# Model Complexity and Out-of-Sample Performance: Evidence from S&P 500 Index Returns

Andreas Kaeck<sup>§</sup> Paulo Rodrigues<sup>‡</sup> Norman J. Seeger<sup>†</sup>

<sup>a</sup>University of Sussex, United Kingdom, E-mail: <u>a.kaeck@sussex.ac.uk</u> <sup>b</sup>Maastricht University, The Netherlands, E-mail: <u>p.rodrigues@maastrichtuniversity.nl</u> <sup>c</sup>Corresponding author: VU University Amsterdam, De Boelelaan 1105, 1081 HV Amsterdam, The Netherlands, E-mail: <u>n.j.seeger@vu.nl</u>, Phone: +31 20 598 1512

 $<sup>^{\</sup>mbox{\sc subscript{subccript{subscript{su$ 

 $<sup>^{\</sup>ddagger}$ Maastricht University, The Netherlands, E-mail: p.rodrigues@maastrichtuniversity.nl, Phone: +31 43 388 36 33

 $<sup>^{\</sup>dagger} \rm Corresponding author. VU University Amsterdam, The Netherlands, E-mail: n.j.seeger@vu.nl, Phone: +31 20 59 81512$ 

Abstract. We apply a range of *out-of-sample* specification tests to more than forty competing stochastic volatility models to address how model complexity affects out-of-sample performance. Using daily S&P 500 index returns, model confidence set estimations provide strong evidence that the most important model feature is the non-affinity of the variance process. Despite testing alternative specifications during the turbulent market regime of the global financial crisis of 2008, we find no evidence that either finite- or infinite-activity jump models or other previously proposed model extensions improve the out-of-sample performance further. Applications to Value-at-Risk demonstrate the economic significance of our results. Furthermore, the out-of-sample results suggest that standard jump diffusion models are misspecified.

**Key Words:** Out-of-sample specification tests; jump-diffusion models; Lévyjump models; non-affine variance models; forecasting

JEL Classifications: G12; G15; C53

## 1 Introduction

In this paper, we analyze continuous-time and discrete-time models for S&P 500 index returns to study the relationship between model complexity and out-of-sample performance. The study of time-series dynamics of major stock market indices, such as the S&P 500, has previously attracted a large number of empirical studies, see e.g. Eraker *et al.* (2003), Christoffersen *et al.* (2010), Bates (2012), or Kou *et al.* (2013). One hotly debated topic, for instance, is the statistical and economic role of sudden price jumps, as the accurate modeling of tail events is of the utmost importance for many risk management and option pricing applications. To this end, Bates (2012) studies the crash risk in the US stock market using daily S&P 500 index returns from 1929 until 2010. Other model features have been highlighted in the literature, and applied research today is faced with the challenge of selecting model dynamics from a huge number of alternative specifications, including non-affine variance processes, multi-factor variance specifications, finite or infinite activity jump processes in discrete or continuous-time modeling frameworks.

Despite the importance of this research area, most papers in the continuous-time literature focus on *in-sample* specification tests or assess model performance by studying option price data. In this paper, we diverge from this approach and provide a range of different *out-of-sample* performance tests. In-sample studies are very helpful to learn about the structural building blocks required to produce stylized facts in the data. However, eventually the out-of-sample performance of a model is crucial for market participants using such a model in finance applications that are affected by uncertain future market scenarios. Our main aim is to understand to what extent the superior performance of sophisticated stochastic models prevails when they are applied outside their estimation period. To this end, we first estimate more than forty different stochastic models and encompass the most widely used model features in the continuous-time literature. Various model specification tests are then applied to an out-of-sample period of S&P 500 index returns, including the turbulent market regime during the onset of the financial market crisis in 2008. To the best of our knowledge, this paper is the first to provide comprehensive out-of-sample evidence for a large set of stochastic models.<sup>1</sup>

The estimation of continuous-time models is challenging, and a range of different estimation and filtering techniques has been developed.<sup>2</sup> At least partly driven by the differences in estimation methodology, there appears to be no standard in the continuoustime literature as far as model evaluation criteria are concerned. Eraker *et al.* (2003) use Bayes Factors and in-sample QQ plots to assess the in-sample fit and provide evidence of the impact of several model features on the shape of implied volatility smiles. Bates (2012) also uses in-sample QQ plots and a comparison of in-sample unconditional distri-

<sup>&</sup>lt;sup>1</sup>Few papers consider the out-of-sample performance of continuous-time models. Yun (2014) conducts a range of density forecasting tests using affine one-factor jump-diffusion models, Shackleton *et al.* (2010) use similar model specifications. This paper differs substantially from the aforementioned papers as we focus on a much broader number of models and specification tests.

 $<sup>^{2}</sup>$ A range of alternative estimation procedures have been proposed, including simulated methods of moments approaches, approximate maximum likelihood estimation, efficient methods of moments and Bayesian MCMC estimation algorithms (see Andersen *et al.*, 2002, Eraker *et al.*, 2003, Bates, 2006 or Johannes *et al.*, 2009).

butions, in addition to implications for option pricing. Andersen *et al.* (2002) provide in-sample specification tests as well as option pricing implications. Kaeck (2013) and Ignatieva *et al.* (2015) rely on the deviance information criterion, an in-sample Bayesian fit statistic developed in Spiegelhalter *et al.* (2002). Christoffersen *et al.* (2010) provide (in-sample) QQ plots as well scatter plots of variance level changes, and conclude that affine variance processes are rejected by the data. Li *et al.* (2008) apply (in-sample) kernel density plots, QQ plots and Kolmogorov-Smirnov (KS) tests to in-sample model residuals. Kou *et al.* (2013) also use KS tests and QQ plots in addition to comparing model autocorrelation functions to those observed in the data. Szerszen (2009) provides QQ plots and Value-at-Risk (VaR) specification tests based on in-sample parameters.

Our approach differs from the aforementioned studies in that the focal point of our paper is the out-of-sample performance. Since we are dealing with a very large number of model specifications, we employ the model confidence set estimation procedure of Hansen et al. (2011). We separate out subsets of models that have a statistically indistinguishable performance according to various different out-of-sample loss functions. In doing so, we accept that one single best performing model might not exist but rather that different modeling approaches may be equally successful. First, we compare likelihood-based outof-sample fit statistics, including sequential likelihoods as proposed by Johannes et al. (2009). This allows us to detect the time-periods during which particular specifications out- or underperform. Secondly, we follow Gneiting and Ranjan (2011) in comparing models using the continuous ranked probability score (CRPS), a criterion that can be used to compare the out-of-sample forecasting performance. CRPS has the advantage that weighted versions of the statistic retain propriety, which is essential for comparing the performance in various areas of the forecasting distributions. It is often argued that jump models in particular provide a better fit to the tails of the return distribution, and weighted CRPS fit statistics are employed to study model performance in the tails as well as the center of the return distribution. Thirdly, we test the economic significance of our results by applying the VaR loss function of González-Rivera *et al.* (2004). And fourthly,

we use a range of absolute model tests suggested by Berkowitz (2001) and others.

We use daily return observations such as in Eraker *et al.* (2003), Bates (2012) and others, although there are other data sources for extracting information about the datagenerating process. The VIX and VXO is used as a proxy of variance in Bakshi et al. (2006), Chourdakis and Dotsis (2011), and Mijatovic and Schneider (2014) to learn about structures of pure stochastic volatility models. Market prices of derivatives, for example, such as S&P 500 index options or VIX options, have been used in several papers (see Bakshi et al., 1997 or Bardgett et al., 2015). The study of option markets can provide valuable insights into the dynamics of latent state variables, but at the cost of requiring further assumptions about the dynamics of the stochastic discount factor. In addition, option pricing applications are more restrictive in terms of the data-generating processes as very few tractable models exist outside the standard affine model class (see Duffie et al., 2000). Our goal is to study asset price dynamics via the a priori imposition of as few restrictions as possible, and hence we focus in this paper on the dynamics under the physical rather than the risk-neutral measure. A second interesting data source is high-frequency returns (see Abhyankar et al., 1997). To compare our results to findings in the literature, we use daily index returns rather than high frequency data. Datagenerating processes for high frequency data are substantially more complex because such models are required to cope with intradaily trading patterns such as seasonality or other market microstructure effects. Stroud and Johannes (2014) provide sophisticated model specifications that can deal with intradaily returns.

To provide reliable evidence, we include a large number of alternative specifications. Starting from the continuous-time benchmark model proposed by Heston (1993), many extensions have been proposed in the literature. One area of research has focused on Poisson jump models, such as Bates (1996), or extensions to double-jumps as in Duffie *et al.* (2000) and Eraker *et al.* (2003). Intuitively, such models allow for occasional spikes in the data (for instance the market crash of October 1987), which are captured by a finite-activity jump process. Variations of these models alter the jump size distribution or

introduce time-varying jump intensities (see Kou, 2002, Pan, 2002 or Kaeck, 2013). More recently, a number of studies have introduced models based on infinite-activity Lévy processes. Li et al. (2008), Szerszen (2009), Bates (2012) and Ornthanalai (2014) provide evidence that suggests that such a modeling approach can be advantageous.<sup>3</sup> Another strand of the literature studies multifactor variance specifications, as these support more erratic variance movements (see Chernov et al., 2003 or Kaeck and Alexander, 2012). The literature also presents convincing evidence in favor of non-affine variance dynamics (see Jones, 2003, Christoffersen et al., 2010, Mijatovic and Schneider, 2014, or Ignatieva et al., 2015), albeit often at the cost of tractability as these models do not allow for closed-form characteristic functions. Finally, discrete-time GARCH models as well as discrete-time stochastic volatility specifications provide alternative modeling frameworks; for recent surveys, we refer to Bauwens *et al.* (2006) and Andersen *et al.* (2009). In this paper, we rely on features that have previously been proposed in the literature: affine vs non-affine models, single-factor vs multi-factor specifications, diffusion models vs jump models, finite activity vs infinite activity, discrete-time vs continuous-time models. The combination of these building blocks leads to a very comprehensive set of competing models. Although not the focus of this paper, we also study some new model specifications such as non-affine time-changed Lévy models.

Our empirical tests provide two main results. First, we find that no model is able to produce out-of-sample predictions in line with the true data-generating process. Using the test statistics developed in Berkowitz (2001) we find that all models analyzed are rejected when tested on the entire out-of-sample period. Second, we find that in terms of relative model performance more parsimonious stochastic volatility models outperform models that include a jump component. This is a surprising result, since numerous papers find that jump models outperform continuous stochastic volatility models *in-sample* (see Eraker *et al.* (2003), Eraker (2004) or Ignatieva *et al.* (2015)).

<sup>&</sup>lt;sup>3</sup>Lee and Hannig (2010) and Aït-Sahalia and Jacod (2011) propose statistical tests to distinguish between finite and infinite-activity jumps in high-frequency data.

