

SVERIGES RIKSBANK
WORKING PAPER SERIES

233



Flexible Modeling of Conditional Distributions Using Smooth Mixtures of Asymmetric Student T Densities

Feng Li, Mattias Villani and Robert Kohn

OCTOBER 2009

WORKING PAPERS ARE OBTAINABLE FROM

Sveriges Riksbank • Information Riksbank • SE-103 37 Stockholm
Fax international: +46 8 787 05 26
Telephone international: +46 8 787 01 00
E-mail: info@riksbank.se

The Working Paper series presents reports on matters in the sphere of activities of the Riksbank that are considered to be of interest to a wider public.

The papers are to be regarded as reports on ongoing studies and the authors will be pleased to receive comments.

The views expressed in Working Papers are solely the responsibility of the authors and should not to be interpreted as reflecting the views of the Executive Board of Sveriges Riksbank.

FLEXIBLE MODELING OF CONDITIONAL DISTRIBUTIONS USING SMOOTH MIXTURES OF ASYMMETRIC STUDENT T DENSITIES

FENG LI, MATTIAS VILLANI, AND ROBERT KOHN

Sveriges Riksbank Working Paper Series No. 233

October 2009

ABSTRACT. A general model is proposed for flexibly estimating the density of a continuous response variable conditional on a possibly high-dimensional set of covariates. The model is a finite mixture of asymmetric student-t densities with covariate dependent mixture weights. The four parameters of the components, the mean, degrees of freedom, scale and skewness, are all modelled as functions of the covariates. Inference is Bayesian and the computation is carried out using Markov chain Monte Carlo simulation. To enable model parsimony, a variable selection prior is used in each set of covariates and among the covariates in the mixing weights. The model is used to analyse the distribution of daily stock market returns, and shown to more accurately forecast the distribution of returns than other widely used models for financial data.

KEYWORDS: Bayesian inference, Markov Chain Monte Carlo, Mixture of Experts, Variable selection, Volatility modeling.

Li: *Department of Statistics, Stockholm University, SE-106 91 Stockholm, Sweden.* e-mail: feng.li@stat.su.se. Villani: *Research Division, Sveriges Riksbank and Department of Statistics, Stockholm University.* Kohn: *Australian School of Business, University of New South Wales.* The views expressed in this paper are solely the responsibility of the authors and should not be interpreted as reflecting the views of the Executive Board of Sveriges Riksbank. Robert Kohn was partially supported by ARC grant DP0667069.

1. INTRODUCTION

This paper is concerned with estimating the conditional predictive distribution $p(y|x)$, where y is a univariate continuous response variable and x is a possibly high-dimensional vector of covariates. Our approach is an exercise in nonparametric regression density estimation since $p(y|x)$ is modeled flexibly both for any given x but also across different covariate values.

Villani, Kohn and Giordani (2008) propose the Smooth Adaptive Gaussian Mixture (SAGM) model as flexible model for regression density estimation. Their model is a finite mixture of Gaussian densities with the mixing probabilities, the component means and component variances modeled as functions of the covariates x , with Bayesian variable selection in all three sets of covariates. See Früewirth-Schnatter (2006) for a comprehensive introduction to mixture models.

Villani et al. (2008) argues in favor of a *complex-and-few* modeling philosophy where enough flexibility is used within the mixture components so that the number of components can be kept to a minimum; see also Wood, Jiang and Tanner (2003). This is in sharp contrast to the *simple-and-many* approach used in the machine learning literature (in particular the mixture-of-experts model introduced in Jacobs, Jordan, Nowlan and Hinton (1991), and Jordan and Jacobs (1994)) where the components are often linear homoscedastic regressions, or even constant functions. Villani et al. (2008) show that a single complex component can often give a better and numerically more stable fit in substantially less computing time than a model with many simpler components. Moreover, simulations and real applications in Villani et al. (2008) show that a simple-and-many approach can fail to fit heteroscedastic data even with a very large number of components, especially in situations with more than one or two covariates. Having heteroscedastic components in the mixture is therefore crucial for accurately modeling heteroscedastic data.

In one of their applications, Villani et al. (2008) model the distribution of daily stock market returns as a function of lagged returns and smooth measures of recent volatility. The best model uses one component to fit the strong heteroscedasticity in the data and the other two or three components to capture the additional kurtosis and/or skewness. The current paper continues the complex-and-few approach and extends the SAGM model by generalizing the Gaussian components to asymmetric student- t densities, thereby making it possible to capture skewness and excess kurtosis within the components. Each component density has four parameters: location, scale, degrees of freedom and skewness, and each of these four parameters are modeled as function of covariates. This makes it possible to have e.g. the degrees of freedom smoothly varying over covariate space in a way dictated by the data. An efficient Markov Chain Monte Carlo (MCMC) simulation method is proposed that allows for Bayesian variable selection in all four parameters of the asymmetric t density, and in the mixture weights. The variable selection makes it possible to handle a large number of covariates. Reducing the number of effective parameters by variable selection mitigates problems with over-fitting and is also beneficial for the convergence of the MCMC algorithm. The methodology is applied to modeling the distribution of daily returns from the S&P500 stock market index. It is shown that a smooth mixture of asymmetric student t components outperforms SAGM and

other commonly used models for financial data in an out-of-sample evaluation of the predictive density during the financial turmoil in the end of year 2008 and beginning of 2009.

2. THE MODEL AND PRIOR

2.1. Smooth mixtures. Our model is a finite mixture density with weights that are smooth functions of the covariates,

$$(2.1) \quad p(y|x) = \sum_{k=1}^K \omega_k(x) p_k(y|x),$$

where $p_k(y|x)$ is the k th component density with weight $\omega_k(x)$. The component densities are asymmetric student t densities described in detail in the next section. The weights are modeled by a multinomial logit function

$$(2.2) \quad \omega_k(x) = \frac{\exp(x'\gamma_k)}{\sum_{r=1}^K \exp(x'\gamma_r)},$$

with $\gamma_1 = 0$ for identification. The covariates in the components can in general be different from the covariates in the mixture weights. Jiang and Tanner (1999a,b) show that smooth mixtures with sufficiently many components can approximate a wide class of densities.

To simplify the MCMC simulation, we express the mixture model in terms of latent variables as in Diebolt and Robert (1994) and Escobar and West (1995). Let s_1, \dots, s_n be unobserved indicator variables for the observations in the sample such that $s_i = k$ means that the i th observation belongs to the k th component, $p_k(y|x)$. The model in (2.1) and (2.2) can then be written

$$\begin{aligned} \Pr(s_i = k|x_i, \gamma) &= \omega_k(x_i) \\ y_i|(s_i = k, x_i) &\sim p_k(y_i|x_i). \end{aligned}$$

Conditional on $s = (s_1, \dots, s_n)'$, the mixture model decomposes into K separate component models $p_1(y|x), \dots, p_K(y|x)$, with each data observation being allocated to one and only one component.

2.2. The component models. The component densities in SAGM are Gaussian with both the mean and variance functions of covariates. Our article extends this model so the component densities belong to an asymmetric student t family. More specifically, the component models are split- t densities (Geweke, 1989) according to the following definition.

Definition 1. The random variable y follows a *split- t distribution* with ν degrees of freedom, $y \sim t(\mu, \phi, \lambda, \nu)$, if its density function is of the form

$$c \cdot \kappa(\mu, \phi, \nu) I(y \leq \mu) + c \cdot \kappa(\mu, \lambda\phi, \nu) I(y > \mu),$$

where

$$\kappa(\mu, \phi, \nu) = \left(\frac{\nu}{\nu + \frac{(y-\mu)^2}{\phi^2}} \right)^{(\nu+1)/2},$$

is the kernel of a student t density with variance $\phi^2\nu/(\nu-2)$ and $c = 2[(1+\lambda)\phi\sqrt{\nu}\text{Beta}(\frac{\nu}{2}, \frac{1}{2})]^{-1}$ is the normalization constant. The location parameter is μ , $\phi > 0$ is the scale parameter, and $\lambda > 0$ is the skewness parameter. When $\lambda < 1$ the distribution is skewed to the left, when $\lambda > 1$ it is skewed to the right, and when $\lambda = 1$ it reduces to the usual symmetric student- t density.

The (one-component) split- t is similar to the ARCD model of Hansen (1994) which he estimates by maximum likelihood to model the conditional density of the U.S. Dollar / Swiss Franc exchange rate.

The next lemma gives the first four central moments of the split- t density. We use the following definition of skewness and excess kurtosis

$$S(y) = \frac{E[y - E(y)]^3}{V(y)^{3/2}}$$

$$K(y) = \frac{E[y - E(y)]^4}{V(y)^2} - 3,$$

where $V(y)$ denotes the variance. The following lemma, which can be proved by straightforward algebra, gives some basic properties of the split- t distribution.

Lemma 2. *If $y \sim t(\mu, \phi, \lambda, \nu)$ then*

$$\begin{aligned} E(y) &= \mu + h \\ V(y) &= \frac{1 + \lambda^3}{1 + \lambda} \frac{\nu}{\nu - 2} \phi^2 - h^2 \\ E[y - E(y)]^3 &= \frac{6\lambda(\lambda^2 - 1) - 2(\lambda^4 - 1)}{(1 + \lambda)(\nu - 1)(\nu - 3)} \text{Beta}\left(\frac{\nu}{2}, \frac{1}{2}\right) \nu^{\frac{3}{2}} \phi^3 + 2h^3 \\ E[y - E(y)]^4 &= \frac{3\nu^2 \phi^4 (1 + \lambda^5)}{(1 + \lambda)(\nu - 2)(\nu - 4)} - 3h^4 + 6h^2 \frac{(1 + \lambda^3) \nu}{(1 + \lambda)(\nu - 2)} \phi^2 \\ &\quad - 16h \frac{(\lambda - 1)(\lambda^2 + 1) \nu^{\frac{3}{2}} \phi^4}{(\nu - 1)(\nu - 3) \phi \text{Beta}\left(\frac{\nu}{2}, \frac{1}{2}\right)}, \end{aligned}$$

where

$$h = \frac{2\sqrt{\nu}\phi(\lambda - 1)}{(\nu - 1) \text{Beta}\left(\frac{\nu}{2}, \frac{1}{2}\right)},$$

and moment of order r exists if $\nu > r$.

