatility model with fat-tails (SV- t) for returns and volatility (see Lopes and Polson, 2010a, for a recent review). The basic SV model is specified by evolution dynamic

$$\begin{aligned} y_t &= \exp\{h_t/2\}\epsilon_t & \epsilon_t &\stackrel{\text{iid}}{\sim} N(0, 1), \\ h_t &= \alpha + \beta h_{t-1} + \tau u_t & u_t &\stackrel{\text{iid}}{\sim} N(0, 1). \end{aligned}$$

The fat-tailed SV- t is obtained by adding an extra random scale parameter λ_t and, as described in the conditionally iid setting, is equivalent to assuming that $\epsilon_t \sim t_v(0, 1)$ (see, for example, JPR, 2004). The model can then be expressed as

$$\begin{aligned} y_t &= \exp\{h_t/2\}\sqrt{\lambda_t}\epsilon_t & \epsilon_t &\stackrel{\text{iid}}{\sim} N(0, 1) \\ h_t &= \alpha + \beta h_{t-1} + \tau u_t & u_t &\stackrel{\text{iid}}{\sim} N(0, 1) \\ \lambda_t &\stackrel{\text{iid}}{\sim} IG(v/2, v/2). \end{aligned}$$

173 The parameter β is the persistence of the volatility process and τ^2 the volatility of the
 174 log-volatility. Estimation of these parameters will be greatly affected by the fat-tail error
 175 assumptions which in turn will affect predicting price and volatility (see, for example,
 176 Jacquier and Polson, 2000).

177 To complete the model specification, we need a prior distribution for the parameters
 178 (α, β, τ^2) given v . For simplicity, we take a conditionally conjugate normal-inverse-
 179 gamma-type prior. Specifically, $(\alpha, \beta) | \tau^2 \sim N(b_0, \tau^2 B_0)$ and $\tau^2 \sim IG(c_0, d_0)$, for known
 180 hyper-parameters b_0, B_0, c_0 , and d_0 . Lack of prior information is achieved when $B_0^{-1} \approx$
 181 0 and $c_0 \approx 0$. The marginal prior distribution for (α, β) is, therefore, a Student's t
 182 distribution. This conditionally conjugate structure will aid in the development of our
 183 particle learning algorithm as it leads to conditional sufficient statistics. Nonconjugate
 184 prior specifications can also be handled in our framework, see Lopes et al. (2010) for
 185 further discussion.

187
 188 **2.3. Priors on v**

189 In the models considered so far, an important modeling assumption is the regularization
 190 penalty $p(v)$ on the tail thickness. A default Jeffreys-style prior was developed by Fonseca
 191 et al. (2008) and, we will see, with a number of desirable properties—particularly when
 192 learning a fat-tail (e.g., a t_4 -distribution) from a finite dataset. The default Jeffreys prior
 193 for v takes the form

$$194 \quad 195 \quad 196 \quad 197 \quad p(v) = \frac{1}{\sigma} \left(\frac{v}{v+3} \right)^{1/2} \left\{ \psi' \left(\frac{v}{2} \right) - \psi' \left(\frac{v+1}{2} \right) - \frac{2(v+3)}{v(v+1)^2} \right\}^{1/2}, \quad (1)$$

198 where $\psi'(a) = d\{\psi(a)\}/da$ and $\psi(a) = d\{\log \Gamma(a)\}/da$ are the trigamma and digamma
 199 functions, respectively. The interesting feature of this prior is its behavior as v goes to
 200 infinity and it has polynomial tails of the form $p(v) \sim v^{-4}$. This is in contrast to commonly
 201 used priors such as Fernandez and Steel (1999) and Geweke (1993) who essentially specify
 202 priors with exponential tails of the form $v \exp\{-\lambda v\}$, for a subjectively chosen hyper-
 203 parameter, λ . In this case, the tail of the prior decays rather fast for large values of v and
 204 assessing the degree of tail thickness can require prohibitively large samples.

205 Table 1 compares Fonseca's robust prior to several exponential priors, including the
 206 exponential prior with mean 20 (rate $\lambda = 0.05$), which was advocated, for instance, by
 207 Geweke (1993). As it can be seen, despite its higher mass for heavy tailed distributions
 208 (small values of v), Fonseca's prior also places higher mass for normality (large values
 209 of v), when compared to the exponential with mean 20. The exponential priors, with
 210 high mean, essentially place zero mass on normality, whereas Fonseca's prior places
 211 approximately 0.01 probability on normality, which although being small can still be
 212 overwhelmed by an informative likelihood.

213 In our empirical analysis, we will show how this prior reduces bias in the posterior
 214 mean $E(v|y')$ and also how it helps discriminate a fat-tailed t_4 -distribution from
 215

TABLE 1
Fonseca's Prior and Geweke's Prior

v_0	10	20	50	150	190
$P_{\mathcal{E}}(v < v_0 \lambda = 0.01)$	0.01	0.20	0.45	0.8959	0.98180
$P_{\mathcal{E}}(v < v_0 \lambda = 0.05)$	0.36	0.61	0.91	0.9995	0.99997
$P_{\mathcal{E}}(v < v_0 \lambda = 0.20)$	0.85	0.98	1.00	1.0000	1.00000
$P_J(v < v_0)$	0.85	0.93	0.98	0.9972	0.99952

normality. On the other hand, the flat uniform prior suffers from placing too much mass on high values of v —which are close to normality—making the inference problem harder for finite samples.

3. PARTICLE LEARNING FOR FAT-TAILS

We now provide a discussion of particle learning with particular reference to estimating fat-tails. Sequential Bayesian computation requires calculation of a set of posterior distributions $p(v|y^t)$, for $t = 1, \dots, T$, where $y^t = (y_1, \dots, y_t)$.

Loosely speaking, particle learning is a sequential Monte Carlo scheme that sequentially learns a low dimensional vector of essential states, usually comprising a combination of a few latent states of the state-space model along with conditional sufficient statistics for fixed, time-invariant parameters. Section 3.1 provides a thorough explanation and implementation of particle learning for the iid- t case. See Carvalho et al. (2010), Lopes et al. (2010), Lopes and Tsay (2011), and Lopes and Carvalho (2013) for extended discussion and several examples of PL in action.

Central to PL is the creation of a *essential state vector* Z_t to be tracked sequentially. We assume that this vector is conditionally sufficient for the parameter of interest; so that $p(v|Z_t)$ is either available in closed-form or can easily be sampled from. More precisely, given samples $\{Z_t^{(i)}\}_{i=1}^N \sim p(Z_t|y^t)$ and a Rao–Blackwellized identity, then a simple mixture approximation to the set of posteriors is given by

$$p^N(v|y^t) = \frac{1}{N} \sum_{i=1}^N p(v|Z_t^{(i)}).$$

Here the conditional posterior $p(v|Z_t^{(i)})$ will include the dependence on σ^2 for the iid- t case and (α, β, τ^2) and the latent volatilities $h^t = (h_1, \dots, h_t)$ for the SV- t case through the essential state vector.

The task of sequential Bayesian computation is then equivalent to a filtering problem for the essential state vector, drawing $\{Z_t^{(i)}\}_{i=1}^N \sim p(Z_t|y^t)$ sequentially from the set of

posterior. To this end, PL exploits the following sequential decomposition of Bayes' rule

$$\begin{aligned}
 p(Z_{t+1}|y^{t+1}) &= \int p(Z_{t+1}|Z_t, y_{t+1}) d\mathbb{P}(Z_t|y^{t+1}) \\
 &\propto \int \underbrace{p(Z_{t+1}|Z_t, y_{t+1})}_{\text{propagate}} \underbrace{p(y_{t+1}|Z_t)}_{\text{resample}} d\mathbb{P}(Z_t|y^t).
 \end{aligned}$$

The distribution $d\mathbb{P}(Z_t|y^{t+1}) \propto p(y_{t+1}|Z_t)d\mathbb{P}(Z_t|y^t)$ is a 1-step smoothing distribution. Here $\mathbb{P}(Z_t|y^t)$ denotes the current distribution of the current state vector and in particle form corresponds to $\frac{1}{N} \sum_{i=1}^N \delta_{Z_t^{(i)}}$, with δ a Dirac measure.

Bayes rule above then gives us a prescription for constructing a sequential simulation-based algorithm: given $\mathbb{P}(Z_t|y^t)$, find the smoothed distribution $\mathbb{P}(Z_t|y^{t+1})$ via resampling and then propagate forward using $p(Z_{t+1}|Z_t, y_{t+1})$. This simply finds draws from the next filtering distribution $\mathbb{P}(Z_{t+1}|y^{t+1})$. Parameter inference is then achieved offline using $p(\theta|Z_{t+1})$.

From a sampling perspective, this leads to a very simple algorithm for updating particles $\{Z_t\}_{i=1}^N$ to $\{Z_{t+1}\}_{i=1}^N$ in the following three steps:

1. *Resample*: with replacement from a multinomial with weights proportional to the predictive distribution $p(y_{t+1}|Z_t^{(i)})$ to obtain $\{Z_t^{\zeta(i)}\}_{i=1}^N$;
2. *Propagate*: with $Z_{t+1}^{(i)} \sim p(Z_{t+1}|Z_t^{\zeta(i)}, y_{t+1})$ to obtain $\{Z_{t+1}^{(i)}\}_{i=1}^N$;
3. *Learning*: v from $p(v|Z_{t+1})$.

The ingredients of particle learning are the essential state vector Z_t , a predictive probability rule $p(y_{t+1}|Z_t^{(i)})$ for resampling $\zeta(i)$, and a propagation rule to update particles: $Z_t^{\zeta(i)} \rightarrow Z_{t+1}^{(i)}$. The essential state vector will include the necessary conditional sufficient statistics for parameter learning given a model specification.

3.1. PL for the iid- t Case

First, we consider the normal location-scale model of Section 2.1 with $\mu = 0$ for simplicity. The model corresponds to a data augmentation scheme $(y_t|\sigma^2, \lambda_t) \sim N(0, \sigma^2 \lambda_t)$ with $(\lambda_t|v) \sim IG(v/2, v/2)$. To complete the model, we assume priors of the form $\sigma^2 \sim IG(n_0/2, n_0 \sigma_0^2/2)$ and Jeffreys prior $p(v)$ for v (Eq. 1).