There are two possible explanations why jump models are outperformed. First, one may interpret this result as evidence for misspecification of the jump component (despite the fact that we use quite sophisticated jump modeling) and *not* as evidence against the importance of modeling jumps in equity returns. Our results may be driven by the fact that jumps are difficult to estimate and jump distributions and intensities may vary strongly over time. For instance, jump parameters may be very different during periods of crisis and this may cause model misspecification. This finding is related to results in Santa-Clara and Yan (2010) who find a weak connection between variance and the jump intensity when both processes are estimated independently.

Second, and more important, our results may provide useful insights of how the global financial crisis unfolded. High returns may either be driven by jumps or stochastic volatility. Jumps are crucial to explain a number of rare events such as the market crash of 1987 (see the discussion in Eraker *et al.* (2003)). On the contrary, periods of high market volatility may render jumps obsolete as stochastic volatility is sufficient to generate a sequence of large returns in times of prolonged high market volatility. The result of pure stochastic volatility during our out-of-sample period was sufficient to model financial crisis returns from 2007 to 2009. The distinction between how shocks are created is important for many applications in finance as rare event models may have very different implications compared to models driven by stochastic volatility. This finding is related to Stroud and Johannes (2014) who draw similar conclusions using high frequency returns.

To study these findings in more detail, we separate the out-of-sample period into two sub-samples to investigate whether jumps have a more pronounced effect during the onset of the financial market crisis from 2007 until 2009. Surprisingly, the main driver of model performance during this period is the non-affine variance process and model confidence sets include all continuous-time GARCH models, whereas affine models perform particularly poorly. Among the non-affine specifications, the simplest diffusive variance model of Heston (1993) performs best and we conclude that the specification of the variance process is far more important than the inclusion (or the distribution) of the jump process. Model simplicity also pays off during the second, calmer sub-sample from 2010 until 2014, for which simple affine and non-affine diffusions are the only models in the 10% model confidence set. Given these strong results, we investigate whether jump models at least provide superior performance in modeling the left tail of the return distribution. Using CRPS with additional weight on the left tail shows that jump models improve their relative performance, but for this loss function the model confidence set consists of all tested specifications, and none of the jump models provides superior performance. Jumps have previously been highlighted to provide a useful modeling tool for asset returns (see for instance Eraker *et al.*, 2003), our results do not imply that jumps do not exist in the data, rather we find that they are less important for out-of-sample forecasting exercises. Applications to Value at Risk confirm that for modeling the left tail, discrete-time GARCH models as well as simple stochastic volatility models without jumps are most successful.

## 2 Model Specifications

For the first model category, we assume that the log asset price  $s_t = \ln S_t$  follows a jump-diffusion process with stochastic variance and a stochastic mean reversion level as proposed by Duffie *et al.* (2000), Egloff *et al.* (2010) and others:

$$ds_t = \left(\mu_c - \frac{1}{2}v_t - \lambda_t \bar{k}\right) dt + \rho_v \sqrt{v_t} dW_t^v + \sqrt{1 - \rho_v^2} \sqrt{v_t} dW_t^s + \xi_t dN_t \tag{1}$$

$$dv_t = \kappa_v \left( m_t - v_t \right) dt + \sigma_v v_t^{\gamma} dW_t^{\nu}$$
(2)

$$dm_t = \kappa_m \left(\theta_m - m_t\right) dt + \sigma_m m_t^{\gamma} dW_t^m, \tag{3}$$

where  $\mu_c$  is the drift and  $v_t$  denotes the stochastic variance process with speed of mean reversion  $\kappa_v$  and volatility parameter  $\sigma_v$ . The stochastic mean-reversion level  $m_t$  is governed by the speed of mean reversion  $\kappa_m$ , the long-term mean-reversion level  $\theta_m$  and the diffusion parameter  $\sigma_m$ . All three Brownian motion processes  $W^v$ ,  $W^s$  and  $W^m$  are uncorrelated, and as a consequence  $\rho_v$  determines the correlation between variance innovations and returns. The parameter  $\gamma$  identifies the dependence of the diffusion functions on the level of variance and long-term variance, respectively.<sup>4</sup> For  $\gamma = \frac{1}{2}$  we obtain an extension of the standard affine specification of Heston (1993) (labeled A); for  $\gamma = 1$  we have a continuous-time GARCH (henceforth CGARCH, see Nelson, 1990) process (G); and finally, if the parameter may take any value between one half and three halves, we obtain a general CEV variance model (C). Jump events occur at random times whenever increments in the Poisson counting process are equal to one, i.e.  $dN_t = 1$ . We assume that N has a state-dependent intensity  $\lambda_t = \lambda_c + \lambda_v v_t$ , where  $\lambda_c$  is the time-independent part of the jump intensity and  $\lambda_v$  measures the dependence of the jump probability on the current variance level. The iid jump size  $\xi_t$  is normally distributed with mean  $\mu_s$  and standard deviation  $\sigma_s$ . Furthermore, we follow the standard convention that jump sizes are independent of all other stochastic variables. The jump compensator of this model is given by  $\bar{k} = \exp\left[\mu_s + \frac{1}{2}\sigma_s^2\right] - 1$ .

For alternative jump specifications, we follow Bates (2012) and focus on jump models driven by CMGY model dynamics (see Carr *et al.*, 2002) and assume that the log asset price dynamics in Equation (1) are replaced by

$$ds_t = \left(\mu_c - \frac{1}{2}v_t\right)dt + \rho_v \sqrt{v_t} dW_t^v + \sqrt{1 - \rho_v^2} \sqrt{v_t} dL_t \tag{4}$$

where  $dL_t$  is the increment of a compensated Lévy process. The logarithm of the characteristic function  $\Psi^{CMGY}(u,t) = \mathbb{E}\left[\exp\left\{uL_t^{CMGY}\right\}\right]$  of the generalized CGMY process of Carr *et al.* (2003) is given by

$$\ln \Psi^{CMGY}(u,t) = (\mu - \omega)ut + tV \left[ w_n \frac{(G+u)^{Y_n} - G^{Y_n}}{Y_n(Y_n - 1)G^{Y_n - 2}} + (1 - w_n) \frac{(M-u)^{Y_p} - M^{Y_p}}{Y_p(Y_p - 1)M^{Y_p - 2}} \right]$$

 $<sup>^4 \</sup>mathrm{For}$  simplicity, we use the same CEV parameter  $\gamma$  for both the variance and the long-term variance process.

where  $\omega$  is a normalizing constant, V is the variance per unit time and  $w_n$  determines the fraction of downward jumps. We further define

$$C_n = \frac{w_n V}{\Gamma(2 - Y_n) G^{Y_n - 2}}$$
 and  $C_p = \frac{(1 - w_n) V}{\Gamma(2 - Y_p) M^{Y_p - 2}},$ 

where  $\Gamma(z)$  denotes the gamma function. With this definition, the parameter range is restricted to  $C_n, C_p, G, M > 0$  and  $Y_p, Y_n < 2$ . For  $Y_p, Y_n < 0$  the process has finite activity, for  $Y_p, Y_n < 1$  the process has finite variation. The model with  $L_t = L_t^{CMGY}$ nests a wide range of models used in the finance literature; we follow Bates (2012) in using the parameter restrictions  $Y_n = Y_p = 1$  for the double exponential (DEXP) jump model of Kou (2002), and  $Y_n = Y_p = 0$  for the variance gamma (VG) model of Madan and Seneta (1990) for which a non time-changed version has been estimated in Li *et al.* (2008). The full CMGY is labeled YY whereas an extension to

$$\ln \Psi^{YYD}(u,t) = (1 - f_{jump}) \frac{1}{2} (u^2 - u)t + f_{jump} \ln \Psi^{CMGY}(u,t)$$

for  $f_{jump} \in [0, 1]$  is called YYD, where D indicates an additional diffusive component. For further details such as Lévy densities or normalizing constants we refer to Bates (2012).

The asset price specifications introduced in this section allow us to distinguish between a wide range of models previously employed in the literature. The main model categories as far as the jump dynamics are concerned distinguish between either no jumps (such as in Heston, 1993), finite-activity jumps (Bates, 1996 or Duffie *et al.*, 2000) and infiniteactivity jumps (Madan and Seneta, 1990, Carr *et al.*, 2002). The variance dynamics are subdivided into affine, GARCH and general non-affine CEV dynamics (Nelson, 1990, Jones, 2003 or Christoffersen *et al.*, 2010) for both one and two-factor variance models (Egloff *et al.*, 2010 or Bates, 2012). Our model setup differs substantially from previous research which has often focused on comparing models along a single-dimension. We also compare continuous-time specifications with popular discrete-time GARCH (henceforth DGARCH) models and introduce these in the robustness section for expositional ease. We provide an overview of the one-factor continuous-time models used in this paper in Table 1, two-factor versions of the models have an additional identifier MF.

[Table 1 about here.]

## **3** Econometric Methodology

#### 3.1 Model Estimation

Our econometric methodology builds on the maximum likelihood estimation proposed by Bates (2006). Under affine model specifications, let the  $\Delta$ -step ahead conditional characteristic function of the asset return  $r_{t+\Delta} = s_{t+\Delta} - s_t$  and latent state variables be given by

$$\Psi_t^{\mathcal{G}}(u_1, u_2, u_3) \equiv \mathbb{E}\left[e^{u_1 r_{t+\Delta} + u_2 v_{t+\Delta} + u_3 m_{t+\Delta}} \middle| \mathcal{G}_t\right]$$
  
= exp { $\mathcal{A}(u_1, u_2, u_3, \Delta) + \mathcal{B}(u_1, u_2, u_3, \Delta) v_t + \mathcal{C}(u_1, u_2, u_3, \Delta) m_t$ }

where  $\mathcal{G}_t = \sigma (\{S_\tau, v_\tau, m_\tau\} : \tau \leq t)$  is the  $\sigma$ -algebra (information) generated by both the asset price and the latent state variables and  $u_1, u_2, u_3 \in \mathbb{C}$  (as long as well-defined). The functional form of the complex-valued functions  $\mathcal{A}$ ,  $\mathcal{B}$  and  $\mathcal{C}$  follows from Duffie *et al.* (2000). These functions satisfy ODEs that can be solved explicitly for one-factor affine variance specifications.<sup>5</sup> For two-factor models, we use the numerical algorithm developed in Bates (2012) to solve for  $\mathcal{A}$ ,  $\mathcal{B}$  and  $\mathcal{C}$ . This approach yields approximations that are highly accurate for applications such as ours.