The CDF of a split t distribution is of the form

$$\frac{1}{1 + \lambda} + \frac{a \cdot \text{Sign}(y - \mu)}{1 + \lambda} \left(1 - \frac{\text{Beta}\left(t; \frac{\nu}{2}, \frac{1}{2}\right)}{\text{Beta}\left(\frac{\nu}{2}, \frac{1}{2}\right)} \right)$$

where

$$t = \frac{\nu a^2 \phi^2}{\nu a^2 \phi^2 + (y - \mu)^2},$$

and $a = \lambda$ if $y > \mu$ and $a = 1$ otherwise, and $\text{Beta}\left(t; \frac{\nu}{2}, \frac{1}{2}\right)$ is the incomplete beta function (Abramowitz and Stegun, 1972).

Each of the four parameters μ, ϕ, λ and ν are connected to covariates as

$$\begin{aligned}
 \mu &= \beta_{\mu 0} + x'_t \beta_{\mu} \\
 \ln \phi &= \beta_{\phi 0} + x'_t \beta_{\phi} \\
 \ln \lambda &= \beta_{\lambda 0} + x'_t \beta_{\lambda} \\
 \ln \nu &= \beta_{\nu 0} + x'_t \beta_{\nu}
 \end{aligned}
 \tag{2.3}$$

but any smooth link function can equally well be used in the MCMC methodology. Additional flexibility can be obtained by letting a subset of the covariates be a non-linear basis expansions, e.g. additive splines or splines surfaces (Ruppert, Wand and Carroll, 2003) as in Villani et al. (2008), but this is not pursued here. A strength of our approach is that the four regression coefficient vectors: $\beta_{\mu}, \beta_{\phi}, \beta_{\nu}$ and β_{λ} are all treated in a unified way in the MCMC algorithm. Whenever we refer to a regression coefficient vector without subscript, β , the argument applies to all of the split- t parameters in (2.3).

This split- t model will often be flexible enough to fit the data, but there are datasets that require a smooth mixture model, for example when the data are multimodal for some covariates values. A second example occurs when the wrong link function is used in one of the split- t parameters, where the mixture can then correct for this erroneous choice. A third example is when there are outliers in the data that cannot be accommodated by a student t density.

A smooth mixture of split- t densities is a model with a large number of parameters, however, and is therefore likely to over-fit the data unless model complexity is controlled effectively. We use Bayesian variable selection in all four split- t parameters, and in the mixing function. This can lead to important simplifications of the split- t components. Not only does this control complexity for a given number of components, but it also simplifies the existing components if an additional component is added to the model (the LIDAR example in Villani, Kohn and Giordani (2007) illustrates this well). Increasing the number of components can therefore even reduce the number of effective parameters in the model.

A more extreme, but often empirically relevant, simplification of the model is to assume that one or more split- t parameters are *common* to the components, that is, only the intercepts in (2.3) are allowed to be different across components. The unrestricted model where the regression coefficients are allowed to differ across components is said to have *separate* components.

2.3. The prior. Although the MCMC methodology (see Section 3.2) allows any prior distribution, we shall now present an easily specified prior that depends only on a few hyper-parameters. First, we standardize the covariates by subtracting the mean and dividing by the standard deviation. This allows us to assume prior independence between the intercept and the remaining regression coefficients, and the intercepts have the interpretation of being the (possibly transformed) split- t parameters at the mean of the original covariates. Since there can be a large number of covariates in the model, our strategy is to incorporate available prior information via the intercepts, and to treat the remaining regression coefficients more informally. Assuming a normal prior for μ implies a normal prior on $\beta_{\mu 0}$. The other three

split- t parameters ϕ , λ and ν are assumed to follow independent log-normal priors with means m^* and s^* , where m^* and s^* are different for the different split- t parameters. This translates into a normal prior on the intercept with mean

$$m_0 = \ln m^* - \frac{1}{2} \ln \left[\left(\frac{s^*}{m^*} \right)^2 + 1 \right]$$

and variance

$$s_0^2 = \ln \left[\left(\frac{s^*}{m^*} \right)^2 + 1 \right].$$

The regression coefficients β_μ , β_ϕ , β_ν and β_λ are assumed to be independent a priori. We allow for Bayesian variable selection by augmenting each parameter vector β by a vector of binary covariate selection indicators $\mathcal{I} = (i_1, \dots, i_p)$ such that $\beta_j = 0$ if $i_j = 0$. Let $\beta_{\mathcal{I}}$ denote the subset of β selected by \mathcal{I} . We assume the following prior for each β vector

$$\beta_{\mathcal{I}} | \mathcal{I} \sim N(0, \tau_\beta^2 I)$$

and $\beta_{\mathcal{I}^c} | \mathcal{I}^c$ is identically zero, where \mathcal{I}^c is the complement of \mathcal{I} . Alternatively, one can use a g -prior (Zellner, 1986) $\beta \sim N \left[0, \tau_\beta^2 (X'X)^{-1} \right]$ and then condition on the restrictions imposed by \mathcal{I} ; Denison, Holmes, Mallick and Smith (2002, p. 80-81) discusses the advantages and disadvantages of these two different priors. The g -prior is less appealing in a mixture context since $(X'X)^{-1}$ may be a bad representation of the covariance between parameters in the smaller components, see Villani et al. (2008) for a discussion, and we will therefore use the identity matrix here. We use $\tau_\beta = 10$ as the default value. Given that the covariates have been standardized to zero mean and unit variance, these priors are vague. We investigate the sensitivity of the posterior inferences and model comparison with respect to τ_β in Section 4.

The variable selection indicators are assumed to be independent Bernoulli with probability ω_β a priori, but more complicated distributions are easily accommodated, see e.g. the extension in Villani et al. (2008) for splines in a mixture context or a prior which is uniform on the variable selection indicators for a given model size in Denison, Holmes, Mallick and Smith (2002). It is also possible to estimate ω_β as proposed in Kohn, Smith and Chan (2001) with an extra Gibbs sampling step. Note that ω_β may be different for each split- t parameter. Our default prior has $\omega_\beta = 0.5$.

The prior on the mixing function decomposes as

$$p(\gamma, \mathcal{Z}, s) = p(s | \gamma, \mathcal{Z}) p(\gamma | \mathcal{Z}) p(\mathcal{Z}),$$

where \mathcal{Z} is the $p \times (K - 1)$ matrix with variable selection indicators for the p covariates in the mixing function (recall that $\gamma_1 = 0$ for identification). The variable indicators in \mathcal{Z} are assumed to be *iid Bernoulli*(ω_γ). The prior on $\gamma = (\gamma'_2, \dots, \gamma'_m)'$ is assumed to be of the form

$$\gamma_{\mathcal{Z}} | \mathcal{Z} \sim N(0, \tau_\gamma^2 I),$$

and $\gamma_{\mathcal{Z}^c} = 0$ with probability one. We use $\tau_\gamma^2 = 10$ as default value. Finally, $p(s | \gamma, \mathcal{Z})$ is given by the multinomial logit model in (2.2). To reduce the number of parameters and to speed up the MCMC algorithm we restrict the columns of \mathcal{Z} to be identical, i.e. make the assumption

that a covariate is either present in the mixing function in all components, or does not appear at all, but the extension to general \mathcal{Z} is straightforward, see Villani et al. (2008).

3. INFERENCE METHODOLOGY

3.1. The general MCMC scheme. We use MCMC methods to sample from the joint posterior distribution, and draw the parameters and variable selection indicators in blocks. Villani et al. (2008) experimented with several different algorithms in a related setting and the algorithm outlined below is similar to their preferred algorithm. The details of the algorithm are given in Appendix A. The method used to select the number of components is discussed in Section 3.3.

The algorithm is a Metropolis-within-Gibbs sampler that draws parameters using the following six blocks:

- (1) $\{(\beta_\mu^{(k)}, \mathcal{I}_\mu^{(k)})\}_{k=1, \dots, K}$
- (2) $\{(\beta_\phi^{(k)}, \mathcal{I}_\phi^{(k)})\}_{k=1, \dots, K}$
- (3) $\{(\beta_\lambda^{(k)}, \mathcal{I}_\lambda^{(k)})\}_{k=1, \dots, K}$
- (4) $\{(\beta_\nu^{(k)}, \mathcal{I}_\nu^{(k)})\}_{k=1, \dots, K}$
- (5) $s = (s_1, \dots, s_n)$
- (6) γ and \mathcal{I}_Z

The parameters in the different components are independent conditional on s . This means that each of the first four blocks split up into K independent updating steps. Each updating step in the first four blocks is sampled using highly efficient tailored MH proposals following a general approach described in the next section. The latent component indicators in s are independent conditional on the model parameters and are drawn jointly from their full conditional posterior. Conditional on s , Step 6 is a multinomial logistic regression with variable selection, and γ and \mathcal{I}_Z are drawn jointly using a generalization of the method used to draw blocks 1-4, see Villani et al. (2008) for details.