Now, the key to our approach is the use of an essential state vector Z_t . The algorithm requires the following distributions: $p(y_{t+1}|Z_t)$, $p(v, \sigma^2|Z_t)$, and $p(\lambda_t|\sigma^2, v, y_t)$. Bayes rule yields

$$p(v|\lambda^t) \equiv p(v|Z_{t1}, Z_{t2}) \propto p(v) \left(\frac{(\frac{v}{2})^{\frac{v}{2}}}{\Gamma(\frac{v}{2})} \right)^t Z_{t1}^{-(v/2+1)} \exp\{-v Z_{t2}/2\} \quad (2)$$

and

$$p(\sigma^2|y', \lambda') \equiv p(\sigma^2|Z_{i3}, Z_{i4}) \sim IG(Z_{i3}/2, Z_{i4}/2) \quad (3)$$

with recursive updates for the parameter sufficient statistics

$$\begin{aligned} Z_{i1} &= Z_{i-1,1}\lambda_t & \text{and} & & Z_{i2} &= Z_{i-1,2} + 1/\lambda_t, \\ Z_{i3} &= Z_{i-1,3} + 1 & \text{and} & & Z_{i4} &= Z_{i-1,4} + y_t^2/\lambda_t, \end{aligned}$$

with initial values $Z_{01} = 1$, $Z_{02} = 0$, $Z_{03} = n_0$, and $Z_{04} = n_0\sigma_0^2$.

Additionally, the predictive distribution for resampling and the latent state conditional posterior for propagation are directly available as

$$p(y_{t+1}|\lambda_{t+1}, Z_t) \sim t_{Z_{i3}+2}\left(0, \frac{Z_{i4}}{Z_{i3} + 2}\lambda_{t+1}\right), \quad (4)$$

$$p(\lambda_t|\sigma^2, v, y_t) \sim IG\left(\frac{v+1}{2}, \frac{v + y_t^2/\sigma^2}{2}\right). \quad (5)$$

Therefore, we use an essential state vector given by $Z_t = (\lambda_{t+1}, Z_{t1}, Z_{t2}, Z_{t3}, Z_{t4})$. We are now ready to outline the steps of the PL scheme (see Panel A).

When $\mu \neq 0$ and a conditionally conjugate prior for location μ is used, say $N(\mu_0, \sigma^2 C_0)$, it follows that Eq. (3) is replaced by $p(\sigma^2|y', \lambda')p(\mu|\sigma^2, y', \lambda')$, while the vector Z_t is expanded accordingly. If instead the prior for μ is $N(\mu_0, C_0)$, independent of σ^2 , then the essential vector Z_t would include one of the two parameters, most likely μ since in practice location parameters are easier to update.

3.2. PL for the SV- t Case

Particle learning for the SV- t model is similar to the iid- t model despite being somewhat more elaborated with the latent state now being the scale mixture λ_t as well as the log-volatilities h_t . In addition, there are three parameters (α, β, τ^2) driving the log-volatility dynamic behavior, as opposed to σ^2 in the iid- t model.

Static Parameters. Let us first deal with $\theta = (\alpha, \beta, \tau^2)$ the vector of fixed parameters driving the log-volatility equation (see Section 2.2). Conditional on the latent volatilities $h' = (h_1, \dots, h_t)$, sampling θ is rather straightforward since it is based on the conjugate Bayesian analysis of the normal linear regression with $x'_t = (1, h_{t-1})$ (Gamerman and Lopes, 2006, Chapter 2), i.e., $(\alpha, \beta|\tau^2) \sim N(b_t, \tau^2 B_t)$ and $\tau^2 \sim IG(c_t, d_t)$. The parameter sufficient statistics are $Z_t^\theta = (b_t, B_t, c_t, d_t)$, and they can be determined recursively as

$$\begin{aligned} B_t^{-1}b_t &= B_{t-1}^{-1}b_{t-1} + h_t x_t, \\ B_t^{-1} &= B_{t-1}^{-1} + x_t x_t', \\ c_t &= c_{t-1} + 1/2, \\ d_t &= d_{t-1} + (h_t - b_t' x_t)h_t/2 + (b_{t-1} - b_t)' B_{t-1}^{-1} b_{t-1}/2. \end{aligned} \quad (6)$$

Start at time $t = 0$ with particle set $\{(\nu, \sigma^2, Z_{01}, Z_{02}, Z_{03}, Z_{04})^{(i)}\}_{i=1}^N$.

Step 1. For $i = 1, \dots, N$,

- Sample $\lambda_{t+1}^{(i)} \sim IG(\nu^{(i)}/2, \nu^{(i)}/2)$,
- Set $Z_t^{(i)} = (\lambda_{t+1}, Z_{t1}, Z_{t2}, Z_{t3}, Z_{t4})^{(i)}$.

Step 2. Resample particles $\{(\tilde{\nu}, \tilde{\sigma}^2, \tilde{Z}_{t1}, \tilde{Z}_{t2}, \tilde{Z}_{t3}, \tilde{Z}_{t4})^{(i)}\}_{i=1}^N$ with weights proportional to $p(y_{t+1}|Z_t^{(i)})$ (equation 4),

Step 3. For $i = 1, \dots, N$,

- Sample $\lambda_{t+1}^{(i)} \sim p(\lambda_{t+1}|\tilde{\sigma}^{2(i)}, \tilde{\nu}^{(i)}, y_{t+1})$ (equation 5),
- Update the essential state vector:

$$\begin{aligned} Z_{t+1,1} &= \tilde{Z}_{t1}^{(i)} \lambda_{t+1}^{(i)} & \text{and} & & Z_{t+1,2} &= \tilde{Z}_{t2}^{(i)} + 1/\lambda_{t+1}^{(i)} \\ Z_{t+1,3} &= \tilde{Z}_{t3}^{(i)} + 1 & \text{and} & & Z_{t+1,4} &= \tilde{Z}_{t4}^{(i)} + y_{t+1}^2/\lambda_{t+1}^{(i)} \end{aligned}$$
- Sample $\nu^{(i)} \sim p(\nu|Z_{t+1}^{(i)})$ (equation 2),
- Sample $\sigma^{2(i)} \sim p(\sigma^2|Z_{t+1}^{(i)})$ (equation 3).

Set $t = t + 1$ and return to step 1.

PANEL A Particle learning for the iid- t model.

Q13

Resampling Step. To sequentially resample the log-volatility h_t and propagate a new volatility state h_{t+1} , we use the Kim, Shephard, and Chib (1998) strategy of approximating the distribution of $\log \tilde{y}_t^2$, where $\tilde{y}_t^2 = y_t^2/\lambda_t$, by a carefully tuned seven-component mixture of normals¹. Then, a standard data augmentation argument allows the mixture of normals to be conditionally transformed in individual normals, i.e., $(\varepsilon_t|k_t) \sim N(\mu_{k_t}, v_{k_t}^2)$, such that $k_t \sim \text{Mult}(\pi)$. Conditionally on k^t , the SV- t model for $z_{k_t} = \log y_t^2 - \log \lambda_t - \mu_{k_t}$ can be rewritten as a standard first order dynamic linear model, i.e.,

Q3

$$\begin{aligned} (z_{k_t}|h_t, \lambda_t, k_t) &\sim N(h_t, v_{k_t}^2), \\ (h_t|h_{t-1}, \theta) &\sim N(\alpha + \beta h_{t-1}, \tau^2), \end{aligned}$$

with conditional state sufficient statistics $Z_t^h = (m_t, C_t)$ given by the standard Kalman recursions (West and Harrison, 1997). More explicitly, the conditional posterior

Q4

¹More precisely, $\log \tilde{y}_t^2 = h_t + \varepsilon_t$, where $\varepsilon_t = \log \varepsilon_t^2$ follows a $\log \chi_1^2$ distribution, a parameter-free left skewed distribution with mean -1.27 and variance 4.94 . They show that the $\log \chi_1^2$ can be well approximated by $\sum_{j=1}^7 \pi_j N(\mu_j, v_j^2)$, where $\pi = (0.0073, 0.1056, 0.00002, 0.044, 0.34, 0.2457, 0.2575)$, $\mu = (-11.4, -5.24, -9.84, 1.51, -0.65, 0.53, -2.36)$, and $v^2 = (5.8, 2.61, 5.18, 0.17, 0.64, 0.34, 1.26)$.

388 $(h_t|Z_t^h, \theta) \sim N(m_t, C_t)$ with moments given by

389
390
$$m_t = (1 - A_t)a_t + A_t z_{k_t} \quad \text{and} \quad C_t = (1 - A_t)R_t, \quad (7)$$

391
392 where $a_t = (\alpha + \beta m_{t-1})$, $A_t = R_t/Q_t$, $R_t = \beta^2 C_{t-1} + \tau^2$ and $Q_t = R_t + v_{k_t}^2$.

393 **Essential State Vector.** We will take advantage of the above Kalman recursions in the
394 resampling step. We use an essential state vector of the form

395
396
$$Z_t = (\lambda_{t+1}, Z_t^0, Z_t^v, Z_t^h),$$

397
398 where the subset (Z_t^0, Z_t^v) of Z_t is essentially the set (Z_{t1}, \dots, Z_{t4}) derived from the iid- t
399 model.

400 There are many efficiencies to be gained with this approach over traditional SMC
401 approaches. For example, we only need to sample h_{t-1} and h_t (Step 2) in order to
402 propagate Z_t^0 and sample θ (Step 4). In other words, PL does not necessarily need to
403 keep track of the log-volatilities. For instance, point-wise evaluations of $p(h_t|y^t)$ can be
404 approximated by the Monte Carlo average of the Kalman filter densities, i.e., $p^N(h_t|y^t) =$
405 $\frac{1}{N} \sum_{i=1}^N p(h_t; m_t^{(i)}, C_t^{(i)})$.

406 For estimation of the fat-tails, we can use a Rao–Blackwellized density estimate. For
407 example in the SV- t case, in order to reduce Monte Carlo error, we use an estimate of
408 the form

409
410
$$p(v|y^t) = \mathbb{E} \{p(v|\lambda_t, h_t, y^t)\} \approx \frac{1}{N} \sum_{i=1}^N p(v|(\lambda^t, h^t)^{(i)}, y^t),$$

411
412 where $\{(\lambda^t, h^t)^{(i)}\}_{i=1}^N$ are draws from $p(\lambda^t, h^t|y^t)$. This leads to efficiency gains as the
413 conditional $p(v|\lambda^t, h^t, y^t)$ and conditional mean $\mathbb{E}(v|\lambda^t, h^t, y^t)$ are known in closed form.
414 We are now ready to outline the steps of the PL scheme for the SV- t model (see Panel B).