To describe the filtering algorithm, the joint characteristic function of the latent state

<sup>&</sup>lt;sup>5</sup>We refer to Heston (1993), Bates (1996), Pan (2002) or Bates (2012) for the exact functional form of these functions.

variables (given the information generated by the asset returns) is defined as

$$\Lambda_t(u_2, u_3) \equiv \mathbb{E}\left[ \left. e^{u_2 v_t + u_3 m_t} \right| \mathcal{Y}_t \right]$$

where  $\mathcal{Y}_t = \sigma(\{S_\tau\} : \tau \leq t)$  is the information generated by observing the asset price only. By the law of iterated conditioning, it follows that

$$\Psi_t^{\mathcal{Y}}(u_1, u_2, u_3) \equiv \mathbb{E} \left[ e^{u_1 r_{t+\Delta} + u_2 v_{t+\Delta} + u_3 m_{t+\Delta}} \middle| \mathcal{Y}_t \right]$$
$$= e^{\mathcal{A}(u_1, u_2, u_3, \Delta)} \Lambda_t(\mathcal{B}(u_1, u_2, u_3, \Delta), \mathcal{C}(u_1, u_2, u_3, \Delta)).$$

As a result, standard Fourier inversion methods provide the probability of a return observation conditional on all past returns:

$$p(r_{t+\Delta}|\mathcal{Y}_t) = \frac{1}{2\pi} \int_{\mathbb{R}} e^{\imath u r_{t+\Delta}} \Psi_t^{\mathcal{Y}}(\imath u, 0, 0) \, du,$$
(5)

where i is the imaginary unit. We apply this numerical procedure to calculate the loglikelihood for the different model specifications. The last step in the filtering algorithm provides the update of  $\Lambda_t$ , and is given by<sup>6</sup>

$$\Lambda_{t+\Delta}(u_2, u_3) = \frac{1}{2\pi p\left(r_{t+\Delta} \mid \mathcal{Y}_t\right)} \int_{\mathbb{R}} e^{i u r_{t+\Delta} + \mathcal{A}(i u, u_2, u_3, \Delta)} \Lambda_t(\mathcal{B}(i u, u_2, u_3, \Delta), \mathcal{C}(i u, u_2, u_3, \Delta)) \, du.$$

To start the procedure  $\Lambda_0(u_2, u_3)$  is set to the unconditional characteristic function.<sup>7</sup>

Non-affine model specifications lack closed-form characteristic functions and hence the method described above cannot be directly applied. We estimate non-affine models by locally approximating them with an affine model specification. More specifically, we approximate the one-day ahead characteristic function by plugging  $\sigma_v^a = \sigma_v \mathbb{E} \left[ v_t^{\gamma-0.5} | \mathcal{Y}_t \right]$  into the respective affine characteristic function. Compared to the standard Euler dis-

 $<sup>^{6}</sup>$ We use this characteristic function to implement a moment-matching procedure, see Bates (1996).

<sup>&</sup>lt;sup>7</sup>As suggested in Bates (2006) the numerical stability of the integrals is improved by calculating the Fourier transform of a "shifted" density and then numerically invert this function. We refer to the appendix of Bates (2006) for more details on this procedure.

cretization applied in the literature (see Eraker *et al.*, 2003), such approximation is likely to be negligible over small time-steps because compared to an Euler discretization (which works well in practice, see Eraker *et al.*, 2003 or Li *et al.*, 2008), only part of the variance dynamics are kept constant. In Appendix A we provide simulation evidence to substantiate this claim and show that our estimation routine can accurately estimate affine and non-affine model parameters.

#### 3.2 Model Confidence Sets

Our empirical results include a large number of models and hence pairwise model comparisons provide only limited insight. To allow for multiple comparisons, we employ the Model Confidence Set (MCS) procedure proposed by Hansen *et al.* (2011). A MCS is defined as a set that contains the best model(s) from a collection of competing models, say  $\mathcal{M}^0$ , with a user-specified level of confidence  $(1 - \alpha)$ , where  $\alpha$  denotes the significance level (typically 10% and 25%).<sup>8</sup> The *best* models are identified based on a user-specified criterion that quantifies the relative performance of the models. Various such criteria are introduced below. A desirable property of the MCS procedure is that it acknowledges the informativeness of the data. Whereas informative data lead to the MCS containing only a few models (or even just one model), less informative data result in the MCS containing more or potentially even all models. The MCS procedure does not make a statement about which model is the true model, as performance is assessed relative to other competing models.

To fix notation, let the competing models in  $\mathcal{M}^0$  be indexed by  $i = 1, \ldots, m_0$ , with  $m_0$ denoting the number of models in  $\mathcal{M}^0$ . A user-specified loss function  $L_{i,t}$  measures the performance of each model i at time t, and the relative performance between model i and j is defined as  $d_{ij,t} \equiv L_{i,t} - L_{j,t}$  for all  $i, j \in \mathcal{M}^0$ . The expected loss of model i is defined

 $<sup>^{8}</sup>$ This interpretation is analogous to that of a classical confidence interval, hence the MCS contains the best models with a chosen confidence level. This procedure does not necessarily thereby identify one *best* model, as the MCS might consist of several models that are not statistically superior to one another.

as  $\mu_{ij} \equiv \mathcal{E}(d_{ij,t})$  according to which models are ranked, hence model *i* is preferred to *j* if  $\mu_{ij} < 0$ . The set of superior models is defined as  $\mathcal{M}^* \equiv \{i \in \mathcal{M}^0 : \mu_{ij} \leq 0 \quad \forall j \in \mathcal{M}^0\}.$ 

The objective of the MCS procedure is to determine  $\mathcal{M}^*$ . To estimate  $\mathcal{M}^*$ , candidate models are evaluated using an equivalence test  $\delta_{\mathcal{M}}$  and inferior models are subsequently removed from the initial model set based on the elimination rule  $e_{\mathcal{M}}$ . That is, a series of iterative hypothesis tests is performed, testing at each step the hypothesis

$$H_{0,\mathcal{M}}: \mu_{ij} = 0 \quad \forall i, j \in \mathcal{M}, \tag{6}$$

where  $\mathcal{M} \subset \mathcal{M}^0$  and the alternative hypothesis,  $H_{A,\mathcal{M}}$ , is given by  $\mu_{ij} \neq 0$  for some  $i, j \in \mathcal{M}$ . The equivalence test  $\delta_{\mathcal{M}}$  is used to test  $H_{0,\mathcal{M}}$  for all  $\mathcal{M} \subset \mathcal{M}^0$ . As long as the hypothesis is rejected, the elimination rule  $e_{\mathcal{M}}$  is applied to determine the most inferior model of  $\mathcal{M}$  which is then eliminated from  $\mathcal{M}$  and by this means a sequence of sets  $\mathcal{M}^0 = \mathcal{M}_1 \supset \mathcal{M}_2 \supset \cdots \supset \mathcal{M}_{m_0}$  is defined, where  $\mathcal{M}_i = \{e_{\mathcal{M}_i}, \ldots, e_{\mathcal{M}_{m_0}}\}$ . The procedure is repeated until  $H_{0,\mathcal{M}}$  cannot be rejected any more. We call the set of all surviving models  $\widehat{\mathcal{M}}^*_{1-\alpha}$ , the model confidence set with confidence level  $(1 - \alpha)$ .

Analogous to classical statistical inference, MCS *p*-values are defined as follows:  $P_{H_{0,\mathcal{M}_i}}$ denotes the *p*-value related to hypothesis  $H_{0,\mathcal{M}_i}$ . The *p*-value  $P_{H_{0,\mathcal{M}_i}}$  is calculated as  $1 - F_i(t_i)$  for  $F_i(t_i)$  being the cdf of the *i*-th test statistic  $t_i$ . A large value for the test statistic leads to small values for  $P_{H_{0,\mathcal{M}_i}}$  with the interpretation that the hypothesis  $H_{0,\mathcal{M}_i}$ , that all models in  $\mathcal{M}_i$  are equal, is likely to be statistically rejected. The MCS *p*-value for the model determined by elimination rule  $e_{\mathcal{M}_j}$  is calculated using  $\hat{p}_{e_{\mathcal{M}_j}} = \max_{i \leq j} P_{H_{0,\mathcal{M}_i}}$ . This makes it easy to determine whether a model belongs to  $\widehat{\mathcal{M}}^*$  or not, as model *i* is an element of  $\widehat{\mathcal{M}}^*_{1-\alpha}$  for a given significance level  $\alpha$  if  $\hat{p}_{e_{\mathcal{M}_j}} \geq \alpha$ . Therefore, the MCS-*p*-value is interpreted such that a model with a small *p*-value being unlikely to be a member of  $\mathcal{M}^*$ .

Specifying equivalence tests and elimination rules requires the choice of a loss function by which model performance is assessed. We use several different loss functions which are defined in Section 3.3 below. To test the performance of model *i* against alternative model specifications, Hansen *et al.* (2011) propose using a multiple *t*-statistics approach based on the test statistic  $T_{R,\mathcal{M}} = \max_{i,j\in\mathcal{M}} |t_{ij}|$ , with  $t_{ij} = \overline{d}_{ij}/\sqrt{\widehat{var}(\overline{d}_{ij})}$  and with  $\overline{d}_{ij} = T^{-1} \sum_{t=1}^{T} d_{ij,t}$ . Since the distribution of the test statistic is non-standard, a bootstrap algorithm is used to estimate the MCS *p*-values (see appendix of Hansen *et al.*, 2011). The natural elimination rule is then given as  $e_{R,\mathcal{M}} = \max_{i\in\mathcal{M}} \sup_{j\in\mathcal{M}} |t_{ij}|$ , i.e., in case of rejection of the null hypothesis, the rule eliminates the model that contributes most to the test statistic.

#### 3.3 Loss Functions

#### 3.3.1 Predictive Likelihood

The first loss function employed in this paper uses the predictive log-likelihood to compute the loss  $L_{i,t}$  (see Amisano and Giacomini, 2007, Bao *et al.*, 2007 or Wilhelmsson, 2013). Let  $f_{i,t}$  denote the predictive density of model *i* from time  $t - \Delta$  to *t*. The relative performance between two models over time is then given by  $\bar{d}_{ij} = T^{-1} \sum_t -\ln(f_{i,t}/f_{j,t})$ , where *T* denotes the number of observations. The minus sign in front of the logarithm converts the log-likelihoods into a loss function, hence model *i* is preferred over model *j* if  $\bar{d}_{ij}$  is negative.