Mixture models have well known identification problems, the most serious one being the so called label switching problem, which means that the likelihood is invariant with respect to permutations of the components in the mixture see e.g. Celeux, Hurn and Robert (2000), Jasra, Holmes and Stephens (2005) and Frühwirth-Schnatter (2006). The aim of our article is to estimate the predictive density, so that label switching is neither a numerical nor conceptual problem (Geweke, 2007). If an interpretation of the mixture components is required, then it is necessary to impose some identification restrictions on some of the model parameters, e.g. an ordering constraint (Jasra, Holmes and Stephens, 2005).

3.2. Updating (β, \mathcal{I}) using variable-dimension finite-step Newton proposals. Nott and Leonte (2004) extend Gamerman's (1997) method for generating MH proposals in a generalized linear model (GLM) to the variable selection case. Villani et al. (2008) extend the algorithm to a general setting not restricted to the exponential family. We first treat the problem without variable selection. The algorithm in Villani et al. (2008) only requires that

the posterior density can be written as

$$(3.1) \quad p(\beta|y) \propto p(y|\beta)p(\beta) = \prod_{i=1}^n p(y_i|\varphi_i)p(\beta),$$

where $\varphi_i = x_i'\beta$ and x_i is a covariate vector for the i th observation. Note that $p(\beta|y)$ may be a conditional posterior density and the algorithm can then be used as a step in a Metropolis-within-Gibbs algorithm. The full conditional posteriors for blocks 1-4 in Section 3.1 are clearly all of the form in (3.1). Newton's method can be used to iterate R steps from the current point β_c in the MCMC sampling toward the mode of $p(\beta|y)$, to obtain $\hat{\beta}$ and the Hessian at $\hat{\beta}$. Note that $\hat{\beta}$ may not be the mode but is typically close to it already after a few Newton iterations, so setting $R = 1, 2$ or 3 is usually sufficient. This makes the algorithm fast, especially when the gradient and Hessian are available in closed form, which is the case here, see Appendix A.

Having obtained good approximations of the posterior mode and covariance matrix from the Newton iterations, the proposal β_p is now drawn from the multivariate t -distribution with $g > 2$ degrees of freedom:

$$\beta_p|\beta_c \sim t \left[\hat{\beta}, - \left(\frac{\partial^2 \ln p(\beta|y)}{\partial \beta \partial \beta'} \right)^{-1} \Big|_{\beta=\hat{\beta}}, g \right],$$

where the second argument of the density is the covariance matrix.

In the variable selection case we propose β and \mathcal{I} simultaneously using the decomposition

$$g(\beta_p, \mathcal{I}_p|\beta_c, \mathcal{I}_c) = g_1(\beta_p|\mathcal{I}_p, \beta_c)g_2(\mathcal{I}_p|\beta_c, \mathcal{I}_c),$$

where g_2 is the proposal distribution for \mathcal{I} and g_1 is the proposal density for β conditional on \mathcal{I}_p . The Metropolis-Hasting acceptance probability is

$$a[(\beta_c, \mathcal{I}_c) \rightarrow (\beta_p, \mathcal{I}_p)] = \min \left(1, \frac{p(y|\beta_p, \mathcal{I}_p)p(\beta_p|\mathcal{I}_p)p(\mathcal{I}_p)g_1(\beta_c|\mathcal{I}_c, \beta_p)g_2(\mathcal{I}_c|\beta_p, \mathcal{I}_p)}{p(y|\beta_c, \mathcal{I}_c)p(\beta_c|\mathcal{I}_c)p(\mathcal{I}_c)g_1(\beta_p|\mathcal{I}_p, \beta_c)g_2(\mathcal{I}_p|\beta_c, \mathcal{I}_c)} \right).$$

The proposal density at the current point $g_1(\beta_c|\mathcal{I}_c, \beta_p)$ is a multivariate t -density with mode $\tilde{\beta}$ and covariance matrix equal to the negative inverse Hessian evaluated at $\tilde{\beta}$, where $\tilde{\beta}$ is the point obtained by iterating R steps with the Newton algorithm, this time starting from β_p . A simple way to propose \mathcal{I}_p is to randomly select a small subset of \mathcal{I}_c and then always propose a change of the selected indicators. This proposal can be refined in many ways, using e.g. the adaptive scheme in Nott and Kohn (2005), where the history of \mathcal{I} -draws is used to adaptively build up a proposal for each indicator. It is important to note that β_c and β_p may now be of different dimensions, so the original Newton iterations no longer apply. We will instead generate β_p using the following generalization of Newton's method. The idea is that when the parameter vector β changes dimensions, the dimension of the functionals $\varphi_c = x'\beta_c$ and $\varphi_p = x'\beta_p$ stay the same, and the two functionals are expected to be quite close. A generalized Newton update is

$$(3.2) \quad \beta_{r+1} = A_r^{-1}(B_r\beta_r - g_r), \quad (r = 0, \dots, R-1),$$

where $\beta_0 = \beta_c$, and the dimension of β_{r+1} equals the dimension of β_p , and

$$\begin{aligned}
g_r &= X'_{r+1}d + \frac{\partial \ln p(\beta)}{\partial \beta} \\
A_r &= X'_{r+1}DX_{r+1} + \frac{\partial^2 \ln p(\beta)}{\partial \beta \partial \beta'} \\
B_r &= X'_{r+1}DX_r + \frac{\partial^2 \ln p(\beta)}{\partial \beta \partial \beta'},
\end{aligned}
\tag{3.3}$$

where d is an n -dimensional vector with gradients $\partial \ln p(y_i|\varphi_i)/\partial \varphi_i$ for each observation currently allocated to the component being updated. Similarly, D is a diagonal matrix with Hessian elements

$$\frac{\partial^2 \ln p(y_i|\varphi_i)}{\partial \varphi_i \partial \varphi_i'}$$

X_r is the matrix with the covariates that have non-zero coefficients in β_r , and all expressions are evaluated at $\beta = \beta_r$. For the prior gradient this means that $\partial \ln p(\beta)/\partial \beta$ is evaluated at β_r , including all zero parameters, and that the sub-vector conformable with β_{r+1} is extracted from the result. The same applies to the prior Hessian (which does not depend on β however, if the prior is Gaussian). Note that we only need to compute the scalar derivatives $\partial \ln p(y_i|\phi_i)/\partial \phi_i$ and $\partial^2 \ln p(y_i|\phi_i)/\partial \phi_i^2$.

After the first Newton iteration the parameter vector no longer changes dimension, and the generalized Newton algorithm in (3.2) reduces to the original Newton algorithm. The proposal density $g_1(\beta_p|\mathcal{I}_p, \beta_c)$ is again taken to be the multivariate t -density in exactly the same way as in the case without covariate selection. Once the simultaneous update of the (β, \mathcal{I}) -pair is completed, we make a final update of the non-zero parameters in β , conditional on the previously accepted \mathcal{I} , using the fixed dimension Newton algorithm.

When a parameter is restricted to be proportional across components (i.e. only the intercept differs between components), the common regression vector β appears in all K components. The updating step for the common β is of the same form as above, but d and D now contain the gradients and Hessians for *all* n observations, where each observation's gradient and Hessian is with respect to the component density that the observation is currently allocated to.

3.3. Model comparison. The key quantity in Bayesian model comparison is the marginal likelihood. The marginal likelihood is sensitive to the choice of prior, however, and this is especially true when the prior is not very informative, see e.g. Kass (1993) for a general discussion and Richardson and Green (1997) in the context of density estimation. By sacrificing a subset of the observations to update/train the vague prior we remove much of the dependence on the prior, and obtain a better assessment of the predictive performance that can be expected for future observations. To deal with the arbitrary choice of which observations to use for estimation and model evaluation, one can use B -fold cross-validation of the log predictive density score (LPDS):

$$B^{-1} \sum_{b=1}^B \ln p(\tilde{y}_b|\tilde{y}_{-b}, x),$$

where \tilde{y}_b is an n_b -dimensional vector containing the n_b observations in the b th test sample and \tilde{y}_{-b} denotes the remaining observations used for estimation. If we assume that the observations

are independent conditional on θ , then

$$p(\tilde{y}_b|\tilde{y}_{-b}, x) = \int \prod_{i \in \mathcal{T}_b} p(y_i|\theta, x_i) p(\theta|\tilde{y}_{-b}) d\theta,$$

where \mathcal{T}_b is the index set for the observations in \tilde{y}_b , and the LPDS is easily computed by averaging $\prod_{i \in \mathcal{T}_b} p(y_i|\theta, x_i)$ over the posterior draws from $p(\theta|\tilde{y}_{-b})$. This requires sampling from each of the B posteriors $p(\theta|\tilde{y}_{-b})$ for $b = 1, \dots, B$, but these MCMC runs can all be run in isolation from each other and are therefore ideal for parallel computing on widely available multi-core processors.

Cross-validation is less appealing in a time series setting, and a more natural approach is to use the most recent observations in a single test sample. Moreover, for time series data it is typically false that the observations are independent conditional on the model parameters, so that the above estimation approach cannot be used. An MCMC estimate of the LPDS of a time series can instead be based on the decomposition

$$p(y_{T+1}, \dots, y_{T+T^*} | y_1, \dots, y_T) = p(y_{T+1} | y_1, \dots, y_T) \cdots p(y_{T+T^*} | y_1, \dots, y_{T+T^*-1}),$$

with each term in the decomposition

$$p(y_t | y_1, \dots, y_{t-1}) = \int p(y_t | y_1, \dots, y_{t-1}, \theta) p(\theta | y_1, \dots, y_{t-1}) d\theta,$$

estimated from a posterior sample of θ 's based on data up to time $t - 1$. The problem is that this requires $T^* - T$ complete runs with the MCMC algorithm, one for each term in the decomposition, which is typically very time-consuming (although computer parallelism can again be exploited). In situations where T is fairly large compared to T^* , we can approximate the LPDS by computing each term $p(y_t | y_1, \dots, y_{t-1})$ using the same posterior sample based on data up to time T . We evaluate the accuracy of this approximation in the empirical application in the next section.