415 **PL and MCMC.** Although direct comparison with MCMC (Verdinelli and Wasserman,
416 1991) is not the focus of this article, we observe that MCMC is inherently a nonsequential
417 procedure. MCMC provides the full joint distribution $p(h^T, \theta, v|y^T)$ including smoothing
418 of the initial volatility states particle learning only computes $p(h_T, \theta|y^T)$ —the distribution
419 of the final state h_T and parameters θ . Another difference is in the assessment of MC
420 errors. MCMC generates a dependent sequence of draws, PL has standard \sqrt{N} MC
421 bounds, but can suffer from accumulation of MC error for larger T . MCMC for learning
422 fat-tails v can exhibit low conductance (Eraker et al., 1998), having difficulty escaping
423 lower values of v in the chain, and can lead to poor convergence. Computationally
424 speaking, the cost of performing MCMC sequentially is prohibitive high when compared
425 to PL. Carvalho et al. (2010, Example 7) compares MCMC to PL in the simple first order
426 normal dynamic linear model and show that, for $T = 1,000$ time periods and $N = 500$
427 particles, PL is roughly one order of magnitude faster than MCMC.
428
429
430

Q5

431 **Step 0.** Sample $\lambda_t^{(i)} \sim IG(\nu^{(i)}/2, \nu^{(i)}/2)$,

432

433 **Step 1.** Resample particles $\{(\tilde{Z}_{t-1}^\theta, \tilde{Z}_{t-1}^h, \tilde{\lambda}_t, \tilde{\theta})\}_{i=1}^N$ with weights

434

435
$$w_t^{(i)} \propto \sum_{k_t=1}^7 \pi_i p_N(z_{k_t}^{(i)}; a_t^{(i)}, Q_t^{(i)}),$$

436

437

438 **Step 2.** Sample (h_{t-1}, h_t) from $p(h_{t-1}, h_t | Z_{t-1}^h, \lambda_t, \theta, y^t)$:

439

440 **Step 2.1.** Sample h_{t-1} from $\sum_{j=1}^7 \pi_j f_N(h_{t-1}; \hat{h}_{t-1,j}, V_{t-1,j})$, where

441
$$\hat{h}_{t-1,j} = V_{t-1,i}(m_{t-1}/C_{t-1} + z_{ti}\beta/(v_i^2 + \tau^2))$$

442
$$V_{t-1,j} = 1/(1/C_{t-1} + \beta^2/(v_j^2 + \tau^2))$$

443

444 for $z_{ti} = \log y_t^2 - \log \lambda_t - \mu_i - \alpha$,

445

446 **Step 2.2.** Sample h_t from $\sum_{j=1}^7 \pi_j f_N(h_t; \tilde{h}_{tj}, W_{tj})$, where

447
$$\tilde{h}_{ti} = W_{ti}(\tilde{z}_{ti}/v_i^2 + (\alpha + \beta h_{t-1})/\tau^2)$$

448
$$W_{ti} = 1/(1/v_i^2 + 1/\tau^2)$$

449

450 for $\tilde{z}_{ti} = \log y_t^2 - \log \lambda_t - \mu_i$,

451

452 **Step 3.** Update $Z_{t+1}^{\nu^{(i)}}$ (equation 4); sample $\nu^{(i)} \sim p(\nu | Z_{t+1}^{\nu^{(i)}})$ (equation 2),

453

454 **Step 4.** Update $Z_t^{\theta^{(i)}}$ (equation 6); sample $\theta \sim p(\theta | Z_t^{\theta^{(i)}})$,

455

456 **Step 5.** Propagate $Z_t^{h^{(i)}}$ (equation 7).

PANEL B Particle learning for the SV- t model

460 4. MODEL ASSESSMENT WITH A SEQUENTIAL BAYES FACTOR

461 Sequential model determination is performed using a Bayes factor \mathcal{B}_T (Jeffreys, 1961;

462 West, 1984). This naturally extends to a sequential version for an infinite sequence of

463 (dependent) data we will still identify the “true” model. A probabilistic approach for

464 determining how quickly you can learn the tail of the error distribution is to use the

465 recursion

466

467

468
$$\mathcal{B}_{T+1} = \frac{p(y_{T+1}|y_1, \dots, y_T)}{q(y_{T+1}|y_1, \dots, y_T)} \mathcal{B}_T.$$

469

470 Blackwell and Dubins (1962) provide a general discussion of the merging of opinions

471 under Bayesian learning. They show that for any two models $p(y_1, \dots, y_T)$ and

472 $q(y_1, \dots, y_T)$ that are absolutely continuous with respect to each other, opinions that

473

474 merge in the following sense. First, \mathcal{B}_T is a martingale, \mathcal{F}_T -measurable and under the true
 475 model Q ,

$$476 \mathbb{E}_Q \left(\frac{p(y_{T+1}|y_1, \dots, y_T)}{q(y_{T+1}|y_1, \dots, y_T)} \middle| \mathcal{F}_T \right) = 1 \text{ so that } \mathbb{E}(\mathcal{B}_{T+1}|\mathcal{F}_T) = \mathcal{B}_T,$$

477 where \mathcal{F}_T represents all the information up to time t .
 480

481 By the martingale convergence theorem, $\mathcal{B}_\infty = \lim_{T \rightarrow \infty} \mathcal{B}_T$ exists almost surely under
 482 Q and in fact $\mathcal{B}_\infty = 0$ a.s. Q . Put simply, the sequential Bayes factor will correctly identify
 483 the “true” model Q under quite general data sequences include the SV- t model we
 484 consider here in detail. Furthermore, by the Shannon–McMillan–Breiman theorem (see,
 485 for example, Cover and Thomas, 2006), we can analyze the rate of learning via the
 486 quantity

$$487 \lim_{T \rightarrow \infty} \frac{1}{T} \ln q(y_1, \dots, y_T) \rightarrow H \quad \text{a.s. } Q,$$

490 where H is the entropy rate defined by $H = \lim_{T \rightarrow \infty} \mathbb{E}_Q(-\ln p(y_{T+1}|y_1, \dots, y_T)) < 0$.
 491 Hence as $H \in [-\infty, 0)$, we have that $\mathcal{B}_\infty = 0$. A similar result for the marginal likelihood
 492 ratio

$$493 \lim_{T \rightarrow \infty} \frac{1}{T} \ln \frac{p(y_1, \dots, y_T)}{q(y_1, \dots, y_T)} \rightarrow \lim_{k \rightarrow \infty} \mathbb{E}_Q \left(\ln \frac{p(y_{k+1}|y_k, \dots, y_1)}{q(y_{k+1}|y_k, \dots, y_1)} \right) < 0 \quad \text{a.s. } Q.$$

494 We will use this in the next subsection.

495 Bayes factors have a number of attractive features as they can be converted into
 496 posterior model probabilities when the model set is exhaustive. Lopes and Tobias (2011)
 497 provide a recent survey including computational strategies based on the Savage–Dickey
 500 density ratio. These results are only asymptotic, and with a finite amount of data, it helps
 501 to analyze the rate of learning using a Kullback–Leibler metric.

502 4.1. Discriminating a t_4 from a Gaussian

503 We can use these theoretical insights (see also Edwards et al., 1963; Lindley, 1956) to
 504 address the question *a priori* of “how long a time series one would have to observe
 505 to have strong evidence of a t_4 versus a Gaussian?” Jeffreys observed that one needs
 506 data sequences of length $T = 500$ to be able to discriminate the tails of an underlying
 507 probability distribution. We now formalize this argument using our sequential Bayes
 508 factor. One is motivated to define *a priori* the “expected” log-Bayes factor for a given
 509 data length, $\overline{\mathcal{B}}_T$, under the Gaussian model

$$510 \frac{1}{T} \ln \overline{\mathcal{B}}_T = \mathbb{E}_{t_\infty} \ln \frac{t_v}{t_\infty} = KL(t_v, t_\infty)$$

Q7

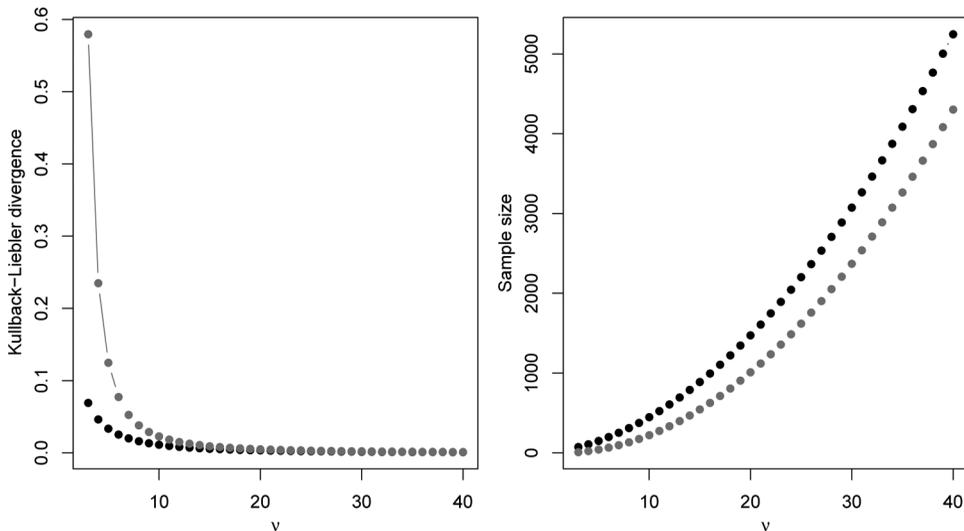
517 under the Gaussian t_∞ -model where KL denotes Kullback–Leibler divergence. Then, a
 518 *priori*, if we are given a level of Bayes factor discrimination $\overline{\mathcal{BF}}_T$, we then have to observe
 519 on average T^* observation to be able to discriminate the two models where

$$520 \quad T^* = \frac{1}{KL(t_\nu, t_\infty)} \ln \overline{\mathcal{BF}}_T.$$

521 This measure is asymmetric, as if the data is generated by a t_ν distribution, the constant
 522 changes to $KL(t_\infty, t_\nu)$. Q8

523 To illustrate the magnitudes of these effects, if we take $\nu = 3$ and $\mathcal{B} = 10$
 524 (strong evidence), for example, this argument would suggest that on average $T = 150$
 525 observations from a standard normal are needed to strongly reject the t_3 model, and
 526 on average $T = 20$ observations from the t_3 to strongly reject the standard normal
 527 distribution. This is borne out in our empirical study. Figure 1 plots the first factor in the
 528 above expression, namely, the Kullback–Leibler divergence between the t_ν -family and the
 529 Gaussian.