#### 3.3.2 Continuous-Ranked Probability Score

We use further loss functions proposed in Gneiting and Ranjan (2011), and in particular we focus on the continuous-ranked probability score (CRPS) which is defined as

$$CRPS(f, y) = \int_{\mathbb{R}} \left( F(z) - \mathbb{1}\left\{ y \le z \right\} \right)^2 dz,$$

where f is the forecasting density, F its corresponding cumulative distribution function and y denotes the realized outcome (the S&P 500 index return in our case).<sup>9</sup> Intuitively, this loss function measures the difference between the forecasting distribution of a model and the *optimal* forecast that would have resulted from perfect foresight. As shown by Laio and Tamea (2006), an equivalent representation of CRPS can be obtained as an integral over quantiles and is given by

$$CRPS(f,y) = 2\int_0^1 \left(\mathbb{1}\left\{y \le F^{-1}(\alpha)\right\} - \alpha\right) \left(F^{-1}(\alpha) - y\right) d\alpha,$$

where  $F^{-1}(\alpha)$  is the  $\alpha$ -quantile of the forecasting distribution.

The advantage of CRPS over other scoring rules (such as the predictive likelihood) is that this loss function can be extended such that particular areas of the distribution function are weighted more heavily while ensuring propriety of the scoring rule. Gneiting and Ranjan (2011) propose weighted versions of CRPS defined as

$$CRPS_w(f,y) = 2\int_0^1 \left(\mathbb{1}\left\{y \le F^{-1}(\alpha)\right\} - \alpha\right) \left(F^{-1}(\alpha) - y\right) w(\alpha) \, d\alpha,$$

where  $w(\alpha)$  is a non-negative weight function on the unit interval. We follow Gneiting and Ranjan (2011) and, in addition to the un-weighted CRPS, use the following weight functions :  $w(\alpha) = \alpha(1-\alpha)$  (center),  $w(\alpha) = (2\alpha - 1)^2$  (tails),  $w(\alpha) = \alpha^2$  (right tail) and  $w(\alpha) = (1-\alpha)^2$  (left tail).

To fix the notation, the average CRSP of model i is defined as

$$\overline{CRPS}_{w,i} = \frac{1}{T} \sum_{t} CRPS_w(f_{i,t}, r_t).$$

Forecasts from densities  $f_{i,t}$  are preferred over forecasts from densities  $f_{j,t}$  if  $\overline{CRPS}_{w,i} < \overline{CRPS}_{w,j}$ . With  $d_{ij,t} = CRPS_w(f_{i,t}, r_t) - CRPS_w(f_{j,t}, r_t)$  and  $\hat{\sigma}_{ij}^2 = \frac{1}{T} \sum_t d_{ij,t}^2$  it can be

<sup>&</sup>lt;sup>9</sup>We use the standard notation  $\mathbb{1}\{A\}$  for the indicator function which takes the value 1 if A is true and zero otherwise.

shown that under the null hypothesis of vanishing expected scores, the test statistic

$$t_T = \sqrt{T} \left( \overline{CRPS}_{w,i} - \overline{CRPS}_{w,j} \right) \hat{\sigma}_{ij}^{-1}$$

asymptotically follows a standard normal distribution (assuming suitable regularity conditions for which we refer to Gneiting and Ranjan, 2011). To use the CRPS loss function for calculation of model confidence sets, the average relative performance between model i and j is defined as  $\bar{d}_{ij} = \overline{CRPS}_{w,i} - \overline{CRPS}_{w,j}$ .

#### 3.3.3 Asymmetric Value-at-Risk Loss Function

The third loss function we employ is proposed by González-Rivera *et al.* (2004) and has been designed specifically for testing the predictive power of models in the context of VaR estimation. The proposed loss function is given by

$$L_{i,t}^{VaR} = (r_t - \operatorname{VaR}_{i,t}^{\alpha}) \times \left(\alpha - \mathbb{1}\left\{r_t < \operatorname{VaR}_{i,t}^{\alpha}\right\}\right)$$

where  $\operatorname{VaR}_{i,t}^{\alpha}$  is the Value at Risk at significance level  $\alpha$  for model *i* estimated at time  $t - \Delta$  for a return horizon of  $\Delta$ . The functional form of the loss function implies that deviations from VaR are weighted more heavily if  $r_t < \operatorname{VaR}_{i,t}^{\alpha}$ , which is in line with the goal of avoiding large losses. The relative performance of two models *i* and *j* is given by  $\bar{d}_{ij} = T^{-1} \sum_t \left( L_{i,t}^{VaR} - L_{j,t}^{VaR} \right).$ 

#### 4 Data

We employ daily log returns of the S&P 500 index for a period from January 2, 1987 until December 31, 2014. This data set overlaps with many previous studies such as Andersen *et al.* (2002) and Eraker *et al.* (2003). We separate the sample into an in-sample period

from 1987 until 2006 and an out-of-sample period from 2007 until 2014. With tranquil and turbulent market regimes in both sub-samples, we have an ideal testing ground for the existence of jumps and the performance of alternative volatility specifications and their relative merits for out-of-sample forecasting. We report all parameters on a yearly basis and set  $\Delta = \frac{1}{252}$ . Table 2 provides summary statistics for the whole sample period, as well as various sub-samples used in this paper.

[Table 2 about here.]

## 5 Empirical Results

In this section, we present in- and out-of-sample results for the one-factor jump-diffusion and Lévy-jump models. We focus on these models first for expositional ease, and discuss multi-factor variance models as well as discrete-time specifications in Section 6.

#### 5.1 Parameter Estimates and In-Sample Performance

[Table 3 about here.]

[Table 4 about here.]

We report parameter estimates for the one-factor jump-diffusion models in Table 3. As most models have been extensively discussed in the literature, we provide only a short interpretation of our estimation results. For the standard Heston model (SV-A) we estimate a long-term variance of 0.029 (which translates into a yearly volatility level of 17.03%) and a vol-of-vol  $\sigma_v$  of 0.435. Eraker *et al.* (2003) for instance find 14.37% for the long-term volatility and 0.3614 for the volatility diffusion parameter in their less turbulent sample which ends before the dot-com bubble bursts. Our correlation estimate of -0.681 and the speed of mean reversion (5.449) are also in line with previous findings. The CGARCH and CEV model parameter estimates portray two patterns: first, a higher  $\gamma$  value leads to a lower speed of mean reversion and secondly, a slightly increased estimate of  $\theta_v$ . Interestingly, for all model classes, the CEV parameter  $\gamma$  is statistically indistinguishable from the CGARCH specification ( $\gamma = 1$ ). Jumps across all different model specifications occur less than once a year in the time-homogeneous jump specifications, but can occur slightly more frequently when the jump probability depends on the prevailing volatility regime. The jump parameter  $\lambda_c$  in all SVSJJ models is estimated to be to zero, and hence our results indicate that time-varying jump probabilities are an important feature of S&P 500 index returns. These results confirm earlier evidence in Bates (2006) and Christoffersen *et al.* (2012). Jump sizes are relatively stable across different specifications with average means of -3% and a standard deviation between 5% and 6%.

The likelihood values at the optimal parameter set indicate that jump models substantially improve the in-sample performance, for instance we find that the log-likelihood increases by a value of between 30 to 40 from SV to SVJ. Time-varying jump probabilities provide further improvements of the log-likelihood, especially in the affine model specification. Consistent with the low parameter estimates for  $\lambda_c$  we find little evidence of an improvement of SVSJJ over SVSJ. The second very consistent result is that the CGARCH models clearly outperform affine models, whereas a free CEV parameter has only a minor effect on the performance measure.

We report the parameter estimates for the one-factor Lévy-jump models in Table 4. For the parameters that govern the stochastic variance process we find similar patterns to those for the one-factor jump-diffusion models. Long-term volatility estimates are quite stable and vary between 17 to 20%. Correlation estimates vary between -0.632 and -0.729 and are increasing slightly with an increasing  $\gamma$ . Estimates for  $\kappa_v$  and  $\theta_v$  have a similar order of magnitude, and decrease and increase respectively with an increasing  $\gamma$  parameter. All affine versions of the one-factor Lévy-jump models (SVYY, SVDEXP, SVVG, SVYYD with  $\gamma = 0.5$ ) are also estimated in Bates (2012) and the reported variance parameters are in line with our estimates: long-term volatility ( $\sqrt{\theta_v}$ ) varies between 15.3 and 17.4%, mean-reversion speed ( $\kappa$ ) varies between 3.961 and 8.318, and correlation ( $\rho$ ) varies between -0.541 and -0.674. Also the parameter estimates for the Lévy-jump models given in Bates (2012) are of the same order of magnitude as ours and show the same structural behavior across model specification. Differences in the estimates can be explained by the different sample periods used in the papers. The weighting parameter  $w_n$  varies between 0.49 and 0.88, and our estimates vary between 0.32 and 0.83. The parameter  $f_{jump}$  which gives the proportion of variance that is driven by the Levy-jump part varies from 0.253 to 0.436, whereas we find values between 0.289 and 0.338. For the SVYYD model we find that  $f_{jump}$  is close to the boundary value 1, which implies that the SVYYD model reduces to SVYY. Therefore, our data does not support an additional diffusion component for this model specification. This finding is confirmed by the log-likelihood values for the SVYYD and the SVYY model, which remain very close even for the different  $\gamma$  specifications.<sup>10</sup>

In terms of in-sample log-likelihood values, we obtain similar model rankings to those in Bates (2012). In particular, we find that the performance of DEXP models is similar to SVVG models and these are outperformed by SVYY model. The additional distribution component in the SVYYD model does not lead to further fit improvements, as log-likelihood values remain very close to the SVYY model.

We briefly discuss the in-sample performance of the models. Let  $\mathcal{L}_{tT}^m$  be the loglikelihood of model m between time t and T. Then  $\mathcal{R}_h^m = (h - t)^{-1} \left(\mathcal{L}_{th}^m - \mathcal{L}_{th}^b\right)$  for  $h = t, \ldots, T$  defines the sequential normalized difference between the log-likelihood function of model m and the benchmark model b between t and h > t. In Figure 1 we compare these relative likelihood sequences over the in-sample period as suggested by Johannes et al. (2009), with SV-A as our benchmark. From the sequential likelihood ratios it is evident that the severe market shock of 1987 plays a crucial role in distinguishing different

<sup>&</sup>lt;sup>10</sup>We have first estimated all Lévy models as in Bates (2012). Since our shorter sample period has significantly fewer large positive return outliers we found M to be unstable in the estimation and fixed the value to estimates in Bates (2012). All empirical results are robust as to whether M is fixed or estimated.

specifications, indeed all models cope with the large -23% return observed on October 19, 1987 far better than the affine SV-A model (see also the discussion in Eraker *et al.*, 2003). Jump models are slightly more successful during this extended period of market turmoil; simple non-linear variance models however also fare relatively well.

## [Figure 1 about here.]

Overall jumps improve the likelihood ratios substantially and for the in-sample period, accounting for these is more important then the choice of variance dynamics. This is evident from the fact that affine jump models out-perform non-affine pure stochastic volatility models. Levy models provide further improvements over jump-diffusion specifications, and roughly half of the difference between jump vs non-jump models results from the October 1987 period.