4. MODELING THE DISTRIBUTION OF DAILY STOCK MARKET RETURNS

4.1. S&P500 data and priors. Modeling the volatility/variability in financial data has been an highly active research area since Engle's (1982) seminal paper introduced the ARCH model (see e.g. Baillie (2006) for a survey of the field), and there are large financial markets for volatility-based instruments. Financial data, such as stock market returns, are typically heavy tailed and subject to volatility clustering, i.e. a time-varying variance that evolves in a very persistent fashion. We here model the entire distribution of daily returns from the S&P500 stock market index, $p(y_t|x_t)$, where $y_t = 100 \ln(p_t/p_{t-1})$ is the daily return at time t , p_t is the closing S&P500 index on day t , and x_t contains the covariate observations at time t . By focusing on the whole distribution of returns we are able to compute e.g. the posterior distribution of the *Value-at-Risk* (VaR), i.e. the 1% quantile of the return distribution, which is of fundamental interest to financial analysts.

We estimate the models using data from 4646 trading days between Jan 1, 1990 and May 29, 2008. The models are then evaluated out-of-sample on the subsequent 199 trading days from May 30, 2008 to March 13, 2009. The data are plotted in the upper left sub-graph of

	μ	ϕ	ν	λ
m^*	0	$\sqrt{(m_\nu^* - 2)/m_\nu^*}$	10	1
s^*	10	1	7	1

TABLE 1. The prior mean and standard deviation on the split- t parameters for the S&P500 stock return data. The prior mean of ϕ is a function of the prior mean of ν such that the variance of returns is unity as in Villani et al. (2008).

Figure 4.1, with the evaluation period marked out in red. To make the results comparable to Geweke and Keane (2007) and Villani et al. (2008), we standardize the covariates to lie in the interval $[-1, 1]$, rather than making them mean zero with unit variance.

Table 1 displays the prior hyper-parameters for the split- t parameters. The prior on ν and λ are fairly vague and the prior on μ and ϕ have been chosen to match the mean and variance in Villani et al. (2008) as closely as possible.

4.2. Models. Geweke and Keane (2007) show that a smooth mixture of homoscedastic Gaussian regressions (the so called Smoothly Mixing Regression, SMR) with two covariates outperforms the typically hard-to-beat t -GARCH(1,1) model (Bollerslev, 1987) in an out-of-sample evaluation based on the LPDS (see Section 3.3). The two covariates are the return yesterday y_{t-1} (`LastDay`) and `CloseAbs95`, a geometrically decaying average of past absolute returns

$$(1 - \rho) \sum_{s=0}^{\infty} \rho^s |y_{t-2-s}|,$$

where $\rho = 0.95$ is the discount factor. Following Geweke and Keane (2007) we assume the mean of each component to be constant since the level of the stock market returns are not expected to be predictable.

Villani et al. (2008) demonstrate that the SAGM model with its heteroscedastic components outperforms the SMR in Geweke and Keane (2007). Villani et al. (2008) also introduce seven additional covariates and show that they substantially improve the out-of-sample performance of the SAGM. We will concentrate on this nine-variable model. The seven additional covariates are: `LastWeek` and `LastMonth`, a moving average of the returns from the previous five and 20 trading days, respectively. The variable `CloseAbs80`, the same variable as `CloseAbs95` but with $\rho = 0.80$, is also added to the covariate set, and so is the square root of $(1 - \rho) \sum_{s=0}^{\infty} \rho^s y_{t-2-s}^2$, for $\rho = 0.80$ and 0.95 (`CloseSqr80` and `CloseSqr95`). Finally, Villani et al. (2008) include a measure of volatility that has been popular in the finance literature: $(1 - \rho) \sum_{s=0}^{\infty} \rho^s (\ln p_{t-1-s}^{(h)} - \ln p_{t-1-s}^{(l)})$, where $p_t^{(h)}$ and $p_t^{(l)}$ are the highest and lowest values of the S&P500 index at day t . This measure has been shown both theoretically and empirically to carry more information on the volatility than changes in closing quotes (Alizadeh, Brandt and Diebold, 2002). We consider both $\rho = 0.8$ (`MaxMin80`) and $\rho = 0.95$ (`MaxMin95`). As in Villani et al. (2008), all variables except `LastDay`, `LastWeek` and `LastMonth` enter the model in logarithmic form.

4.3. Results. We generated 30,000 draws from the posterior, and used the last 25,000 draws for inference. This is more than sufficient for convergence of the parameter estimates, the

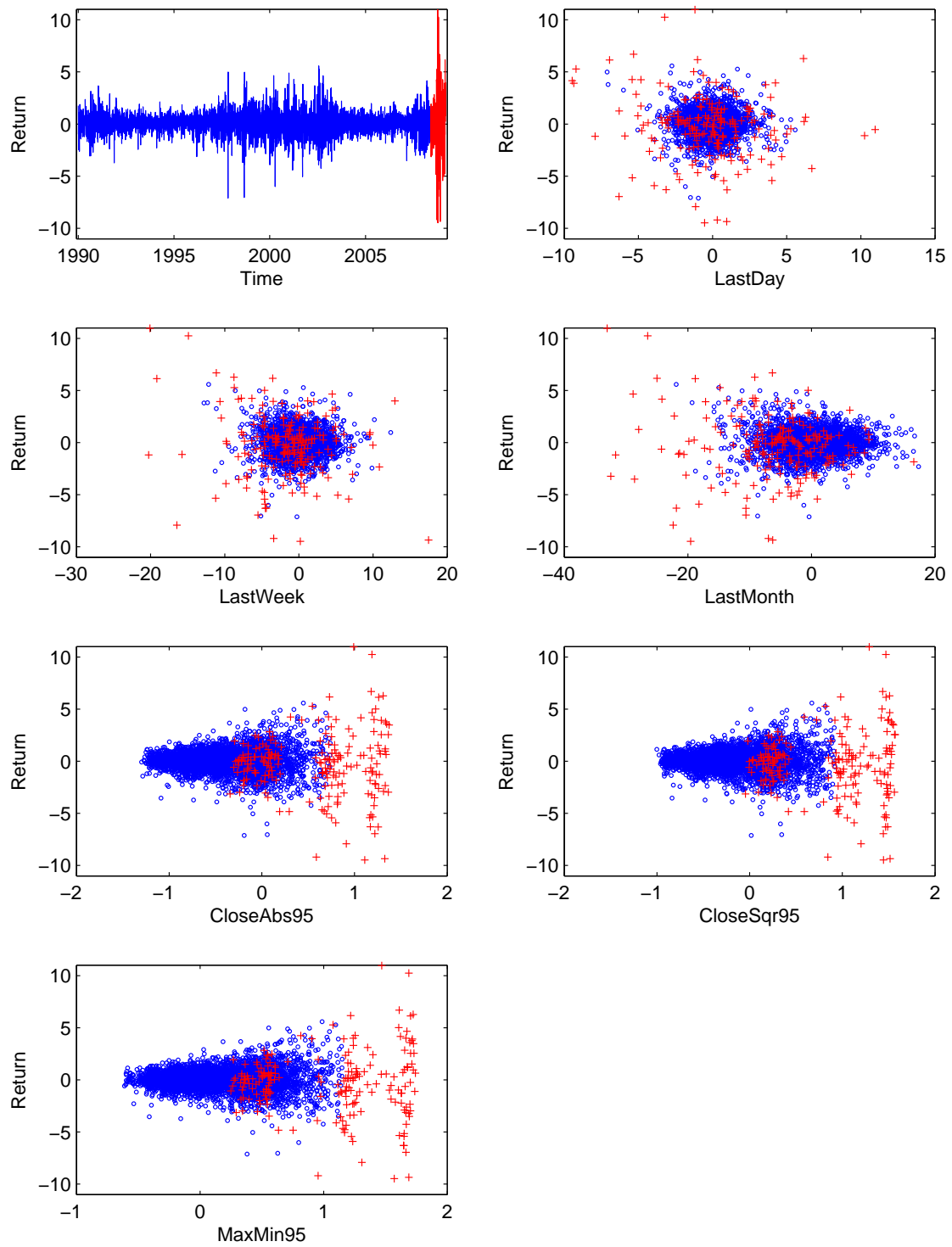


FIGURE 4.1. Graphical display of the S&P500 data from January 1, 1990 to May 29, 2008 (blue lines and circles) and May 30, 2008 to March 13, 2009 (red lines and crosses). The subgraph in the upper left position is a time series plot of Return, the other subgraphs are scatter plots of Return against a covariate.