530 This also confirms the analysis in Gramacy and Pantaleo (2010). In a multivariate
 531 regression setting, they perform a Monte Carlo experiment where T and ν varies with $T \in$
 532 $\{30, 75, 100, 200, 500, 1000\}$ and $\nu \in \{3, 5, 7, 10, \infty\}$. They observed the frequency of time
 533 the \mathcal{B} indicated *strong* preference ($\mathcal{B} > 10$) for a model. Under normal errors, $\nu = 3$ could
 534 be determined with high accuracy for $T \leq 200$, $\nu = 5$ took $T \leq 1,000$, and for $10 \leq \nu < \infty$
 535
 536
 537
 538
 539



557 FIGURE 1 *i.i.d. model*. Discriminating a t_ν from a Gaussian. $KL(t_\nu, t_\infty)$ (black) and $KL(t_\infty, t_\nu)$ (grey).
 558 For $\nu = 4, 10, 20$, theoretical sample sizes are $T^* = 108, 446, 1,473$ for strong evidence against normality and
 559 $T^* = 22, 220, 1,009$ for strong evidence against t_ν .

560 very large samples would be required to discriminate the tails with any degree of posterior
 561 accuracy. Of course, for a given dataset, the Bayes factor might provide strong evidence
 562 even for small samples. The Jeffreys prior then has the nice property (by definition of
 563 the inverse of the Fisher information matrix) of down-weighting these regions of the
 564 parameter space where it is hard to learn the parameters.

565 It is also interesting to address the asymptotic behavior of the fat-tailed posterior
 566 distribution when the true model is not in the set of models under consideration.
 567 Berk (1966, 1970) assumes that the data generating process comes from $y_t \sim q(y)$ —
 568 a model outside our current consideration. Given our fat-tailed model $p(y|\theta, \nu)$, Berk
 569 shows that under mild regularity conditions the posterior distribution $p(\theta, \nu|y)$ will
 570 asymptotically concentrate with probability one on the subset of parameter values where
 571 the Kullback–Leibler divergence between $p(y|\theta, \nu)$ and $q(y)$ is minimized or equivalently
 572 $\int \log p(y|\theta, \nu)q(y)dy$ is maximized.
 573

574 5. EMPIRICAL RESULTS

575 We now illustrate our methodology for iid SV-Student's t error distributions (see Sections
 576 2.1 and 2.2 for the specifications). The iid- t model illustration will serve the additional and
 577 important purpose of showing that the uniform prior is not necessarily always a harmless
 578 prior. The SV- t model will be estimated sequentially on the British pound/U.S. dollar
 579 daily exchange rate series and daily returns on the S&P500 from a period in 2007–2010
 580 that includes the credit crisis. Resulting inferences will be compared with MCMC at the
 581 end of the sample.
 582

583 5.1. The iid- t Model

584 To illustrate the efficiency of our approach, we simulate a sample of size $T = 200$ from
 585 a Student's t_4 distribution, centered at zero and unit scale, i.e., $\sigma^2 = 1$. Figures 2 and 3
 586 show the joint posterior distributions of $p(\sigma^2, \nu|y^t)$ for $t = 50, 100, 150,$ and 200 under,
 587 respectively, the uniform prior and the Jeffreys prior of Fonseca et al. (2008). As the
 588 model implies that $\text{Var}(y_t) = \sigma^2\nu/(\nu - 2)$, one should not be too surprised that there is a
 589 posterior correlation between σ^2 and ν for small values of ν .

590 It is clear that the posterior provides fairly accurate sequential estimates for the joint as
 591 well as the marginal distributions (the exact posterior probabilities are computed on a fine
 592 bivariate grid). On the one hand, the Jeffreys prior, as anticipated, penalizes larger values
 593 of ν with the penalization slightly decreasing as the sample size increases. On the other
 594 hand, the uniform prior is impartial with respect to the number of degrees of freedom,
 595 so any information regarding ν comes exclusively from the likelihood which, in turn, is
 596 fairly uninformative about ν for $t = 50, 100,$ and 150 . Even when $t = 200$, there is still no
 597
 598
 599
 600
 601
 602

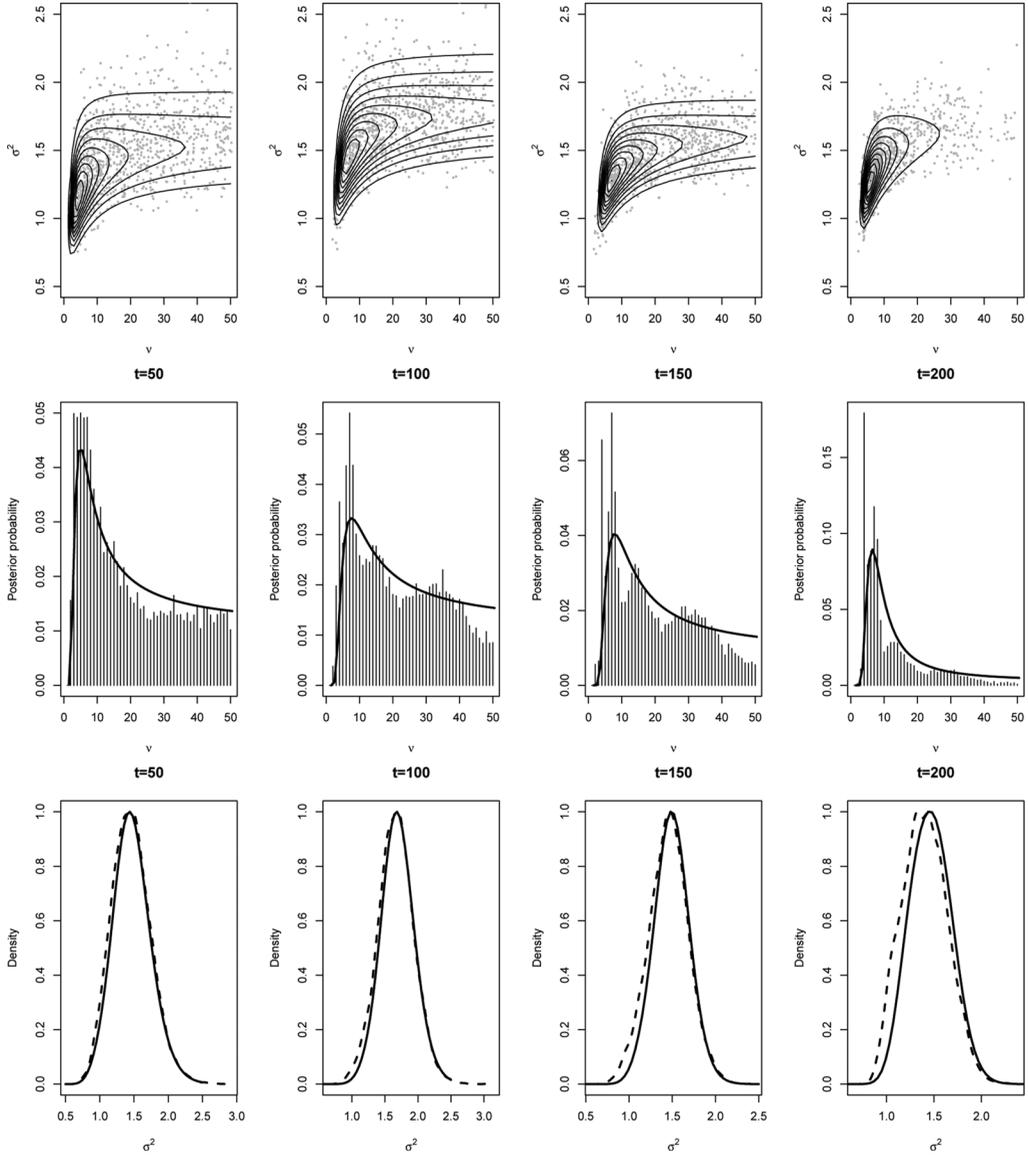


FIGURE 2 *i.i.d. model*. Sequential posterior inference for (σ^2, ν) based on PL for $T = 200$ iid observations drawn from t_4 with uniform prior for ν . PL is based on $N = 10,000$ particles.

negligible mass for values $\nu > 10$. Figure 4 shows that PL estimates are still accurate when $n = 1,000$. It also shows that the marginal posterior of ν is highly concentrated around the true value for $t > 500$, as theoretically predictive in Section 4.1 and Fig. 1.