#### 5.2 Out-of-Sample Forecasting Performance – Log-Likelihood

In Figure 2, we present sequential likelihood ratios for the out-of-sample period which are calculated using Equation (5), fixing structural parameters to those estimated during the estimation window (as in Eraker, 2004). Interestingly, this figure highlights a striking difference from our in-sample results, as the simple SV-G model outperforms all other specifications by roughly two log-likelihood points per year.

## [Figure 2 about here.]

Jump models, although performing well in-sample, do not exhibit major improvements even over the simple affine stochastic volatility model SV-A. In addition, we find that the underperformance of jump models is gradually accumulated over the out-of-sample period rather than being the result of a single outlier. By contrast, the excess likelihood of non-linear variance models is accumulated predominantly during the outbreak of the financial crisis in 2008. We return to this finding further below.

[Table 5 about here.]

To add statistical rigor to our graphical results, in Table 5 we report model confidence set estimations using the out-of-sample negative log-likelihood as a loss function. We focus on affine and CGARCH models, and remove models of the CEV and SVSJJ class as their parameters (and out-of-sample results) are indistinguishable from other model specifications.<sup>11</sup> For the model confidence set estimation, we choose the block length of the bootstrap as follows. For each model, we estimate simple autoregressive (AR) models and determine the optimal lag length according to AIC and BIC fit criteria. We then select the bootstrap block length equal to the maximum lag length of all models in  $\mathcal{M}_0$ . It is evident from these results that the difference in out-of-sample likelihood between SV-G and all other specifications is statistically significant. We find that the MCS consists solely of SV-G at the 25% level, which provides strong evidence in favor of the simple non-affine stochastic volatility model. The first models eliminated from the initial model set are affine jump specifications. After this, jump models with CGARCH variance dynamics are excluded, and interestingly we find virtually no difference between the performance of finite and infinite-activity jump models. The final exclusion is SV-A which, although dominated by the SV-G model, performs much better than all affine jump models. Our results confirm statistically the superiority of the simple CGARCH volatility specification and the fact that the MCS is a singleton can be interpreted as strong evidence that the out-of-sample period is informative with regard to the different model features. These findings also provide the first evidence that the choice of volatility dynamics is more important than modeling jumps.

#### [Table 6 about here.]

The graphical analysis in Figure 2 indicates that there may be two distinct regimes during the out-of-sample period: a first turbulent regime during the international financial crisis (2007-2009), and a second more stable regime until the end of the sample period (2010-2014). As it appears that most of the outperformance of the SV-G model stems from the credit crisis period, we rerun the model confidence set estimation for both sub-

<sup>&</sup>lt;sup>11</sup>Results for these models are available from the authors upon request.

periods separately to understand how the model ranking is affected by different market environments. Results in Panel A of Table 6 confirm that the CGARCH model class performs significantly better during the crisis period, and the model confidence set at the 25% level consists of all CGARCH specifications (with or without jumps), whereas all affine models perform weakly. This confirms the graphical findings that the variance dynamics are very important for adequately modeling market crashes. In Panel B, we focus on the calmer sub-period and find that the model confidence set at the 10% level consists only of the two stochastic volatility models, with insignificant performance differences between SV-A and SV-G. Taken together, the out-of-sample log-likelihood tests provide evidence that simple stochastic volatility models outperform more advanced jump specifications and that the dynamics of the variance process matter, particularly during turbulent market regimes.

#### 5.3 Out-of-Sample Forecasting Performance – Continuous Ranked Probability Score

We now provide out-of-sample results for a loss function that focuses on the forecasting performance of alternative models. In addition, we aim to test whether jump specifications provide superior performance in forecasting tail events, as one advantage of jump models is that they provide additional flexibility to fit the tails of the return distribution. To this end, we follow the framework of Gneiting and Ranjan (2011) and base our assessment of the forecasting performance on the continuous ranked probability score (for formal definitions, see Section 3.3).

#### [Table 7 about here.]

We report empirical results for the unweighted CRPS tests in Table 7.<sup>12</sup> The best performing model is SV-G, in line with the empirical results for the log-likelihood loss function above. In pairwise comparisons, this specification outperforms all other compet-

<sup>&</sup>lt;sup>12</sup>For expositional ease we do not report the result for SVDEXP and SVSJ in these tables as they do not provide additional insights. We nevertheless include them in the model confidence set estimations below to ensure all empirical results are based on the same initial model set  $\mathcal{M}_0$ .

ing models at all conventional significance levels. The lowest absolute pairwise t-statistic results from the comparison with SV-A (with a t-statistic of -2.39). Given the large number of alternative models, this finding presents very strong support for non-linear variance dynamics. Furthermore, we find that diffusion models significantly outperform their jump extensions out-of-sample. In the affine model class the simple SV-A model outperforms all affine jump extensions, with the lowest significance level arising for a t-statistic of -3.28 for the comparison with the variance-gamma jump model. The same finding can be seen for the CGARCH model class where the smallest significance level results from the comparison between SV-G and SVJ-G (with a t-statistic of -2.50). Overall, CGARCH models offer a significant and very consistent improvement over affine models, with t-statistics ranging from 2.39 to 2.89 when comparing the same jump specification with either an affine or CGARCH-type variance process.

## [Table 8 about here.]

We restrict the detailed discussion of CRPS results for alternative weight functions to the left tail of the return distribution as this part of the distribution is most interesting for financial applications such as VaR. In addition, the left tail of the distribution of S&P 500 index returns benefits the most from the addition of jumps and hence weighting the left tail more heavily may uncover potential shortcomings of simple SV specifications. Our test results in Table 8 show that the model ranking is surprisingly little changed after altering the weight function. In particular, SV-G is still the overall best performing model and dominates all other specifications in pairwise model comparisons. However, jump models close the gap to simple SV models, and pairwise CRPS *t*-tests now indicate no statistically significant differences between the forecasting ability of competing model specifications. The empirical results for the forecasting tests with center, tails and right tail weight functions are available upon request, while for ease of exposition we restrict the discussion of these additional weight functions to the model confidence set estimation below.<sup>13</sup>

<sup>&</sup>lt;sup>13</sup>Results for these tests, similar to Table 8, are available upon request.

#### [Table 9 about here.]

We extend our previous findings and add MCS estimations to the pairwise model comparisons (see Table 9). Results for the different weight functions are in general supportive of the model rankings presented above and provide further strong evidence in favor of SV-G. Unsurprisingly, given the strong pairwise outperformance in Table 7, for all test statistics except for the left tail discussed above, SV-G is the only model in the 10% confidence set and it is also the best performing specification for all five test statistics. SV-G is particularly successful in the center and the right tail of the return distribution. It is also notable that the SV-A model provides a poorer performance compared to the log-likelihood loss function, where it was included in the MCS. Confirming our earlier findings, the model confidence set for the left tail includes all models, hence we are not able to distinguish between the forecasting performance in the left tail, at least as far as our out-of-sample period is concerned. Nevertheless, the SV-G model still performs best in this category, albeit at no conventional significance level.

#### [Table 10 about here.]

In Table 10, we report model confidence set results for the two distinct out-of-sample regimes (January 2007 to December 2009 and January 2010 to December 2014). The (unweighted) Gneiting-Ranjan tests confirm that during the financial crisis period models with CGARCH variance dynamics outperform affine specifications, whereas the addition of jumps does not lead to further improvements. The best-performing affine model is, as before, the simple diffusion specification, and this is the only affine model in the 25% model confidence set. By contrast, the second calmer period (January 2010 to December 2014) provides strong evidence that the SV-G and SV-A models dominate jump specifications (they are the only two models in the MCS at the 25% level) and SVJ-G also provides acceptable forecasting performance with a p-value of 0.1077. The results in Panel A and B of Table 10 provide supporting evidence that the forecasting performance in the left tail of the distribution is not improved with jumps of either finite or infinite activity as

model confidence sets include all initial models  $\mathcal{M}_0$ . For completeness, we also report model confidence sets for alternative weight functions.

#### 5.4 Out-of-Sample and Forecasting Performance – Berkowitz

Berkowitz (2001) proposes an alternative method for testing the forecasting performance, building on work from Diebold *et al.* (1998) and others. It is assumed that a forecasting model with density f is employed, whereas the true (unknown) density is given by p. It can be shown that the density of the integral transform z, defined as

$$z = \int_{-\infty}^{y} f(z) \, dz = F(y),$$

is given by  $p(F^{-1}(z))/f(F^{-1}(z))$ . Therefore, under the null hypothesis that the forecasting model is equal to the true data-generating process, the variable z is uniformly distributed and  $\tilde{z} = \Phi^{-1}(z)$  follows a standard normal distribution ( $\Phi^{-1}$  denotes the inverse cumulative distribution function of a standard normal random variable). Furthermore, it can be shown that in a time-series framework, the realizations  $\tilde{z}_t$  need to be iid. Berkowitz (2001) proposes to test this hypothesis using  $\tilde{z}_t - \mu = \rho(\tilde{z}_{t-1} - \mu) + \varepsilon_t$  and the corresponding log-likelihood function  $L(\mu, \sigma, \rho)$ . This implies three possible tests, one for iid, one for independence and one for the joint hypothesis. The likelihood ratio test statistics are given by  $LR_{ind} = -2[L(\hat{\mu}, \hat{\sigma}, 0) - L(\hat{\mu}, \hat{\sigma}, \hat{\rho})]$ ,  $LR_{iid} = -2[L(0, 1, 0) - L(\hat{\mu}, \hat{\sigma}, 0)]$  and  $LR = -2[L(0, 1, 0) - L(\hat{\mu}, \hat{\sigma}, \hat{\rho})]$ . The main advantage of this test procedure is that it provides an absolute test of the forecasting performance.

### [Table 11 about here.]

Table 11 presents results for the LR statistic as this provides the most general analysis.<sup>14</sup> The overall results are further divided into out-of-sample tests for each individual year in the out-of-sample period. Overall, we find that all models are rejected by the

<sup>&</sup>lt;sup>14</sup>Results for other statistics are available upon request.

data. This is unsurprising as most models struggle with explaining the returns during the onset of the financial crisis in 2008, a year for which the null hypothesis is rejected by all models. Similar to our findings for the relative forecasting performance, we find that the statistics for CGARCH models are, however, more than half of the value for affine models and hence provide additional confirmation of the superiority of CGARCH models during the financial crisis period. Interestingly, during all years following 2008, for all of the models, the null hypothesis cannot be rejected at the 5% level. Modeling differences are merely relevant during the crisis period, a finding that reinforces the conclusions above.

### 6 Further Results

In this section, we extend our analysis in various directions and discuss the performance of two-factor jump diffusion models (introduced in Section 2) as well as simple DGARCH models. To focus on our main findings, we report model comparisons with a representative subset of one-factor models, namely SV-A, SV-G, SVSJ-A, SVSJ-G, SVYYD-A and SVYYD-G.