Model	$K = 1$	$K = 2$	$K = 3$	$K = 4$	$K = 5$	Max n.s.e.
SMR	-1044.78	-638.89	-505.74	-487.11	-489.19	0.98 (3)
+ Skew	-540.91	-525.07	-513.85	-506.68	-506.13	0.82 (2)
+ DF	-544.00	-518.71	-498.93	-500.14	-494.29	0.89 (1)
+ Skew + DF	-530.86	-504.63	-498.03	-498.83	-496.87	0.88 (5)
SAGM Common	-477.73	-473.10	-473.12	-470.30	-472.86	0.26 (2)
+ Skew	-474.18	-467.29	-468.75	-467.93	-467.22	0.35 (4)
+ DF	-474.74	-472.92	-470.51	-469.40	-468.87	0.34 (4)
+ Skew + DF	-472.37	-468.92	-469.30	-466.21	-465.86	0.53 (4)
SAGM Separate		-469.21	-469.50	-470.53	-471.02	0.49 (3)
+ Skew		-468.48	-466.93	-467.48	-468.02	0.58 (4)
+ DF		-469.08	-469.24	-462.03	-467.78	0.72 (5)
+ Skew + DF		-466.84	-462.56	-462.47	-474.58	0.74 (5)
GARCH(1,1)	-479.03					
t -GARCH(1,1)	-477.39					

TABLE 2. Evaluating the out-of-sample log predictive density score (LPDS) on the 199 daily returns in the period May 30, 2008 - March 13, 2009. The posterior distribution is computed using data until May 29, 2008, and not updated thereafter, except for the two GARCH models which are based on continuously updated maximum likelihood estimates. The LPDS of the best model for a given number of components is in bold font. The last column gives the maximal numerical standard error of the LPDS for each model with the number of components for which the maximum was obtained in parenthesis. The notation for the models is such that e.g. + Skew means that covariate-dependent skewness is added to the model.

posterior inclusion probabilities and the LPDS; see also Villani et al. (2008) for details regarding convergence in the SAGM model. Three Newton steps were used for all parameters, but experiments with a single Newton step gave essentially the same numerical efficiency. The numerical efficiency of the algorithm is documented in some detail below.

Table 2 presents the LPDS evaluated on the 199 trading days from May 30, 2008 to March 13, 2009, a period covering the financial crisis with an unprecedented volatility. Figure 4.1 shows that prediction in the evaluation period is a tough test of the models because it extrapolates outside the sample used for estimation. The posterior distributions of the models are not updated during the evaluation period (see Section 3.3). With the exception of some of the more poorly fitting models, this approximation of the LPDS is quite accurate. This is documented in Villani et al. (2008) and additional evidence on this issue is provided below.

We observe from Table 2 that the SMR model does poorly, even with a large number of components, and is outperformed by the GARCH(1, 1) and t -GARCH(1, 1) models. A smooth

Model	$K = 1$	$K = 2$	$K = 3$	$K = 4$	$K = 5$
SMR	-982.02	-597.47	-498.87	-484.42	-495.66
SAGM	-477.50	-472.94	-471.28	-471.53	-469.72

TABLE 3. Evaluating the out-of-sample log predictive density score (LPDS) on the 199 daily returns in the period May 30, 2008 - March 13, 2009. The posterior distribution is updated every 10th observation throughout the evaluation sample.

mixture of homoscedastic components can generate some heteroscedasticity in-sample, but will necessarily fail in extrapolating heteroscedastic data outside the estimation sample. The subsequent rows of Table 2 show that adding covariate-dependent skewness and/or student t components (with degrees of freedom a function of covariates) to the SMR improves the LPDS substantially when the number of mixture components is small, but the SMR performs better in its standard form with Gaussian components when K is large. This reinforces the conclusion stressed in Villani et al. (2008) that having heteroscedastic components is crucial for modeling heteroscedastic data.

Table 2 also shows that SAGM is on par with the popular t -GARCH(1,1) already with a single component, outperforms it when $K \geq 2$, and is more than 7 LPDS units better than t -GARCH(1,1) at its maximum when $K = 4$. This is a substantial increase in LPDS since we are only using 199 observation in the evaluation sample.

To ensure that our short cut of keeping the posterior distribution fixed as we go through the evaluation sample does not invalidate the conclusions from the LPDS, we re-computed the LPDS for the SMR and the SAGM with a common variance function, this time updating the posterior at every tenth observation. The results are given in Table 3. A comparison of Table 2 and 3 shows that there are fairly large differences for the most poorly fitting versions of SMR, but that the LPDS values for SAGM do not change much when the posterior is updated continuously.

Table 2 shows that for the one component models, adding either covariate-dependent skewness or degrees of freedom to the SAGM model increases the LPDS by roughly 3 points, and adding them both increases the LPDs by a further 2 points. The split- t with covariate dependent scale, skewness and degrees of freedom is the best one component model, and its performance is close to that of the best SAGM model with four components. The one-component split- t (SAGM + Skew + DF) is similar to the ARCD model of Hansen (1994) which he uses to model the conditional density of the U.S. Dollar / Swiss Franc exchange rate.

If we restrict the scale, skewness and degrees of freedom to be common across components (up to a proportionality constant) we see that adding components to the split- t model improves its forecasting performance. However, we can get an even better LPDS by using separate components. Note that adding components in this case introduces as much as 41 new parameters to the model for every newly added component, and still we do not seem to over-fit even when

the number of components is fairly large. This is because of the self-adjustment mechanism emphasized in Villani et al. (2008): when an additional component is added to the mixture, the variable selection simplifies not only the new component but also the already existing components. The number of effective parameter can therefore even decrease as components are added. But there is a limit to what variable selection can do, and there are clear signs of over-fitting when $K = 5$. Also, the MCMC algorithm struggles when we use $K > 3$ separate components in the split- t model, with lower acceptable probabilities and higher risk of getting stuck in a local mode. Moreover, the split- t model with separate components has one dominant component which is very similar to the one-component model, except for the five-component model which seems to pick up a more complicated structure. We will describe the estimation results for the one-component model in detail below.

Figure 4.2 displays normalized residuals in the evaluation sample for some selected models. A normalized residual is defined as $\Phi^{-1}[F(y_t)]$, where $F(\cdot)$ is the cumulative predictive distribution, where the parameter have been integrated out with respect to the posterior distribution based on the estimation sample, so the residuals in Figure 4.2 are therefore out-of-sample. If the model is correct, the normalized residuals should be *iid* $N(0, 1)$. It is clear from Figure 4.2 that even the SMR with largest LPDS produces much to large residuals during the most volatile period. As indicated in the graph, 19.5% of the normalized residuals from the SMR(4) lie outside a 95% probability interval according to the $N(0, 1)$ reference distribution. The SAGM(1) does better than the SMR, but this model also generates to many outliers: 3.5% of the residuals are outside the 99% reference interval. The remaining four models in Figure 4.2 have rather similar seemingly homoscedastic and independent residuals, and they all have close to the right coverage. The one-component split- t model is doing remarkably well during this very difficult time period.

We now take a more detailed look at the inferences from the one-component split- t model. Table 4 presents summaries of the posterior distribution. The results from the variable selection in the scale parameter is very similar to the results for the variance function in Villani et al. (2008): the covariates `MaxMin95`, `LastWeek` and `LastMonth` have a posterior inclusion probability close to one, and all other covariates are essentially excluded from the scale parameter. There is support for some small skewness in the model, but no covariates enter λ . The degrees of freedom at the posterior mean is $\exp(2.482) = 11.96$, (assuming all other covariates at their mean) which is not very heavy tailed, but `LastWeek` enters the model with probability 0.638 and with a large negative coefficient, so the degrees of freedom is very small for the largest values of `LastWeek` (recall that `LastWeek` $\in [-1, 1]$). The last column of Table 4 gives the inefficiency factor (IF) for all parameters with inclusion probabilities larger than 0.02. It is clear that the MCMC algorithm is very efficient, almost all parameters have IFs smaller than 10. The MH acceptance probabilities for the regression coefficients in μ , ϕ , ν and λ are as high as 95%, 81%, 75% and 94%, respectively.

To explore the sensitivity to variations in the rather arbitrarily set prior parameter τ_β^2 (see Section 2.3), we compute the LPDS for the one-component split- t model using $\tau_\beta^2 = 1, 10$ and 100 (the default), obtaining an LPDS of $-472.89, -472.61$ and -472.37 , respectively. Since the LPDS is based on the posterior distribution from a large sample (unlike the marginal