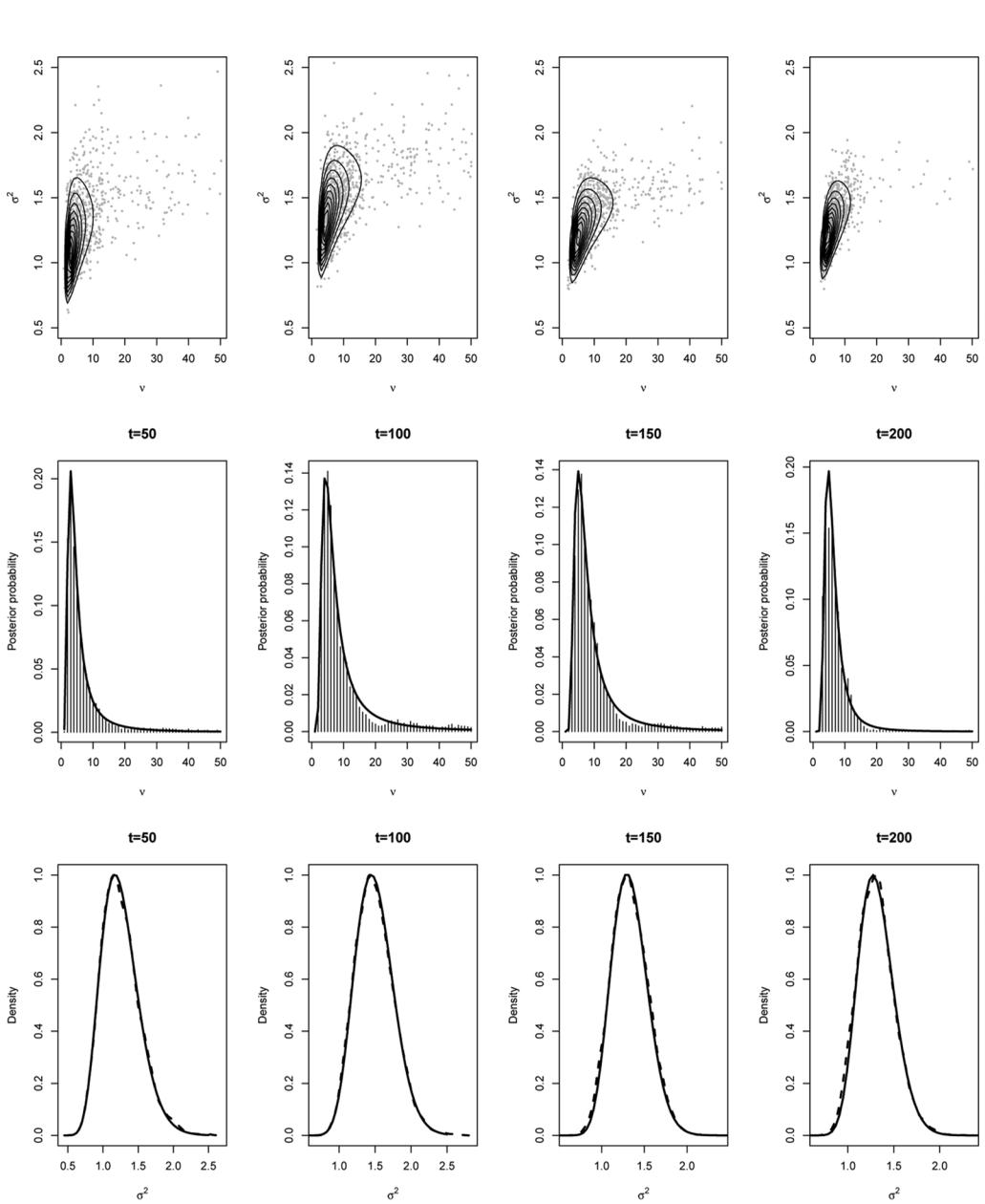
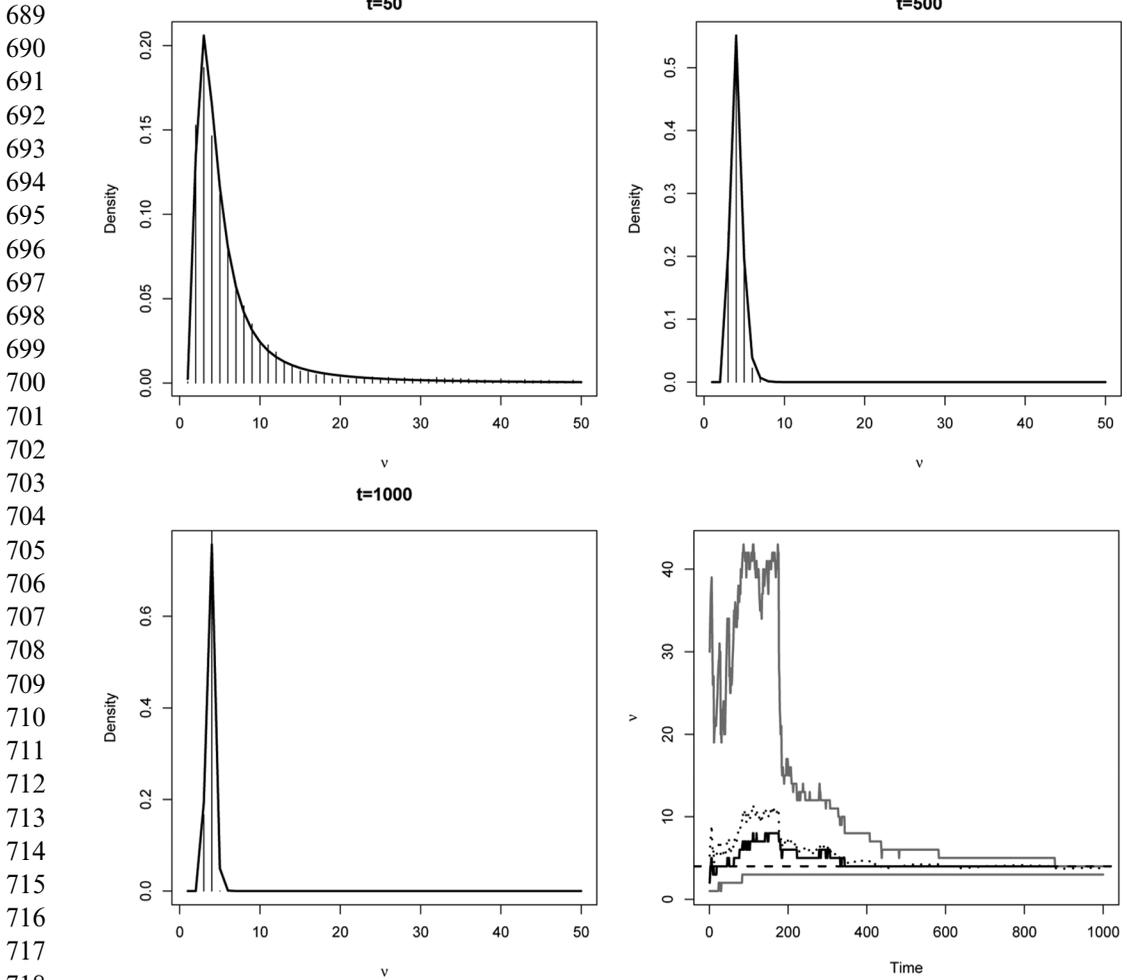


FIGURE 3 *i.i.d. model*. Sequential posterior inference for (σ^2, ν) based on PL for $T = 200$ iid observations drawn from t_4 with Jeffreys prior for ν . PL is based on $N = 10,000$ particles.



719 FIGURE 4 *i.i.d. model*. Sequential posterior inference for ν based on PL for $T = 1,000$ iid observations
 720 drawn from t_4 with Jeffreys prior for ν . PL is based on $N = 10,000$ particles.

721
 722
 723 The undesirable bias of the not-so-harmless uniform prior is highlighted in the Monte
 724 Carlo exercise summarized by Figs. 5 and 6. The posterior means, medians, and modes of
 725 ν based on $p(\nu|y^t)$, $t = 30, 50, 100, 300, 400$, and 500 are compared across $R = 50$ samples.
 726 As it can be seen, the bias of the uniform prior is striking for samples of size up to
 727 $T = 100$, when compared to those of the Jeffreys prior. For samples of size $T = 400$ and
 728 $T = 500$, the bias is much smaller, but a closer look reveals its presence. For example, the
 729 25th percentiles of the mean, median, and mode box-plots when $T = 500$ are all above
 730 the true value $\nu = 4$ for the uniform prior.
 731