#### 6.1 Multi-factor Variance Models

[Table 12 about here.]

Parameter estimates for two-factor jump-diffusion models are based on the same insample period from January 2, 1987 until December 29, 2006 and are reported in Table 12. Our main results can be summarized as follows. First, the stochastic process  $m_t$  is slowly mean-reverting (estimates for  $\kappa_m$  range between 0.516 and 1.408) and it exhibits a relatively low diffusive volatility parameter  $\sigma_m$ . Secondly, the addition of the timevarying mean reversion level significantly alters the dynamics of stochastic variance. The process  $v_t$  is now much faster mean-reverting to  $m_t$  than it is in one-factor models and it is also significantly more volatile. In the SV-A model class, for instance, the mean reversion speed  $\kappa_v$  for one- and two-factor models is 5.449 and 26.928 respectively, whereas estimates for  $\sigma_v$  increase from 0.435 to 0.633. A high value for  $\kappa_v$  implies that  $v_t$  varies erratically around the long-term variance  $m_t$ . Thirdly, we find that the estimate for  $\gamma$  is slightly higher than in one-factor models with values of between 1.057 and 1.176. This is likely due to the fact that variance itself moves more violently around  $m_t$  and a higher CEV parameter facilitates such fast-moving behavior. Jump parameter estimates are comparable to the one-factor specifications discussed above. Given our previous findings regarding the minor importance of jumps for out-of-sample forecasting, we refrain from extending the analysis to Lévy-jump models and restrict our results to jump-diffusions to capture jump-like behavior.

### [Figure 3 about here.]

The left part of Figure 3 shows in-sample sequential likelihood ratios for all two-factor jump-diffusion models (using SV-A as benchmark model). The overall evolution of these statistics is comparable to that for one-factor models; in particular, we find that jump models out-perform simple diffusion specifications and non-affine stochastic variance models also provide further improvements. The right-hand graph in Figure 3 documents out-of-sample sequential likelihood ratios. It is evident that the start of the global financial crisis in 2008 is an important time period for distinguishing the performance of alternative models, and non-affine specifications perform substantially better than affine models during this market regime. Interestingly, affine multi-factor models are particularly unsuccessful at explaining S&P 500 index returns during the crisis period. A possible explanation for this finding is that while the variance process in affine two-factor models is more erratic, its rapid mean-reverting behavior forces  $m_t$  to drive the overall variance level. Since  $v_t$  in one-factor models is more volatile than  $m_t$  it is possible that affine two-factor models are less successful at modeling more substantial variance changes. Unreported results confirm that in the affine models, the spot variance of one-factor models exceeds the variance levels of two-factor models during the peak of the financial market crisis in

#### [Table 13 about here.]

Table 13 presents out-of-sample model confidence set estimates for the negative predictive log likelihood loss function. These results complement the graphical results presented earlier and further compare the model performance of the two-factor specifications with the most successful specification of each model class of Section 5. The MCS estimates confirm that SV-A and SV-G are the best performing one-factor models, and there are only two additional two-factor models in the 25%-level confidence set, namely MF-SV-G and MF-SVSJJ-G. The best performing model is MF-SV-G, followed by SV-G which exhibits a MCS p-value of 0.9147. Although slightly less extreme than in the case of the one-factor specifications, two-factor models do not benefit from adding additional jumps to capture large outliers either, and it is more important to account for non-linear variance dynamics as affine multi-factor models perform particularly poorly. These findings suggest that using a multi-factor model with non-affine variance dynamics provides a similar performance to a simpler one-factor non-affine specification. We therefore conclude that two-factor models do not add significant gains for our out-of-sample data set. However, we do not find any evidence that more complex models lead to a deterioration in performance either.

#### [Table 14 about here.]

Table 14 provides further out-of-sample results, using the Gneiting and Ranjan (2011) test procedure with weighted CRPS test statistics as our loss function. The results for an unweighted objective function, similarly to the predictive likelihood results, suggest that non-affine model dynamics are important and that non-affine two-factor models provide similar out-of-sample performance to simple SV-G and SV-A models. As before, there is substantially less evidence in multi-factor models that jumps have a negative effect on the forecasting performance, and all non-affine MF models are included in the 25% model confidence set. As before, it proves very difficult to distinguish between the forecasting

performances in the left tail of the return distribution, where all models but MF-SV-A are included in the MCS. Non-affine models are superior at forecasting the right tail, where affine models perform poorly. While most of the attention in the literature is devoted to the left tail of the return distribution, the right tail may be of particular interest to investors with short positions.

#### 6.2 Discrete-Time GARCH Models

In order to compare our results to simpler DGARCH models, we estimate further specifications that have been found to perform well in the discrete-time literature.<sup>15</sup> Our benchmark model is given by a multi-factor GJR-GARCH model (see Glosten *et al.*, 1993), written in a form that explicitly highlights the long-term variance level:

$$r_{t+1} = \mu + \varepsilon_{t+1} = \mu + \sqrt{h_{t+1}} z_{t+1} + I_{t+1} \xi_{t+1} - \lambda \mu_j$$
(7)

$$h_{t+1} = q_{t+1} - \left(\alpha_h + \frac{1}{2}\gamma_h\right)\left(q_t + \psi\right) + \beta_h\left(h_t - q_t\right) + \left(\alpha_h + \gamma_h \mathbb{1}_{\varepsilon_t < 0}\right)\varepsilon_t^2 \tag{8}$$

$$q_{t+1} = q_q - \left(\alpha_q + \frac{1}{2}\gamma_q\right)\left(q_q + \psi\right) + \beta_q\left(q_t - q_q\right) + \left(\alpha_q + \gamma_q \mathbb{1}_{\varepsilon_t < 0}\right)\varepsilon_t^2 \tag{9}$$

where  $h_t$  is the diffusive variance,  $q_t$  is the long-term variance level,  $\alpha_h$  and  $\beta_h$  are model parameters determining the speed of mean reversion and how quickly the variance changes in response to a return shock, and  $\gamma_h$  determines the leverage effect. The long-term variance itself follows a GJR specification with parameters  $q_q$ ,  $\beta_q$ ,  $\alpha_q$  and  $\gamma_q$ . Jumps in the asset price process are driven by iid Bernoulli variables  $I_t$  with probability  $P(I_t = 1) = \lambda$ and  $\xi_t$  is normally distributed with mean  $\mu_j$  and standard deviation  $\sigma_j$ . The variance of the jump component is given by  $\psi = \operatorname{Var}(I_t\xi_t) = \lambda \left(\mu_j^2 + \sigma_j^2\right) - \lambda^2 \mu_j^2$ . The error term  $z_t$  is iid with zero mean and unit variance, driven by either a normal distribution or a standardized Student-t distribution with degree of freedom parameter  $\eta$ .

The choice of this general discrete-time model is driven by several considerations. First,

<sup>&</sup>lt;sup>15</sup>See, e.g., Bauwens *et al.* (2006) and Engle and Ng (1993).

the model in its unrestricted form includes all the features studied for continuous-time models, namely a two-factor variance process, jumps and fat-tailed (non-Gaussian) error term distributions. And secondly, the model allows us to study nested, more parsimonious model specifications to test which features of discrete-time models are important in outof-sample exercises. We label the single factor models GJR-N and GJR-t, depending on the distribution of the error term. For these two specifications, we apply the restriction  $\lambda = 0$  and  $q_t = \bar{q}$  where  $\bar{q}$  is a constant. The corresponding two-factor DGARCH models are labeled MF-GJR-N and MF-GJR-t. For models with Gaussian error terms we also include models with normally distributed jumps, and we add an additional *J*-identifier for these jump specifications.

#### [Table 15 about here.]

Table 15 reports parameter estimates for the single- and two-factor DGARCH models. Following the discrete-time literature, parameters are estimated on daily percentage returns  $r_{t+\Delta}^p = 100 \times (s_{t+\Delta} - s_t)$  during the in-sample period from January 2, 1987 until December 29, 2006. For ease of comparison, we scale the log-likelihood at the optimal parameter set to be comparable with previously reported continuous-time models. For the sake of brevity, we do not discuss parameter estimates in detail; they are consistent overall with earlier results and values reported in the literature. Models with *t*-distributed error terms notably perform best in-sample, with large log-likelihood improvements over Gaussian models. Interestingly, single- and multi-factor models with normally distributed error terms and jumps are also outperformed by simpler models with fat-tailed error terms.

[Table 16 about here.]

In Tables 16 and 17 we compare the DGARCH model performance to various continuoustime benchmark models. For the predictive log-likelihood loss function in Table 16, model confidence sets at both the 10% and 25% level are not affected by the addition of DGARCH models, and the outperformance of SV-G is still significant. This highlights the superiority of continuous-time specifications over sophisticated DGARCH models. The best models from the DGARCH model class are GJR-t and MF-GJR-t models, and hence our results imply that the most important out-of-sample feature to consider is a fat-tailed error term. The Gneiting-Ranjan tests are summarized in Table 17 and also add further support to earlier findings. SV-G outperforms all other specifications and the 10% model confidence set is a singleton for three of the five weight functions (no weight, center, right tail). For modeling the left tail of the return distribution, we find that the 10% model confidence set includes all models, whereas the 25% set includes all but the MF-GJR-N specification. Driven by this finding, the results for the tail weight function provide evidence in favor of SV-G, with simple one-factor GJR models also providing adequate performance. To shed further light on the model ranking within the DGARCH class, we run a separate set of model confidence set estimations (unreported). These results confirm that GJR-t is the most successful discrete-time specification, being the only model in the 25% confidence set for the unweighted, center, right-tail and left-tail loss function.

## 7 Implications for Value at Risk

In this section, we provide out-of-sample tests using a VaR-based loss function. Our aim is to understand the role of complex models for a standard application in financial risk management. To this end, we base out-of-sample tests on the asymmetric VaR loss function of González-Rivera *et al.* (2004). This function penalizes return observations below VaR more than return observations that are above VaR. For the details see Section 3.3.

#### [Table 18 about here.]

We first present MCS estimations for a VaR loss function with a significance level  $\alpha = 1\%$ , as this is the most common level used for financial applications. We report esti-

mation results for all model classes in Table 18. The 25% model confidence set consists of five models, DGARCH models with fat-tailed error terms (GJR-t, MF-GJR-t, GJR-N-J) and two simple stochastic volatility models (SV-A, SV-G). All remaining models are contained in the 10% model confidence set. These findings can be interpreted as follows. First, complex continuous-time models do not provide any improvement over simpler DGARCH specifications as far as VaR estimations are concerned. Interestingly, simple DGARCH specifications (in particular GJR-t) outperform all jump-augmented continuous-time specifications. Secondly, the best-performing continuous-time models are SV-A and SV-G, a finding that supports earlier evidence in favor of these two specifications.