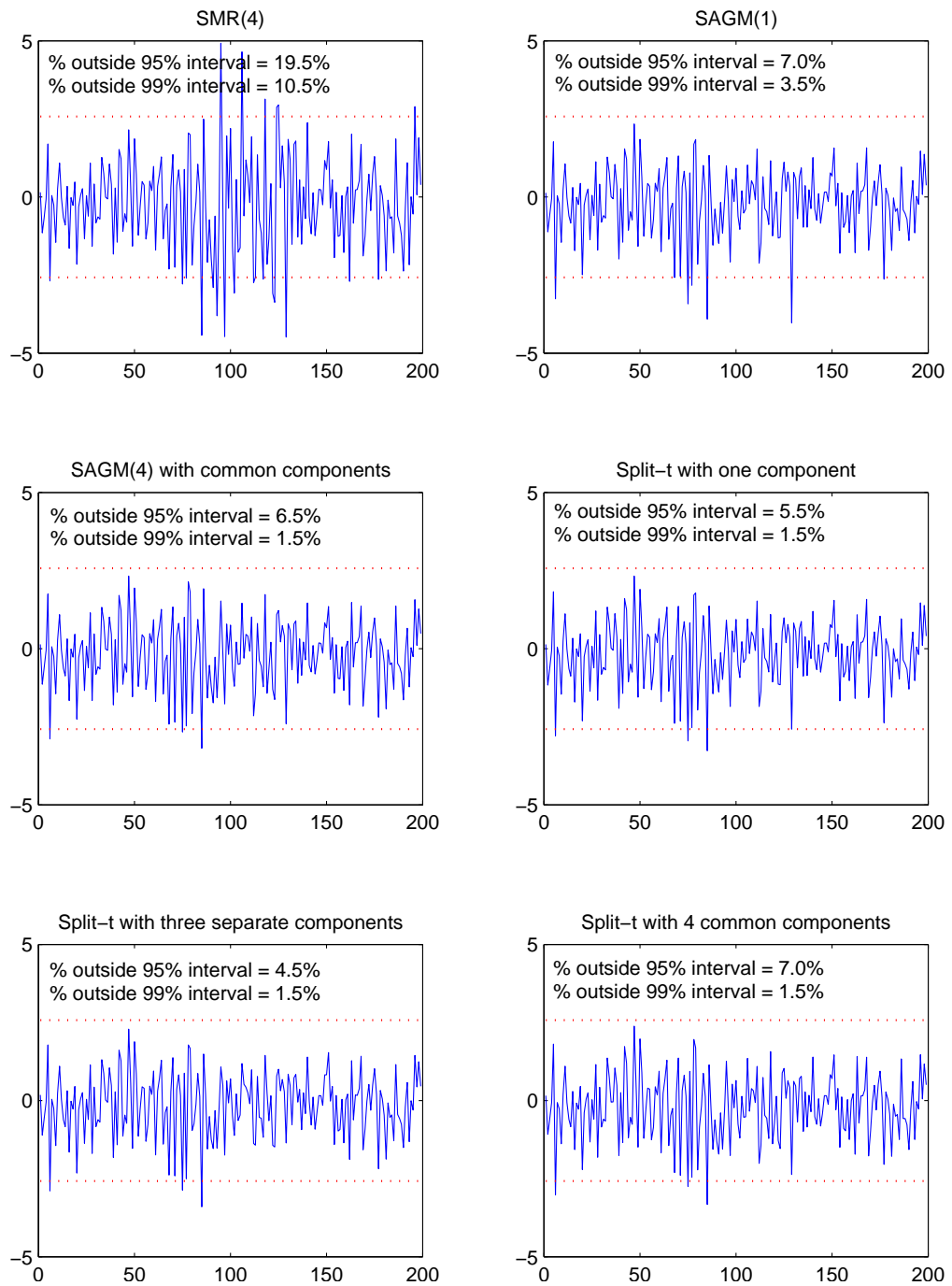


FIGURE 4.2. Plot of the 199 normalized residuals in the evaluation sample over time. The red dotted lines are the 99% probability intervals under the $N(0, 1)$ reference distribution. Each sub-graph displays the percentage of normalized residuals outside the 95% and 99% probability intervals of the $N(0, 1)$ reference distribution.

Parameters	Mean	Stdev	Post.Incl.	IF
Location μ				
Const	0.084	0.019	–	9.919
Scale ϕ				
Const	0.402	0.035	–	7.125
LastDay	-0.190	0.120	0.036	0.903
LastWeek	-0.738	0.193	0.985	18.519
LastMonth	-0.444	0.086	0.999	4.133
CloseAbs95	0.194	0.233	0.035	1.445
CloseSqr95	0.107	0.226	0.023	2.715
MaxMin95	1.124	0.086	1.000	6.012
CloseAbs80	0.097	0.153	0.013	–
CloseSqr80	0.143	0.143	0.021	–
MaxMin80	-0.022	0.200	0.017	–
Degrees of freedom ν				
Const	2.482	0.238	–	5.708
LastDay	0.504	0.997	0.112	2.899
LastWeek	-2.158	0.926	0.638	5.463
LastMonth	0.307	0.833	0.089	5.560
CloseAbs95	0.718	1.437	0.229	3.020
CloseSqr95	1.350	1.280	0.279	2.758
MaxMin95	1.130	1.488	0.222	6.564
CloseAbs80	0.035	1.205	0.101	2.789
CloseSqr80	0.363	1.211	0.112	3.330
MaxMin80	-1.672	1.172	0.254	4.178
Skewness λ				
Const	-0.104	0.033	–	10.423
LastDay	-0.159	0.140	0.027	1.170
LastWeek	-0.341	0.170	0.135	8.909
LastMonth	-0.076	0.112	0.016	–
CloseAbs95	-0.021	0.096	0.008	–
CloseSqr95	-0.003	0.108	0.006	–
MaxMin95	0.016	0.075	0.008	–
CloseAbs80	0.060	0.115	0.009	–
CloseSqr80	0.059	0.111	0.010	–
MaxMin80	0.093	0.096	0.013	–

TABLE 4. Posterior summary of the one-component split- t model. The posterior mean, standard deviation and inefficiency factors (IF) are computed conditional on a covariate being in the model. The IFs are not computed for parameters with posterior probabilities smaller than 0.02.

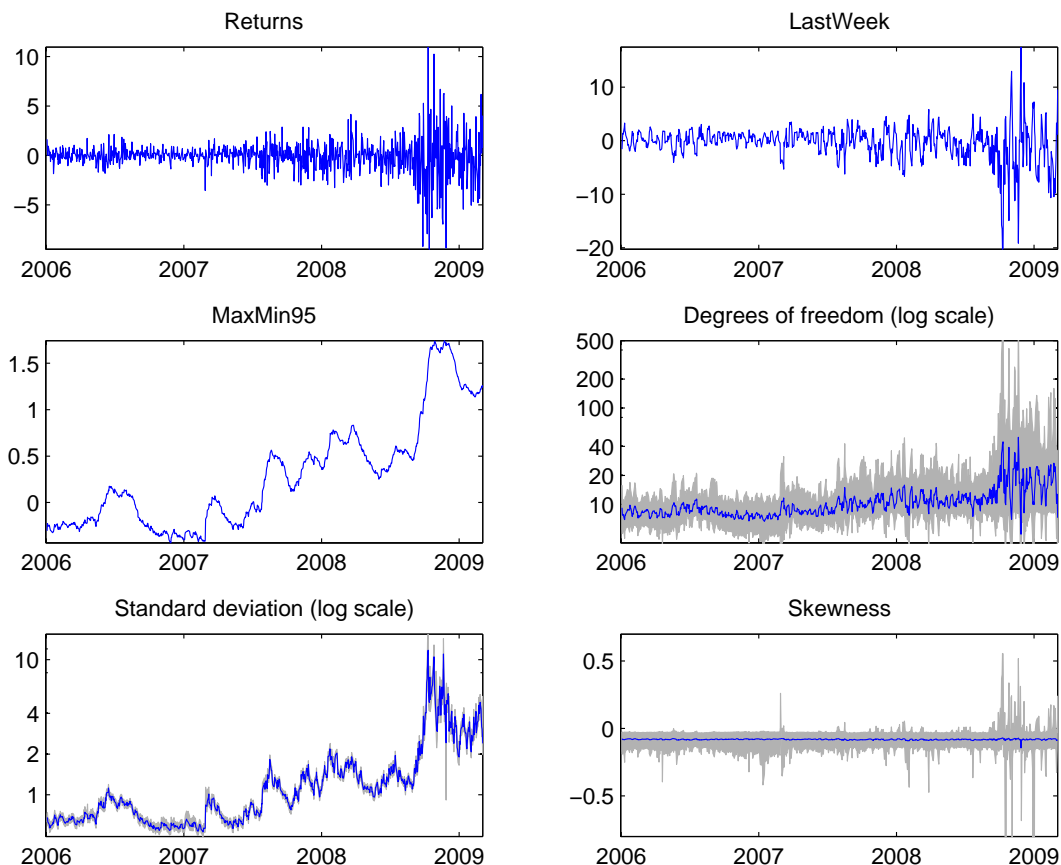


FIGURE 4.3. Time series plot of the posterior median and 95% probability intervals for some moments of the return distribution. The time series of returns and two of the key covariates are also plotted. The posterior distribution is based on the full sample up to March 13, 2009. The distribution of the standard deviation and the skewness are conditioned on $\nu > 2$ and $\nu > 3$, respectively.

likelihood which is based on the prior), this insensitivity to the prior is reassuring but not surprising. We also compare the posterior inference on the regression coefficients for the same three values of τ_β^2 . The posterior means and standard deviations are very insensitive to changes in τ_β^2 while the posterior inclusion probabilities generally decrease with τ_β^2 , but not to the extent of overturning the results about the importance of individual covariates. The effect of the prior on the inclusion probabilities is smaller for the covariates that almost certainly enter the model. As an example, the posterior inclusion probabilities for `LastDay` in ϕ is 0.290, 0.110 and 0.036 for $\tau_\beta^2 = 1, 10$ and 100, respectively, while for `MaxMin95` they are 1.000, 0.999 and 1.000 for the same three priors. Interestingly, the only significant covariate in the degrees of freedom function, `LastWeek`, has posterior inclusion probabilities of 0.66, 0.76 and 0.64 in ν for the three different values of τ_β^2 .

Finally, Figure 4.3 presents some posterior moments, such as the standard deviation and skewness, for the one-component split- t model over the latter part of the sample (including the evaluation sample). The model is estimated on all available data up March 13, 2009. Figure 4.3 shows that the median of the degrees of freedom actually increased during the most

volatile part of the financial crisis (but at the same time the scale parameter rose dramatically to bring about a very large boost in standard deviation of returns), but, during some spells, the posterior distribution of ν also has a long left tail with substantial probability mass on very small values of ν .

5. CONCLUSIONS

A general model is presented for estimating the distribution of a continuous variable conditional on a set of covariates. The model is a mixture of asymmetric student t densities with the mixture weights and all four component parameters, location, scale, degrees of freedom and skewness, being functions of covariates. We take a Bayesian approach to inference and estimate the model by an efficient MCMC simulation method. Bayesian variable selection is carried out to obtain model parsimony and guard against overfitting. The model is applied to analyse the distribution of daily stock market returns conditional on nine covariates and outperforms widely used GARCH models and other recently proposed mixture models in an out-of-sample evaluation of returns during the recent financial crisis.