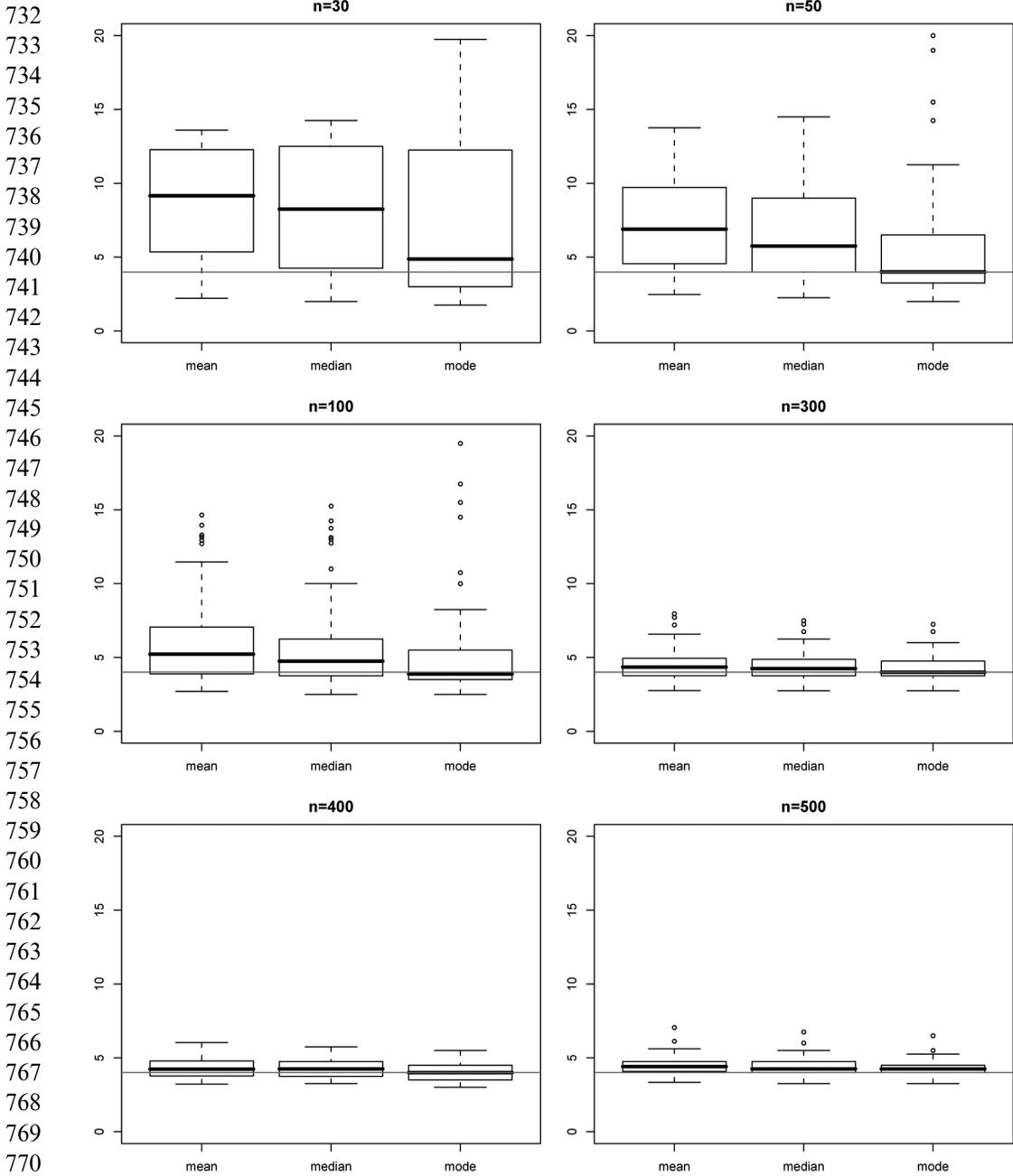
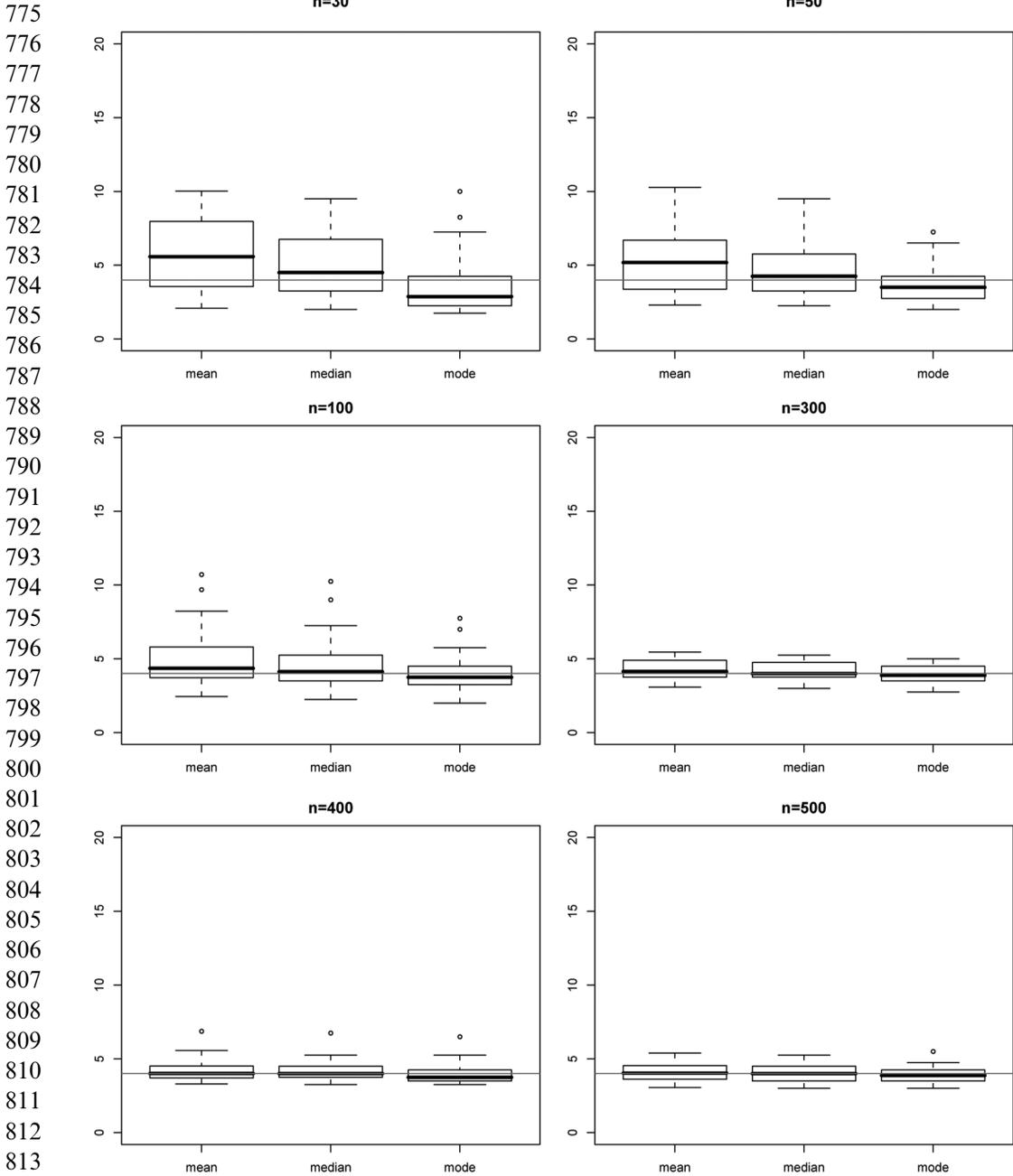


FIGURE 5 *i.i.d. model*. Posterior mean, median, and mode for the number of degrees of freedom ν under the uniform prior, for different sample sizes and based on a Gibbs sampler of length $M = 1,000$ after a burn-in period of M_0 draws. Boxplots are based on $R = 50$ datasets.

774



775
 776
 777
 778
 779
 780
 781
 782
 783
 784
 785
 786
 787
 788
 789
 790
 791
 792
 793
 794
 795
 796
 797
 798
 799
 800
 801
 802
 803
 804
 805
 806
 807
 808
 809
 810
 811
 812
 813
 814
 815
 816
 817

FIGURE 6 *i.i.d. model*. Posterior mean, median, and mode for the number of degrees of freedom ν under the Jeffreys prior, for different sample sizes and based on a Gibbs sampler of length $M = 1,000$ after a burn-in period of M_0 draws. Boxplots are based on $R = 50$ datasets.

5.2. The SV- t Model

We now revisit the well-known British pound versus U.S. dollar exchange rate data of Jacquier et al. (2004). The data consists of $T = 937$ daily rates from October 1st, 1981 to June, 28th 1985. For illustration purposes, we simulated data with exactly the same length from a SV- t_4 model with parameters $(\nu, \alpha, \beta, \tau^2) = (4, -0.202, 0.980, 0.018)$ and initial value $h_0 = -8.053$. Both simulated and real data sets are presented in Fig. 7.

The prior distribution of ν is given by the discretized version of Fonseca et al.'s (2008) Jeffreys prior, similar to the approach taken in Section 5.1 (see Eq. 1). The vector log-volatility parameters (α, β, τ^2) are independent, *a priori*, of ν and its prior distribution is given by $(\alpha, \beta) | \tau^2 \sim N(b_0, \tau^2 B_0)$ and $\tau^2 \sim IG(\eta_0/2, \eta_0 \tau_0^2/2)$, while the posterior for the log-volatility at time $t = 0$ is given by $h_0 \sim N(m_0, C_0)$. The hyper-parameters are set at

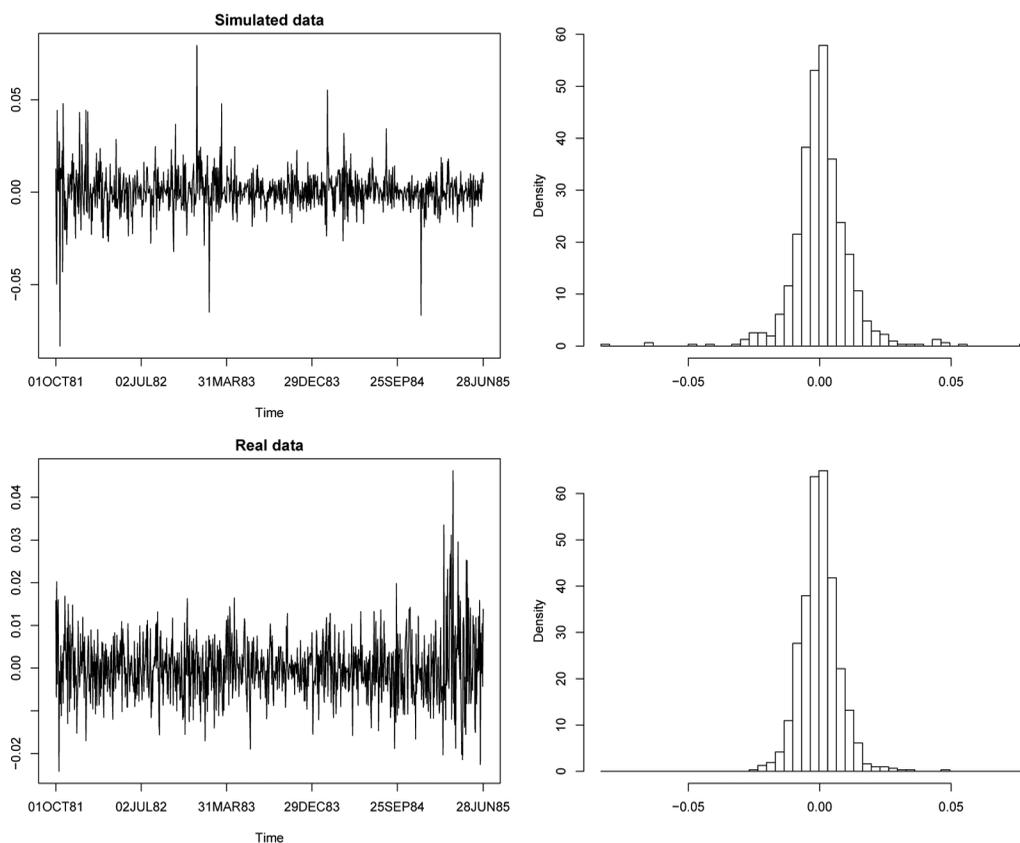


FIGURE 7 SV- t model. The top row corresponds to simulated data ($T = 937$) from the SV- t_ν model with parameters $\nu = 4$, $\alpha = -0.202$, $\beta = 0.980$, $\tau^2 = 0.018$, and $x_0 = -8.053$. The bottom row corresponds to JPR's (1994) British pound vs. U.S. dollar exchange ($T = 937$) daily rates from go from October 1, 1981 to June 28, 1985.