In order to test these results for robustness, we rerun the analysis for two further significance levels  $\alpha = 0.5\%$  and  $\alpha = 2\%$  (unreported).<sup>16</sup> Interestingly, the smaller the significance level, the more significant is the outperformance of the DGARCH specifications. For  $\alpha = 0.5\%$ , the 25% model confidence set consists of GJR-t, MF-GJR-t and GJR-N-J and the only additional model in the 10% model confidence set is GJR-N. Therefore for small significance levels, we find that simple DGARCH models significantly outperform continuous-time models. For higher  $\alpha$ -levels, the choice of model is less important, as for  $\alpha = 2\%$ , we find that all but two models (MF-GJR-N-J, MF-SV-A) are included in the 25% MCS. Overall, this finding suggests that while continuous-time models provide significant improvements when the loss function takes into account the whole density (such as the predictive log-likelihood or the unweighted CRPS statistic), simple DGARCH models with fat error terms are superior for applications that focus on the performance of the left tail only.

<sup>&</sup>lt;sup>16</sup>Detailed results for these tests are available upon request.

## 8 Conclusion

This paper studies the out-of-sample performance of several popular time-series models for S&P 500 index returns. We use an in-sample data set from 1987 to 2006 for model estimation and test how well alternative models fare in explaining index returns during an out-of-sample period starting in 2007. We test a plethora of models, including finiteand infinite-activity jumps, non-affine variance and multi-factor variance specifications, in discrete- and continuous-time. Model specification tests include likelihood-based statistics and weighted and unweighted continuous-ranked probability scores which are combined with the model confidence set procedure of Hansen *et al.* (2011).

We find that despite the highly turbulent out-of-sample market regime, simple stochastic volatility diffusions outperform more advanced jump specifications. The most important model feature is the non-affinity of the variance process; other model features are found to provide no further improvement during the out-of-sample period of this paper. Furthermore, we find that jump-diffusion models with a constant intensity parameter are misspecified for out-of-sample prediction. Our results in combination with findings in Santa-Clara and Yan (2010) suggest that improving the modeling of the time-variation in jump distributions and jump intensities are promising directions for future research.

## References

- Abhyankar, A., Copeland, L., and Wong, W. (1997). Uncovering nonlinear structure in real-time stock-market indexes: the S&P 500, the DAX, the Nikkei 225, and the FTSE-100. Journal of Business & Economic Statistics, 15(1), 37–41.
- Aït-Sahalia, Y. and Jacod, J. (2011). Testing whether jumps have finite or infinite activity. Annals of Statistics, 39(3), 1689–1719.

- Amisano, G. and Giacomini, R. (2007). Comparing Density Forecasts via Weighted Likelihood Ratio Tests. Journal of Business & Economic Statistics, 25(2), 177–190.
- Andersen, T., Benzoni, L., and Lund, J. (2002). An Empirical Investigation of Continuous-Time Equity Return Models. *The Journal of Finance*, 57(3), 1239–1284.
- Andersen, T. G., Davis, R. A., Krei
  ß, J. P., and Mikosch, T. (2009). Handbook of financial time series. Springer Berlin Heidelberg.
- Bakshi, G., Ju, N., and Ou-Yang, H. (2006). Estimation of Continuous-Time Models with an Application to Equity Volatility Dynamics. *Journal of Financial Economics*, 82(1), 227–249.
- Bakshi, G. S., Cao, C., and Chen, Z. (1997). Empirical Performance of Alternative Option Pricing Models. The Journal of Finance, 52(5), 2003–2049.
- Bao, Y., Lee, T., and Saltoglu, B. (2007). Comparing Density Forecast Models. Journal of Forecasting, 26, 203–225.
- Bardgett, C., Gourier, E., and Leippold, M. (2015). Inferring volatility dynamics and risk premia from the S&P 500 and VIX markets. *Working Paper*.
- Bates, D. S. (1996). Jumps and Stochastic Volatility: Exchange Rate Processes Implicit in Deutsche Mark Options. *Review of Financial Studies*, 9(1), 69–107.
- Bates, D. S. (2006). Maximum Likelihood Estimation of Latent Affine Processes. Review of Financial Studies, 19(3), 909–965.
- Bates, D. S. (2012). U.S. stock market crash risk, 1926-2010. Journal of Financial Economics, 105(2), 229–259.
- Bauwens, L., Laurent, S., and Rombouts, J. V. K. (2006). Multivariate GARCH models: a survey. Journal of Applied Econometrics, 21(1), 79–109.

- Berkowitz, J. (2001). Testing density forecasts, with applications to risk management. Journal of Business & Economic Statistics, **19**(4), 465–474.
- Carr, P., Geman, H., Madan, D. B., and Yor, M. (2002). The Fine Structure of Asset Returns: An Empirical Investigation. *The Journal of Business*, **75**(2), 305–333.
- Carr, P., Geman, H., Madan, D., and Yor, M. (2003). Stochastic volatility for Lévy processes. *Mathematical Finance*, 13(3), 345–382.
- Chernov, M., Gallant, A. R., Ghysels, E., and Tauchen, G. (2003). Alternative models for stock price dynamics. *Journal of Econometrics*, **116**(1-2), 225–257.
- Chourdakis, K. and Dotsis, G. (2011). Maximum Likelihood Estimation of Non-Affine Volatility Processes. Journal of Empirical Finance, 18(3), 533–545.
- Christoffersen, P., Jacobs, K., and Mimouni, K. (2010). Volatility Dynamics for the S&P500: Evidence from Realized Volatility, Daily Returns, and Option Prices. *Review* of Financial Studies, 23(8), 3141–3189.
- Christoffersen, P., Jacobs, K., and Ornthanalai, C. (2012). Dynamic jump intensities and risk premiums: Evidence from S&P500 returns and options. *Journal of Financial Economics*, **106**(3), 447–472.
- Diebold, F. C., Gunther, T. A., and Tay, A. S. (1998). Evaluating Density Forecasts with Applications to Financial Risk Management. *International Economic Review*, **39**(4), 863–883.
- Duffie, D., Pan, J., and Singleton, K. (2000). Transform Analysis and Asset Pricing for Affine Jump-Diffusions. *Econometrica*, 68(6), 1343–1376.
- Egloff, D., Leippold, M., and Wu, L. (2010). The Term Structure of Variance Swap Rates and Optimal Variance Swap Investments. *Journal of Financial and Quantitative Analysis*, 45(05), 1279–1310.

- Engle, R. F. and Ng, V. K. (1993). Measuring and testing the impact of news on volatility. Journal of Finance, 48(5), 1749–1778.
- Eraker, B. (2004). Do stock prices and volatility jump? Reconciling evidence from spot and option prices. *The Journal of Finance*, **59**(3), 1–37.
- Eraker, B., Johannes, M., and Polson, N. (2003). The Impact of Jumps in Volatility and Returns. The Journal of Finance, 58(3), 1269–1300.
- Glosten, L. R., Jagannathan, R., and Runkle, D. E. (1993). On the Relation between the Expected Value and the Volatility of the Nominal Excess Return on Stocks. *Journal of Finance*, 48(5), 1779–1801.
- Gneiting, T. and Ranjan, R. (2011). Comparing Density Forecasts Using Threshold- and Quantile-Weighted Scoring Rules. Journal of Business & Economic Statistics, 29(3), 411–422.
- González-Rivera, G., Lee, T.-H., and Mishra, S. (2004). Forecasting volatility: A reality check based on option pricing, utility function, value-at-risk, and predictive likelihood. *International Journal of Forecasting*, **20**(4), 629–645.
- Hansen, P. R., Lunde, A., and Nason, J. M. (2011). The Model Confidence Set. Econometrica, 79(2), 453–497.
- Heston, S. L. (1993). A Closed-Form Solution for Options with Stochastic Volatility with Applications to Bond and Currency Options. *Review of Financial Studies*, 6(2), 327–343.
- Ignatieva, K., Rodrigues, P., and Seeger, N. (2015). Empirical Analysis of Affine Versus Nonaffine Variance Specifications in Jump-Diffusion Models for Equity Indices. *Journal* of Business & Economic Statistics, **33**(1), 68–75.

- Johannes, M. S., Polson, N. G., and Stroud, J. R. (2009). Optimal Filtering of Jump Diffusions: Extracting Latent States from Asset Prices. *Review of Financial Studies*, 22(7), 2759–2799.
- Jones, C. S. (2003). The dynamics of stochastic volatility: evidence from underlying and options markets. *Journal of Econometrics*, **116**(1-2), 181–224.
- Kaeck, A. (2013). Asymmetry in the jump-size distribution of the S&P 500: Evidence from equity and option markets. *Journal of Economic Dynamics & Control*, **37**(9), 1872–1888.
- Kaeck, A. and Alexander, C. (2012). Volatility dynamics for the S&P 500: Further evidence from non-affine, multi-factor jump diffusions. *Journal of Banking and Finance*, 36(11), 3110–3121.
- Kou, S., Yu, C., and Zhong, H. (2013). Jumps in Equity Index Returns Before and During the Recent Financial Crisis: A Bayesian Analysis. Working Paper.
- Kou, S. G. (2002). A Jump-Diffusion Model for Option Pricing. Management Science, 48(8), 1086–1101.
- Laio, F. and Tamea, S. (2006). Verification tools for probabilistic forecasts of continuous hydrological variables. *Hydrology and Earth System Sciences Discussions*, 3(4), 2145– 2173.
- Lee, S. S. and Hannig, J. (2010). Detecting jumps from Lévy jump diffusion processes. Journal of Financial Economics, 96(2), 271–290.
- Li, H., Wells, M. T., and Yu, C. L. (2008). A Bayesian Analysis of Return Dynamics with Lévy Jumps. *Review of Financial Studies*, **21**(5), 2345–2378.
- Madan, D. B. and Seneta, E. (1990). The Variance Gamma (V.G.) Model for Share Market Returns. The Journal of Business, 63(4), 511–524.

- Mijatovic, A. and Schneider, P. (2014). Empirical Asset Pricing with Nonlinear Risk Premia. *Journal of Financial Econometrics*, **12**(3), 479–506.
- Nelson, D. B. (1990). ARCH models as diffusion approximations. *Journal of Economet*rics, **45**(1-2), 7–38.
- Ornthanalai, C. (2014). Levy jump risk: Evidence from options and returns. Journal of Financial Economics, 112, 69–90.
- Pan, J. (2002). The jump-risk premia implicit in options: evidence from an integrated time-series study. *Journal of Financial Economics*, 63(1), 3–50.
- Santa-Clara, P. and Yan, S. (2010). Crashes, Volatility, and the Equity Premium: Lessons from S&P 500 Options. *Review of Economics and Statistics*, 92(2), 435–451.
- Shackleton, M. B., Taylor, S. J., and Yu, P. (2010). A multi-horizon comparison of density forecasts for the S&P 500 using index returns and option prices. *Journal of Banking* and Finance, **34**(11), 2678–2693.
- Spiegelhalter, D. J., Best, N. G., Carlin, B. P., and van der Linde, A. (2002). Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **64**(4), 583–639.
- Stroud, J. R. and Johannes, M. S. (2014). Bayesian Modeling and Forecasting of 24-Hour High-Frequency Volatility. Journal of the American Statistical Association, 109(508), 1368–1384.
- Szerszen, P. J. (2009). Bayesian Analysis of Stochastic Volatility Models with Lévy Jumps: Application to Risk Analysis. Working Paper.
- Wilhelmsson, A. (2013). Density Forecasting with Time-Varying Higher Moments: A model Confidence Set Approach. Journal of Forecasting, 31(September), 19–31.