6. APPENDIX - MCMC IMPLEMENTATION

To implement the MCMC algorithm we need the gradient and Hessian matrix of the conditional posteriors for each of the four split- t parameters. Since the priors on the regression coefficients in each split- t parameter is a multivariate normal density, the prior gradient and Hessian matrix are

$$\frac{\partial \ln p(\beta)}{\partial \beta} = -\Sigma_{\beta}^{-1}(\beta - \mu_{\beta}) \text{ and } \frac{\partial^2 \ln p(\beta)}{\partial \beta \partial \beta'} = -\Sigma_{\beta}^{-1}.$$

To derive the gradient and Hessian matrix with respect to the likelihood, we write the likelihood as

$$p(y|x, \mu, \phi, \nu, \lambda) = \prod_{\mathcal{S}_1} t(y|\mu, \phi, \nu) \prod_{\mathcal{S}_2} t(y|\mu, \lambda\phi, \nu),$$

where $t(y|\mu, \phi, \nu)$ denotes the student- t density

$$\frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})} \left(\frac{\nu}{\nu + \frac{(y-\mu)^2}{\phi^2}} \right)^{(\nu+1)/2},$$

\mathcal{S}_1 is the set of observations such that $y \leq \mu$ and \mathcal{S}_2 denotes the observations $y > \mu$. It is convenient to define the indicator function

$$I_{\mu} = \begin{cases} 1 & \text{if } y > \mu \\ 0 & \text{if } y \leq \mu \end{cases},$$

and $a = \lambda^{I_{\mu}}$.

The following subsections present the gradient and the Hessian for each split- t parameter.

Gradient and Hessian wrt μ

$$\begin{aligned} \frac{\partial}{\partial \mu} \ln p(y|\mu, v, \phi, \lambda) &= \frac{(1 + \nu)(y - \mu)}{\nu a^2 \phi^2 + (y - \mu)^2} \\ \frac{\partial^2}{\partial \mu^2} \ln p(y|\mu, v, \phi, \lambda) &= \frac{(1 + \nu) \left[(y - \mu)^2 - a^2 \phi^2 \nu \right]}{\left[(y - \mu)^2 + a^2 \phi^2 \nu \right]^2}. \end{aligned}$$

Gradient and Hessian wrt ϕ

$$\begin{aligned} \frac{\partial}{\partial \phi} \ln p(y|\mu, v, \phi, \lambda) &= \frac{\nu \left[(y - \mu)^2 - a^2 \phi^2 \right]}{\phi \left[(y - \mu)^2 + \nu a^2 \phi^2 \right]} \\ \frac{\partial^2}{\partial \phi^2} \ln p(y|\mu, v, \phi, \lambda) &= \frac{2\nu^2 a^2 \left[(y - \mu)^2 - a^2 \phi^2 \right]}{\left[(y - \mu)^2 + \nu a^2 \phi^2 \right]} + \frac{3\nu a^2 - \frac{\nu}{\phi^2} (y - \mu)^2}{(y - \mu)^2 + \nu a^2 \phi^2} \end{aligned}$$

Gradient and Hessian wrt ν

$$\begin{aligned}\frac{\partial}{\partial \nu} \ln p(y|\mu, v, \phi, \lambda) &= \frac{(y - \mu)^2 - v\phi^2 a}{2 \left[(y - \mu)^2 + v\phi^2 a \right]} + \frac{1}{2} \ln \left(\frac{\nu}{v + \frac{(y - \mu)^2}{\phi^2 a}} \right) \\ &\quad + \frac{1}{2} \left[\psi \left(\frac{\nu}{2} + 1 \right) - \psi \left(\frac{\nu}{2} \right) \right] \\ \frac{\partial^2}{\partial \nu^2} \ln p(y|\mu, v, \phi, \lambda) &= \frac{(y - \mu)^4 + \nu\phi^4 a}{2\nu \left((y - \mu)^2 + v\phi^2 a \right)^2} + \frac{1}{4} \left[\psi_1 \left(\frac{\nu}{2} + 1 \right) - \psi_1 \left(\frac{\nu}{2} \right) \right]\end{aligned}$$

where $\psi(\cdot)$ is the digamma function and $\psi_1(\cdot)$ is the trigamma function.

Gradient and Hessian wrt λ

$$\begin{aligned}\frac{\partial}{\partial \lambda} \ln p(y|\mu, v, \phi, \lambda) &= -\frac{1}{1 + \lambda} + \frac{(1 + v)(y - \mu)^2 I_y}{(y - \mu)^2 \lambda + v\phi^2 \lambda^3} \\ \frac{\partial^2}{\partial \lambda^2} \ln p(y|\mu, v, \phi, \lambda) &= \frac{1}{(1 + \lambda)^2} - \frac{(1 + v)(y - \mu)^2 \left[(y - \mu)^2 + 3v\phi^2 \lambda^2 \right] I_y}{\left[(y - \mu)^2 \lambda + v\phi^2 \lambda^3 \right]^2}\end{aligned}$$

REFERENCES

- [1] Abramowitz, M. and Stegun, I. A., eds. (1972), *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*, New York: Dover Publications.
- [2] Baillie, R. T. (2006). Modeling volatility, in *Palgrave Handbook of Econometrics, Vol. 1. Econometric Theory*, Mills, eds. T. C. and Patterson, K., Palgrave Macmillan, New York.
- [3] Bollerslev, T. (1987). A conditional heteroskedastic time series model for speculative prices and rates of return, *Review of Economics and Statistics*, 69, 542-547.
- [4] Celeux, G., Hurn, M. and Robert, C. P. (2000). Computational and inferential difficulties with mixture distributions, *Journal of the American Statistical Association*, **95**, 957-970.
- [5] Denison, D. G. T., Holmes, C. C., Mallick, B. K., and Smith, A. F. M. (2002). *Bayesian Methods for Nonlinear Regression and Classification*, Wiley, Chichester.
- [6] Diebolt, J. and Robert, C. P. (1994). Estimation of finite mixture distributions through Bayesian sampling, *Journal of the Royal Statistical Society B*, **56**, 163-175.
- [7] Escobar, M. D. and West, M. (1995). Bayesian density estimation and inference using mixtures, *Journal of the American Statistical Association*, **90**, 577-588.
- [8] Frühwirth-Schnatter (2006). *Finite Mixture and Markov Switching Models*, Springer, New York.
- [9] Gamerman, D. (1997). Sampling from the posterior distribution in generalized linear mixed models, *Statistics and Computing*, **7**, 57-68.
- [10] Geweke, J. (1989). Bayesian inference in econometric models using Monte Carlo integration, *Econometrica*, 57, 1317-1339.
- [11] Geweke, J. (2007). Interpretation and inference in mixture models: simple MCMC works, *Computational Statistics and Data Analysis*, **51**, 3529-3550.
- [12] Geweke, J. and Keane, M. (2007). Smoothly mixing regressions, *Journal of Econometrics*, **138**, 252-290.
- [13] Hansen, B. E. (1994). Autoregressive conditional density models, *International Economic Review*, **35**, 705-730.
- [14] Jasra, A., Holmes, C. C., and Stephens, D. A. (2005). Markov Chain Monte Carlo methods and the label switching problem in Bayesian mixture modeling, *Statistical Science*, **20**, 50-67.
- [15] Jiang W., and Tanner, M. A. (1999a). Hierarchical mixtures-of-experts for exponential family regression models: approximation and maximum likelihood estimation, *Annals of Statistics*, **27**, 987-1011.

- [16] Jiang W., and Tanner, M. A. (1999b). On the approximation rate of hierarchical mixture-of-experts for generalized linear models, *Neural Computation*, **11**, 1183-1198.
- [17] Jordan, M. and Jacobs, R. (1994). Hierarchical mixtures of experts and the EM algorithm, *Neural Computation*, **6**, 181-214.
- [18] Kass, R. E. (1993). Bayes factors in practice, *The Statistician*, **42**, 551-560.
- [19] Kohn, R., Smith, M., and Chan, D. (2001). Nonparametric regression using linear combinations of basis functions, *Statistics and Computing*, 313-322.
- [20] Nott, D. J., and Kohn, R. (2005). Adaptive sampling for Bayesian variable selection, *Biometrika*, **92**, 747-763.
- [21] Nott, D. J., and Leonte, D. (2004). Sampling schemes for Bayesian variables selection in generalized linear models, *Journal of Computational and Graphical Statistics*, **13**, 362-382.
- [22] Richardson, S. and Green, P. J. (1997). On Bayesian analysis of mixtures with an unknown number of components (with discussion), *Journal of the Royal Statistical Society, B*, **59**, 731-792.
- [23] Ruppert, D., Wand, M. P., and Carroll, R. J. (2003). *Semiparametric Regression*, Cambridge University Press, Cambridge.
- [24] Smith, M., and Kohn, R. (1996). Nonparametric regression using Bayesian variable selection, *Journal of Econometrics*, **75**, 317-344.
- [25] Villani, M., Kohn, R., and Giordani, P. (2007). Nonparametric regression density estimation using smoothly varying normal mixtures, Sveriges Riksbank Working Paper Series, no. 211. Available at www.riksbank.com.
- [26] Villani, M., Kohn, R., and Giordani, P. (2008). Regression density estimation using smooth adaptive Gaussian mixtures, forthcoming in *Journal of Econometrics*.
- [27] Wood, S., Jiang, W. and Tanner, M. A. (2002). Bayesian mixture of splines for spatially adaptive non-parametric regression, *Biometrika*, **89**, 513-528.
- [28] Zellner, A. (1986), On assessing prior distributions and Bayesian regression analysis with g-prior distributions. In *Bayesian Inference and Decision Techniques: Essays in Honor of Bruno de Finetti*, 233-243. North-Holland/Elsevier.