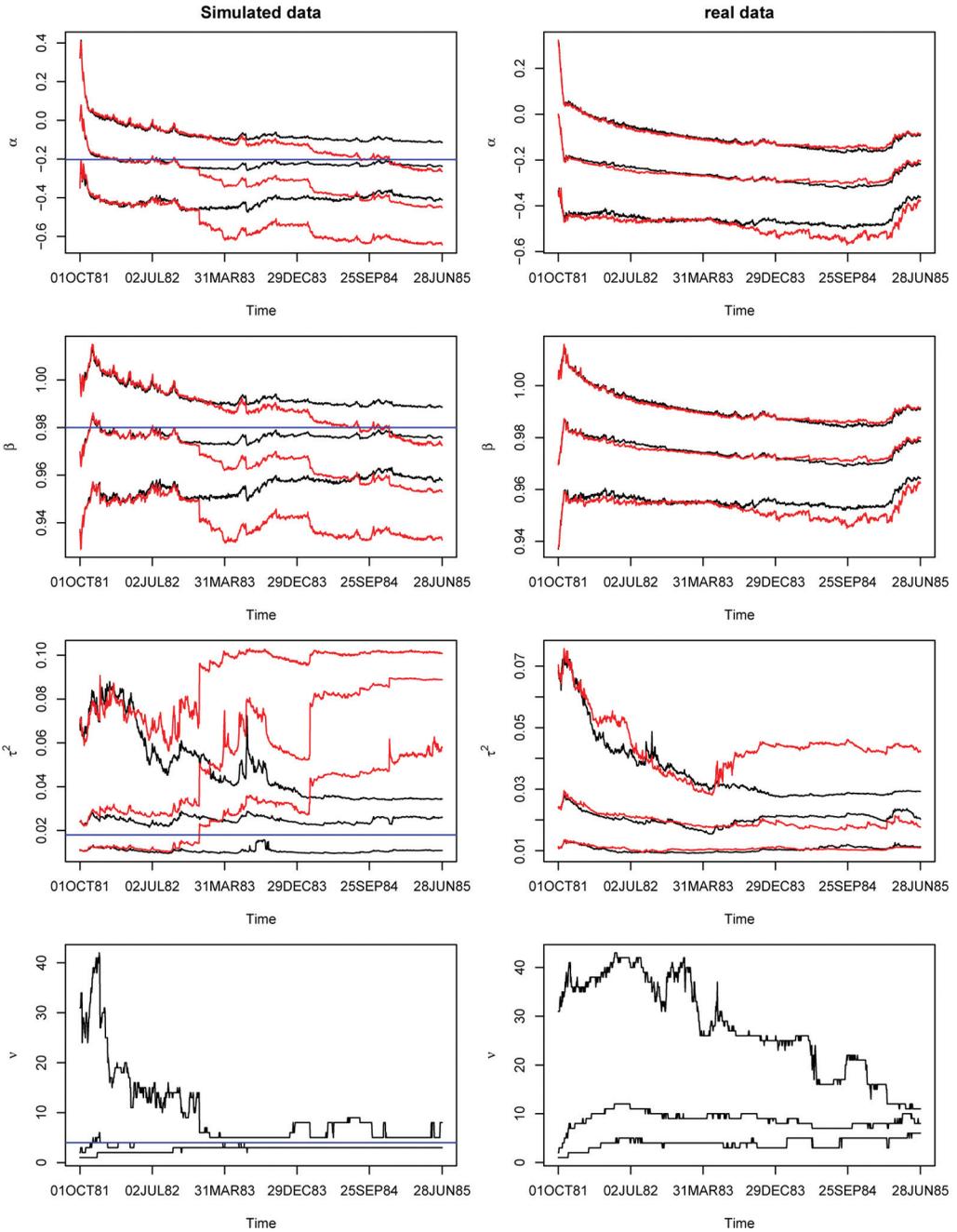


FIGURE 8 *SV-t model*. (2.5, 50, 97.5)th percentiles of the sequential marginal posterior distributions of α , β , τ^2 , and ν for the normal (red lines) and Student's t (black lines) models.

861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903

904 the values $m_0 = \log y_1^2$, $C_0 = 1.0$, $b_0 = (-0.002, 0.97)$, $B_0 = \text{diag}(1.0, 0.01)$, $c_0 = 5.0$, and
 905 $d_0 = 0.1125$.

906 Posterior inference is based on PL with $N = 10,000$ particles. Figures 8 presents 2.5th,
 907 50th and 97.5th percentiles of the sequential marginal distributions of α , β , τ^2 , and ν for
 908 both simulated and real data sets. For the simulated data, the posterior distribution of ν
 909 concentrates around the true value $\nu = 4$ after about 350 observations. For the real data,
 910 ν is highly concentrated with around ten degrees of freedom at the end of the sample;
 911 however, the right tail of the distribution, i.e., large degrees of freedom, is fairly long for
 912 most of the sample. Another interesting fact is that both normal and Student's t model
 913 learn about α and β in a similar manner, while the same cannot be said for the volatility
 914 of the log-volatility parameter, τ^2 . This is perhaps not surprising as the normal model
 915 overestimates the volatility of log-volatility to accommodate the fact that daily rates
 916 violate the plain normality assumption. The same behavior is present in our simulated
 917 data exercise. In fact, the posterior distribution for the log-volatilities, $p(h_t|y^t)$, for the
 918 simulated data based on the normal model has larger uncertainty than for the t_ν model
 919 (figure not shown here). Finally, at the end of the sample we can calculate the marginal
 920 posterior on the tail-thickness $p(\nu|y^T)$, our sequential particle approach agrees with the
 921 MCMC analysis of Jacquier et al. (2004). This suggests that the MC accumulation error
 922 inherent in our particle algorithm is small for these types of data length and models.
 923

924 5.2.1. S&P500: Credit Crisis 2008–2009

925 To study the effect of the credit crisis on stock returns, we revisit daily S&P500 returns
 926 previously studied, amongst many others, by Abanto-Valle et al. (2010) and Lopes and
 927 Polson (2010b). The former article estimates SV models with errors in the class of
 928 symmetric scale mixtures of normal distributions and also base their illustration on the
 929 S&P500 index from January 1999 to September 2008, therefore missing most of the credit
 930 crunch crisis and its aftermath. We concentrate our analysis on the period starting on
 931 January 3, 2007 and ending on October 14, 2010 ($T = 954$ observations). We sequentially
 932 fit the normal model to this data set as well as the t_ν model for $\nu \in \{5, 10, 50\}$. Figure 9
 933 summarizes our findings. The three Student's t models have higher predictive power
 934 than the normal model when measured in terms of log-Bayes factors. This distinction
 935 is particularly strong when comparing the t_5 (or t_{10}) model with the normal model.
 936 Interestingly, the t_5 model becomes gradually closer to the normal model from July 2008
 937 to July 2010, when again it distances itself from normality.

941 Before the onset of the credit crisis in July 2008, the model with the largest Bayes
 942 factor (relative to a normal), and hence the largest posterior model probability (under a
 943 uniform prior on ν) is the t_5 -distribution. This is not surprising as the previous time period
 944 consisted of little stochastic volatility and the occasional outlying return—which is nicely
 945 accommodated by a t_5 error distribution, in the spirit of Jeffreys initial observation about
 946 “real” data. The interesting aspects of Bayesian learning occur in the period of the crisis

947
 948
 949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971
 972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989

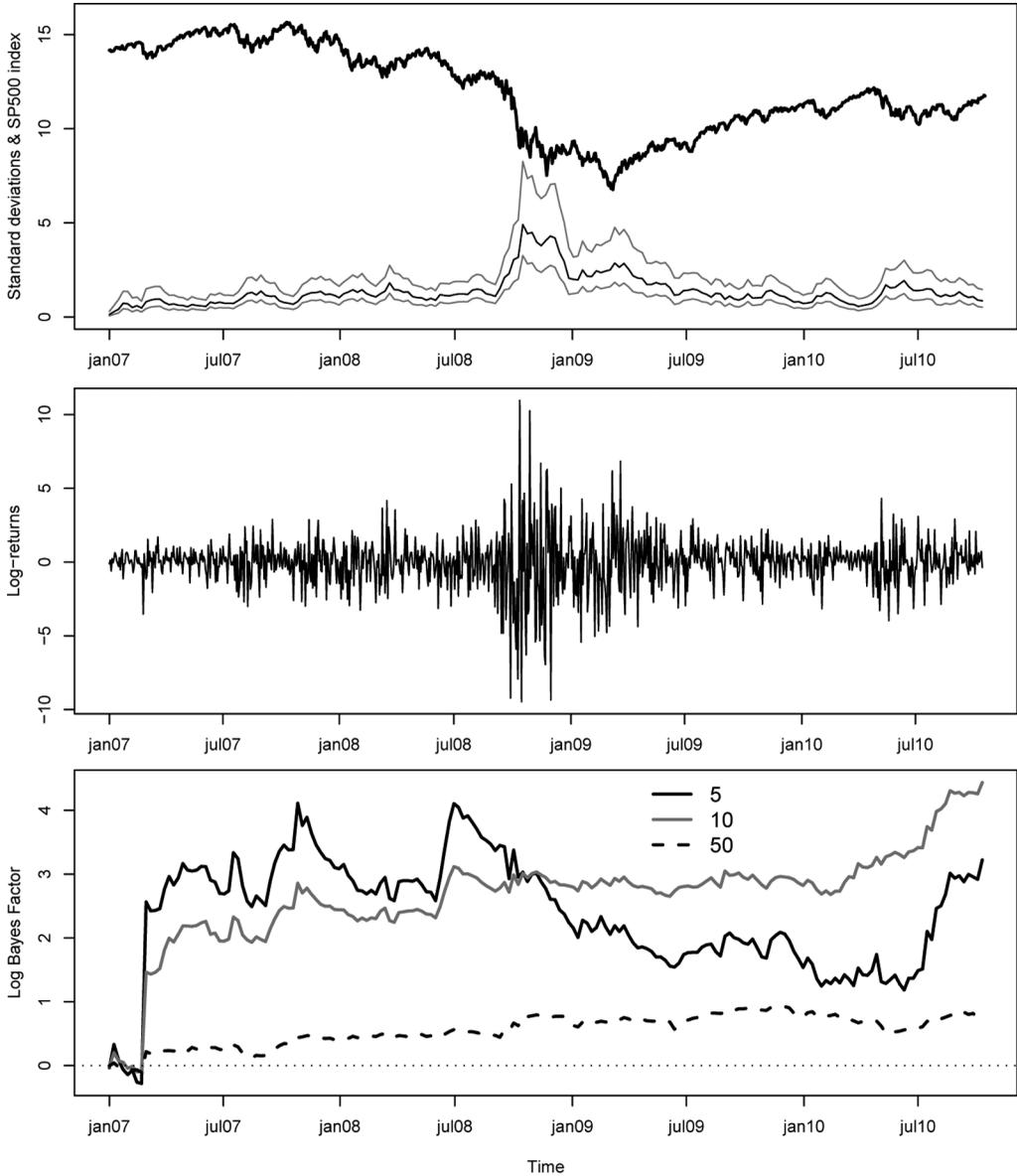


FIGURE 9 *SV-t model for S&P500 returns.* Top frame: S&P500 daily closing price (divided by 100: solid thick line) along with PL approximations to the (2.5, 50, 97.5)th percentiles of the posterior distributions of the time-varying standard deviations $p(\exp\{x_t/2\}|y^t)$, for $t = 1, \dots, T$. Middle frame: Log returns. Bottom frame: Logarithm of the Bayes factors of t_ν against normality for $\nu \in \{5, 10, 50\}$.

990 from July 2008 to March 2009. One immediately sees a dramatic increase in the stochastic
 991 volatility component of the model and the clustering of a high period of volatility. In and
 992 of itself, this is sufficient to “explain” the extreme moves in the market. Correspondingly,
 993 in terms of online estimation of the fat-tails, the Bayes factor quickly moves to favor the
 994 model with light tails, here the t_{10} -distribution. Finally, as the crisis subsides, the volatility
 995 mean reverts and the returns again look like they exhibit some outlying behavior (relative
 996 to the level of volatility) and the sequential Bayes again starts to move to favor the fatter-
 997 tailed t_5 -distribution.
 998
 999

1000 6. DISCUSSION

1001 Estimating tail-thickness of the error distribution of an economic or financial time series
 1002 is an important problem as estimates and forecasts are very sensitive to the tail behavior.
 1003 Moreover, we would like an on-line estimation methodology that can adaptively learn the
 1004 tail-thickness and provide parameter estimates that update as new data arrives. We model
 1005 the error distribution as a t_ν -distribution where $\nu \sim p(\nu)$, and we adopt a default Jeffreys
 1006 prior on the tail-thickness parameter ν . We show that this has a number of desirable
 1007 properties when performing inference with a finite amount of data. We use the sequential
 1008 Bayes factor to provide an on-line test of normality versus fat-tails, and we derive its
 1009 optimality properties asymptotically and in finite sample using a Kullback–Leibler metric.
 1010 We illustrate these effects in the credit crisis of 2008–2009 with daily S&P500 stock return
 1011 data. Our analysis shows how quickly an agent can dynamically learn the tail of the error
 1012 distribution whilst still accounting for parameter uncertainty and time-varying stochastic
 1013 volatility. Figures 2–4 and 8 all show that estimating ν is in fact rather difficult. Figure 8,
 1014 in particular, shows that when the data is not normal it takes several time periods for the
 1015 parameter ν be stably estimated.
 1016

1017 Whilst MCMC is computationally slow for solving the online problem, it does also
 1018 provide the full smoothing distribution at the end of the sampler. This would require
 1019 $O(N^2)$ particles in our approach (see Carvalho et al., 2010, for further discussion),
 1020 and therefore, if smoothed states are required, we recommend filtering forward with
 1021 particles and smoothing with MCMC. Other estimation methods such as nested Laplace
 1022 approximation (Smith, 2000) seem unable to identify the true error structure due to
 1023 the multimodalities present in the posterior and particle methods provide a natural
 1024 alternative. Clearly, there are a number of extensions of our approach, for example, to
 1025 multivariate and dynamic panel data.
 1026
 1027

1028 REFERENCES

- 1029
 1030 Abanto-Valle, C. A., Bandyopadhyay, D., Lachos, V. H., Enriquez, I. (2010). Robust Bayesian analysis
 1031 of heavy-tailed stochastic volatility models using scale mixtures of normal distributions. *Computational*
 1032 *Statistics and Data Analysis* 54:2883–2898.

- 1033 Andrews, D. F., Mallows, C. L. (1974). Scale mixtures of normal distributions. *Journal of Royal Statistical*
 1034 *Society, Series B* 36:99–102. Q9
- 1035 Berk, R. H. (1966). Limiting behavior of posterior distributions when the model is incorrect. *Annals of*
 1036 *Mathematical Statistics* 37:51–58.
- 1037 Berk, R. H. (1970). Consistency a posteriori. *Annals of Mathematical Statistics* 41:894–906.
- 1038 Blackwell, D., Dubins, L. (1962). Merging of opinions with increasing information. *Annals of Mathematical*
 1039 *Statistics* 33:882–886.
- 1040 Carlin, B. P., Polson, N. G., Stoffer, D. S. (1992). A Monte Carlo approach to nonlinear and non-normal
 1041 state space models. *Journal of the American Statistical Association* 87:493–500.
- 1042 Carvalho, C. M., Johannes, M. S., Lopes, H. F., Polson, N.G. (2010). Particle learning and smoothing.
 1043 *Statistical Science* 25:88–106.
- 1044 Chib, S., Nardari, F., Shephard, N. (2002). Markov chain Monte Carlo methods for stochastic volatility
 1045 models. *Journal of Econometrics* 108:281–316.
- 1046 Cover, T. M., Thomas, J. A. (2006). *Elements of Information Theory*. 2nd ed. New York: Wiley.
- 1047 Edgeworth, F. Y. (1888). On a new method of reducing observations relating to several quantities.
 1048 *Philosophical Magazine* 25:184–191.
- 1049 Edwards, W., Lindman, H., Savage, L. J. (1963). Bayesian statistical inference for psychological research.
 1050 *Psychological Review* 70:193–242.
- 1051 Eraker, B., Jacquier, E., Polson, N. G. (1998). The pitfalls of MCMC algorithms. Technical Report, The
 1052 University of Chicago Booth School of Business.
- 1053 Fernandez, C., Steel, M. F. J. (1998). On Bayesian modeling of fat tails and skewness. *Journal of the*
 1054 *American Statistical Association* 93:359–371.
- 1055 Fonseca, T., Ferreira, M. A. R., Migon, H. S. (2008). Objective Bayesian analysis for the Student-*t* regression
 1056 model. *Biometrika* 95:325–333.
- 1057 Gamerman, D., Lopes, H. F. (2006). *Markov Chain Monte Carlo: Stochastic Simulation for Bayesian*
 1058 *Inference*. Baton Rouge: Chapman & Hall/CRC.
- 1059 Geweke, J. (1993). Bayesian treatment of the independent Student-*t* linear linear model. *Journal of Applied*
 1060 *Econometrics* 8:19–40.
- 1061 Gordon, N., Salmond, D., Smith, A. F. M. (1993). Novel approach to nonlinear/non-Gaussian Bayesian state
 1062 estimation. *IEEE Proceedings* F-140:107–113. Q10
- 1063 Gordon, N., Smith, A. F. M. (1993). Approximate non-Gaussian Bayesian estimation and modal consistency.
 1064 *Journal of Royal Statistical Society, Series B* 55:913–918. Q10
- 1065 Gramacy, R., Pantaleo, E. (2010). Shrinkage regression for multivariate inference with missing data, and an
 1066 application to portfolio balancing. *Bayesian Analysis* 5:237–262.
- 1067 Jacquier, E., Polson, N. G. (2000). Discussion of “Time series analysis of non-Gaussian observations”. *Journal*
 1068 *of Royal Statistical Society, B* 62:44–45.
- 1069 Jacquier E., Polson, N. G., Rossi, P. E. (2004). Bayesian analysis of stochastic volatility with fat tails and
 1070 leverage effect. *Journal of Econometrics* 122:185–212.
- 1071 Jeffreys, H. (1961). *Theory of Probability*. New York: Oxford University Press.
- 1072 Lindley, D. V. (1956). On a measure of the information provided by an experiment. *Annals of Mathematical*
 1073 *Statistics* 27:986–1005.
- 1074 Lopes, H. F., Carvalho, C. M. (2013). Online Bayesian learning in dynamic models: An illustrative
 1075 introduction to particle methods. In: West, M., Damien, P., Dellaportas, P., Polson, N. G., Stephens, D.
 A., eds. *Bayesian Theory and Applications*. Clarendon: Oxford University Press, pp. 203–228.
- 1076 Lopes, H. F., Carvalho, C. M., Johannes, M. S., Polson, N. G. (2010). Particle learning for sequential
 1077 Bayesian computation (with discussion). In: Bernardo, J. M., Bayarri, M. J., Berger, J. O., Dawid, A. P.,
 1078 Heckerman, D., Smith, A. F. M., West, M., eds. *Bayesian Statistics*, Vol. 9. Oxford: Oxford University
 1079 Press. To appear. Q11
- 1080 Lopes, H. F., Polson, N. G. (2010a). Bayesian inference for stochastic volatility modeling. In: Böcker, K., ed.
 1081 *Re-Thinking Risk Measurement, Management and Reporting Measurement Uncertainty, Bayesian Analysis*
 1082 *and Expert Elicitation*. Riskbooks, pp. 515–551.

- 1076 Lopes H. F., Polson, N. G. (2010b). Extracting SP500 and NASDAQ volatility: The credit crisis of 2007-2008.
1077 In: O'Hagan, A., West, M. eds. *Handbook of Applied Bayesian Analysis*. Oxford: Oxford University
1078 Press, pp. 319–342.
- 1079 Lopes, H. F., Tobias, J. (2011). Confronting prior convictions: On issues of prior and likelihood sensitivity
1080 in Bayesian analysis. *Annual Review of Economics* 3:107–131.
- 1081 Lopes, H. F., Tsay, R. (2011). Particle filters and Bayesian inference in financial econometrics. *Journal of*
1082 *Forecasting* 30:168–209.
- 1083 McCausland, W. (2012). The HESSIAN method: Highly efficient simulation smoothing, in a nutshell. *Journal*
1084 *of Econometrics* 168:189–206.
- 1085 Smith, A. F. M. (1983). Bayesian approaches to outliers and robustness. In: Florens, J. P., Mouchart,
1086 M., Raoult, J. P., Simar, L., Smith, A. F. M. eds. *Specifying Statistical Models: From Parametric to*
1087 *Nonparametric, Using Bayesian or Non-Bayesian Approaches*. New York: Springer-Verlag, pp. 13–35.
- 1088 Smith, J. Q. (2000). In discussion of “Time series analysis of non-Gaussian observations”. *Journal of Royal*
1089 *Statistical Society, B* 62:29–20.
- 1090 Verdinelli, I., Wasserman, L. (1995). Computing Bayes factors by using a generalization of the Savage-Dickey
1091 density ratio. *Journal of the American Statistical Association* 90:614–618.
- 1092 West, M. (1981). Robust sequential approximate Bayesian estimation. *Journal of Royal Statistical Society,*
1093 *Series B* 43:157–166.
- 1094 West, M. (1984). Bayesian model monitoring. *Journal of Royal Statistical Society, Series B* 48:70–78.
- 1095
- 1096
- 1097
- 1098
- 1099
- 1100
- 1101
- 1102
- 1103
- 1104
- 1105
- 1106
- 1107
- 1108
- 1109
- 1110
- 1111
- 1112
- 1113
- 1114
- 1115
- 1116
- 1117
- 1118

Q12

Q12