Earlier Working Papers:

For a complete list of Working Papers published by Sveriges Riksbank, see www.riksbank.se

Estimation of an Adaptive Stock Market Model with Heterogeneous Agents by <i>Henrik Amilon</i>	2005:177
Some Further Evidence on Interest-Rate Smoothing: The Role of Measurement Errors in the Output Gap by <i>Mikael Apel</i> and <i>Per Jansson</i>	2005:178
Bayesian Estimation of an Open Economy DSGE Model with Incomplete Pass-Through by <i>Malin Adolfson</i> , <i>Stefan Laséen</i> , <i>Jesper Lindé</i> and <i>Mattias Villani</i>	2005:179
Are Constant Interest Rate Forecasts Modest Interventions? Evidence from an Estimated Open Economy DSGE Model of the Euro Area by <i>Malin Adolfson</i> , <i>Stefan Laséen</i> , <i>Jesper Lindé</i> and <i>Mattias Villani</i>	2005:180
Inference in Vector Autoregressive Models with an Informative Prior on the Steady State by <i>Mattias Villani</i>	2005:181
Bank Mergers, Competition and Liquidity by <i>Elena Carletti</i> , <i>Philipp Hartmann</i> and <i>Giancarlo Spagnolo</i>	2005:182
Testing Near-Rationality using Detailed Survey Data by <i>Michael F. Bryan</i> and <i>Stefan Palmqvist</i>	2005:183
Exploring Interactions between Real Activity and the Financial Stance by <i>Tor Jacobson</i> , <i>Jesper Lindé</i> and <i>Kasper Roszbach</i>	2005:184
Two-Sided Network Effects, Bank Interchange Fees, and the Allocation of Fixed Costs by <i>Mats A. Bergman</i>	2005:185
Trade Deficits in the Baltic States: How Long Will the Party Last? by <i>Rudolfs Bems</i> and <i>Kristian Jönsson</i>	2005:186
Real Exchange Rate and Consumption Fluctuations following Trade Liberalization by <i>Kristian Jönsson</i>	2005:187
Modern Forecasting Models in Action: Improving Macroeconomic Analyses at Central Banks by <i>Malin Adolfson</i> , <i>Michael K. Andersson</i> , <i>Jesper Lindé</i> , <i>Mattias Villani</i> and <i>Anders Vredin</i>	2005:188
Bayesian Inference of General Linear Restrictions on the Cointegration Space by <i>Mattias Villani</i>	2005:189
Forecasting Performance of an Open Economy Dynamic Stochastic General Equilibrium Model by <i>Malin Adolfson</i> , <i>Stefan Laséen</i> , <i>Jesper Lindé</i> and <i>Mattias Villani</i>	2005:190
Forecast Combination and Model Averaging using Predictive Measures by <i>Jana Eklund</i> and <i>Sune Karlsson</i>	2005:191
Swedish Intervention and the Krona Float, 1993-2002 by <i>Owen F. Humpage</i> and <i>Javiera Ragnartz</i>	2006:192
A Simultaneous Model of the Swedish Krona, the US Dollar and the Euro by <i>Hans Lindblad</i> and <i>Peter Sellin</i>	2006:193
Testing Theories of Job Creation: Does Supply Create Its Own Demand? by <i>Mikael Carlsson</i> , <i>Stefan Eriksson</i> and <i>Nils Gottfries</i>	2006:194
Down or Out: Assessing The Welfare Costs of Household Investment Mistakes by <i>Laurent E. Calvet</i> , <i>John Y. Campbell</i> and <i>Paolo Sodini</i>	2006:195
Efficient Bayesian Inference for Multiple Change-Point and Mixture Innovation Models by <i>Paolo Giordani</i> and <i>Robert Kohn</i>	2006:196
Derivation and Estimation of a New Keynesian Phillips Curve in a Small Open Economy by <i>Karolina Holmberg</i>	2006:197
Technology Shocks and the Labour-Input Response: Evidence from Firm-Level Data by <i>Mikael Carlsson</i> and <i>Jon Smedsaas</i>	2006:198
Monetary Policy and Staggered Wage Bargaining when Prices are Sticky by <i>Mikael Carlsson</i> and <i>Andreas Westermark</i>	2006:199
The Swedish External Position and the Krona by <i>Philip R. Lane</i>	2006:200
Price Setting Transactions and the Role of Denominating Currency in FX Markets by <i>Richard Friberg</i> and <i>Fredrik Wilander</i>	2007:201
The geography of asset holdings: Evidence from Sweden by <i>Nicolas Coeurdacier</i> and <i>Philippe Martin</i>	2007:202
Evaluating An Estimated New Keynesian Small Open Economy Model by <i>Malin Adolfson</i> , <i>Stefan Laséen</i> , <i>Jesper Lindé</i> and <i>Mattias Villani</i>	2007:203
The Use of Cash and the Size of the Shadow Economy in Sweden by <i>Gabriela Guibourg</i> and <i>Björn Segendorf</i>	2007:204
Bank supervision Russian style: Evidence of conflicts between micro- and macro-prudential concerns by <i>Sophie Claeys</i> and <i>Koen Schoors</i>	2007:205

Optimal Monetary Policy under Downward Nominal Wage Rigidity by <i>Mikael Carlsson</i> and <i>Andreas Westermark</i>	2007:206
Financial Structure, Managerial Compensation and Monitoring by <i>Vittoria Cerasi</i> and <i>Sonja Daltung</i>	2007:207
Financial Frictions, Investment and Tobin's q by <i>Guido Lorenzoni</i> and <i>Karl Walentin</i>	2007:208
Sticky Information vs. Sticky Prices: A Horse Race in a DSGE Framework by <i>Mathias Trabandt</i>	2007:209
Acquisition versus greenfield: The impact of the mode of foreign bank entry on information and bank lending rates by <i>Sophie Claeys</i> and <i>Christa Hainz</i>	2007:210
Nonparametric Regression Density Estimation Using Smoothly Varying Normal Mixtures by <i>Mattias Villani</i> , <i>Robert Kohn</i> and <i>Paolo Giordani</i>	2007:211
The Costs of Paying – Private and Social Costs of Cash and Card by <i>Mats Bergman</i> , <i>Gabriella Guibourg</i> and <i>Björn Segendorf</i>	2007:212
Using a New Open Economy Macroeconomics model to make real nominal exchange rate forecasts by <i>Peter Sellin</i>	2007:213
Introducing Financial Frictions and Unemployment into a Small Open Economy Model by <i>Lawrence J. Christiano</i> , <i>Mathias Trabandt</i> and <i>Karl Walentin</i>	2007:214
Earnings Inequality and the Equity Premium by <i>Karl Walentin</i>	2007:215
Bayesian forecast combination for VAR models by <i>Michael K Andersson</i> and <i>Sune Karlsson</i>	2007:216
Do Central Banks React to House Prices? by <i>Daria Finocchiaro</i> and <i>Virginia Queijo von Heideken</i>	2007:217
The Riksbank's Forecasting Performance by <i>Michael K. Andersson</i> , <i>Gustav Karlsson</i> and <i>Josef Svensson</i>	2007:218
Macroeconomic Impact on Expected Default Frequency by <i>Per Åsberg</i> and <i>Hovick Shahnazarian</i>	2008:219
Monetary Policy Regimes and the Volatility of Long-Term Interest Rates by <i>Virginia Queijo von Heideken</i>	2008:220
Governing the Governors: A Clinical Study of Central Banks by <i>Lars Frisell</i> , <i>Kasper Roszbach</i> and <i>Giancarlo Spagnolo</i>	2008:221
The Monetary Policy Decision-Making Process and the Term Structure of Interest Rates by <i>Hans Dillén</i>	2008:222
How Important are Financial Frictions in the U.S. and the Euro Area by <i>Virginia Queijo von Heideken</i>	2008:223
Block Kalman filtering for large-scale DSGE models by <i>Ingvar Strid</i> and <i>Karl Walentin</i>	2008:224
Optimal Monetary Policy in an Operational Medium-Sized DSGE Model by <i>Malin Adolfson</i> , <i>Stefan Laséen</i> , <i>Jesper Lindé</i> and <i>Lars E.O. Svensson</i>	2008:225
Firm Default and Aggregate Fluctuations by <i>Tor Jacobson</i> , <i>Rikard Kindell</i> , <i>Jesper Lindé</i> and <i>Kasper Roszbach</i>	2008:226
Re-Evaluating Swedish Membership in EMU: Evidence from an Estimated Model by <i>Ulf Söderström</i>	2008:227
The Effect of Cash Flow on Investment: An Empirical Test of the Balance Sheet Channel by <i>Ola Melander</i>	2009:228
Expectation Driven Business Cycles with Limited Enforcement by <i>Karl Walentin</i>	2009:229
Effects of Organizational Change on Firm Productivity by <i>Christina Håkanson</i>	2009:230
Evaluating Microfoundations for Aggregate Price Rigidities: Evidence from Matched Firm-Level Data on Product Prices and Unit Labor Cost by <i>Mikael Carlsson</i> and <i>Oskar Nordström Skans</i>	2009:231
Monetary Policy Trade-Offs in an Estimated Open-Economy DSGE Model by <i>Malin Adolfson</i> , <i>Stefan Laséen</i> , <i>Jesper Lindé</i> and <i>Lars E.O. Svensson</i>	2009:232



Sveriges Riksbank

Visiting address: Brunkebergs torg 11

Mail address: se-103 37 Stockholm

Website: www.riksbank.se

Telephone: +46 8 787 00 00, Fax: +46 8 21 05 31

E-mail: registratorn@riksbank.se