Yun, J. (2014). Out-of-sample density forecasts with affine jump diffusion models. Journal of Banking & Finance, 47, 74–87.

# A Simulation Study

We test the ability of the approximate maximum likelihood method to estimate the parameters of affine and non-affine specifications. This procedure extends simulation results in Bates (2006). We focus on a sample size of 4000 daily returns and simulate processes with 100 intra-daily time steps with an Euler discretization as in Eraker *et al.* (2003). We provide results for the standard stochastic volatility specification and an extension with state-depended jump probabilities.

Tables 19 and 20 report results for a small Monte Carol study with 100 random sample paths. The results indicate that the maximum likelihood method of Bates (2006) very accurately identifies the parameters of the stochastic variance and jump specifications. The local approximation to non-affine specifications leads to a minor loss in the precision of estimated parameters but the estimation methodology is still able to identify the parameters accurately.

[Table 19 about here.]

[Table 20 about here.]

### Table 1: One-Factor Continuous-Time Models

This table provides an overview of the one-factor continuous-time models used in this paper. Panel A lists all jump-diffusion models, whereas Panel B provides specifications built from the CGMY process of Carr *et al.* (2002). Column 1 provides the model number, column 2 the acronym used throughout the paper and column 3 provides a short description of the main model features.

**Panel A: One-factor jump diffusion models.** Models 1 to 12 are nested in the following SDEs (where  $\lambda_t$  is the intensity of N):

$$ds_t = \left(\mu_c - \frac{1}{2}v_t - \lambda_t \bar{k}\right) dt + \rho_v \sqrt{v_t} dW_t^v + \sqrt{1 - \rho_v^2} \sqrt{v_t} dW_t^s + \xi_t dN_t$$
  
$$dv_t = \kappa_v \left(\theta_v - v_t\right) dt + \sigma_v v_t^\gamma dW_t^v.$$

1	SV-A	Stochastic volatility model of Heston (1993), $\lambda_t = 0$ for all $t, \gamma = \frac{1}{2}$
2	SV-G	Continuous-time GARCH model with $\lambda_t = 0$ for all $t, \gamma = 1$
3	SV-C	CEV stochastic volatility model with with $\lambda_t = 0$ for all $t, \gamma \in [0.5, 1.5]$
4	SVJ-A	As model 1 with jump intensity $\lambda_t = \lambda_c$ , normally distributed jump size $\xi_t$
5	SVJ-G	As model 2 with $\lambda_t = \lambda_c$ , normally distributed jump size $\xi_t$
6	SVJ-C	As model 3 with $\lambda_t = \lambda_c$ , normally distributed jump size $\xi_t$
7	SVSJ-A	As model 1 with $\lambda_t = \lambda_v v_t$ , normally distributed jump size $\xi_t$
8	SVSJ-G	As model 2 with $\lambda_t = \lambda_v v_t$ , normally distributed jump size $\xi_t$
9	SVSJ-C	As model 3 with $\lambda_t = \lambda_v v_t$ , normally distributed jump size $\xi_t$
10	SVSJJ-A	As model 1 with $\lambda_t = \lambda_c + \lambda_v v_t$ normally distributed jump size $\xi_t$
11	SVSJJ-G	As model 2 with $\lambda_t = \lambda_c + \lambda_v v_t$ , normally distributed jump size $\xi_t$
12	SVSJJ-C	As model 3 with $\lambda_t = \lambda_c + \lambda_v v_t$ , normally distributed jump size $\xi_t$

**Panel B: One-factor Levy-jump models.** Models 13 to 24 are described by the following SDEs:

$$ds_t = \left(\mu_c - \frac{1}{2}v_t\right)dt + \rho_v \sqrt{v_t} dW_t^v + \sqrt{1 - \rho_v^2} \sqrt{v_t} dL_t$$
$$dv_t = \kappa_v \left(\theta_v - v_t\right)dt + \sigma_v v_t^{\gamma} dW_t^v.$$

13	SVYY-A	$L_t$ driven by CGMY process of Carr <i>et al.</i> (2003), $\gamma = \frac{1}{2}$
14	SVYY-G	$L_t$ driven by CGMY process of Carr <i>et al.</i> (2003), $\gamma = \overline{1}$
15	SVYY-C	$L_t$ driven by CGMY process of Carr <i>et al.</i> (2003), $\gamma \in [0.5, 1.5]$
16	SVDEXP-A	$L_t$ driven by double exponential jumps as in Kou (2002), $\gamma = \frac{1}{2}$
17	SVDEXP-G	$L_t$ driven by double exponential jumps as in Kou (2002), $\gamma = \overline{1}$
18	SVDEXP-C	$L_t$ driven by double exponential jumps as in Kou (2002), $\gamma \in [0.5, 1.5]$
19	SVVG-A	$L_t$ driven by VG process of Madan and Seneta (1990), $\gamma = \frac{1}{2}$
20	SVVG-G	$L_t$ driven by VG process of Madan and Seneta (1990), $\gamma = \overline{1}$
21	SVVG-C	$L_t$ driven by VG process of Madan and Seneta (1990), $\gamma \in [0.5, 1.5]$
22	SVYYD-A	As model 13 with additional diffusive component in $L_t$
23	SVYYD-G	As model 14 with additional diffusive component in $L_t$
24	SVYYD-C	As model 15 with additional diffusive component in $L_t$

### Table 2: Data Statistics

This table provides summary statistics for daily log returns of the S&P 500 index for the whole sample period from January 2, 1987 to December 31, 2014 as well as various sub-samples. In particular the sample in column 4 (with sample start 1987<sup>\*</sup>) excludes the observation on October 16, 1987, a market crash with a log return of -22.9%.

Sample start	1987	1987	$1987^{*}$	2007	2006	2010
Sample end	2014	2006	2006	2014	2009	2014
Observations	7057	5043	5042	2014	757	1257
Mean	0.0003	0.0003	0.0004	0.0002	-0.0003	0.0005
Standard deviation	0.0118	0.0108	0.0103	0.0140	0.0189	0.0101
Skewness	-1.2943	-2.0918	-0.2124	-0.3180	-0.1737	-0.4771
Kurtosis	31.0345	48.3124	8.9613	12.3626	9.0731	7.6493
Percentile $0.5\%$	-0.0426	-0.0330	-0.0321	-0.0542	-0.0764	-0.0355
Percentile 1%	-0.0315	-0.0273	-0.0272	-0.0453	-0.0588	-0.0292
Percentile $2\%$	-0.0251	-0.0226	-0.0226	-0.0323	-0.0479	-0.0238
Percentile $5\%$	-0.0175	-0.0161	-0.0161	-0.0223	-0.0300	-0.0163
Percentile $50\%$	0.0006	0.0005	0.0005	0.0008	0.0009	0.0007
Percentile $95\%$	0.0167	0.0158	0.0158	0.0193	0.0263	0.0153
Percentile $98\%$	0.0240	0.0223	0.0223	0.0292	0.0400	0.0217
Percentile $99\%$	0.0317	0.0276	0.0276	0.0396	0.0526	0.0288
Percentile $99.5\%$	0.0392	0.0347	0.0347	0.0464	0.0657	0.0336

jump diffusions)
(One-factor
results
estimation
parameter
In-sample
Table 3:

notation for percentages. The estimation period is from January 2, 1987 to December 29, 2006. The estimation is performed using an extension of the This table reports the parameter estimation results for the one-factor jump-diffusion models. Parameter estimates correspond to annual units and decimal

maxim parentl	um likeliho lesis. Log-li	od method pr kelihood valu	roposed in I les for each 1	3ates (2006). nodel are giv	. For each <sub>I</sub> ven in the la	barameter, v st row. For	we report th exact model	e maximum definitions s	likelihood es ee Section (2	stimates and 2).	l the standa	rd errors in
Paran	leters	SV models			SVJ models		51	SVSJ models		S	VSJJ models	
3	0.500	1.000	1.000 (0.047)	0.500	1.000	1.115 (0.082)	0.500	1.000	0.995 $(0.087)$	0.500	1.000	1.095 (0.093)
$\pi^{\circ}$	0.041	0.055	0.057	0.048	0.063	0.065	0.049	0.062	0.063	0.051	0.064	0.068
ے ل	(0.006)	(0.016)	(0.011)	(0.031)	(0.023)	(0.024)	(0.022)	(0.023)	(0.023)	(0.029)	(0.026)	(0.023)
$\kappa_v$	5.449	2.549	2.610	3.343	1.949	1.554	3.632	2.020	2.102	3.586	2.143	1.722
	(0.722)	(0.685)	(0.810)	(0.638)	(0.642)	(0.725)	(0.664)	(0.655)	(0.761)	(0.667)	(0.674)	(0.771)
$\theta_v$	0.029	0.036	0.035	0.026	0.031	0.034	0.025	0.030	0.029	0.025	0.029	0.031
	(0.002)	(0.007)	(0.006)	(0.003)	(0.008)	(0.011)	(0.003)	(0.007)	(0.007)	(0.003)	(0.006)	(0.00)
$\sigma_v$	0.435	2.963	2.948	0.309	2.355	3.624	0.309	2.274	2.242	0.307	2.292	3.266
	(0.016)	(0.101)	(0.420)	(0.021)	(0.147)	(1.075)	(0.022)	(0.151)	(0.734)	(0.022)	(0.153)	(1.126)
$ ho_v$	-0.681	-0.782	-0.774	-0.697	-0.773	-0.783	-0.700	-0.774	-0.772	-0.701	-0.774	-0.784
	(0.034)	(0.030)	(0.030)	(0.040)	(0.036)	(0.036)	(0.039)	(0.036)	(0.036)	(0.040)	(0.036)	(0.036)
$\lambda_c$				0.901	0.750	0.742				0.001	0.003	0.004
				(0.290)	(0.243)	(0.250)				(0.399)	(0.358)	(0.377)
$\lambda_v$							54.909	49.414	49.155	55.134	48.720	48.572
							(16.162)	(14.888)	(15.225)	(26.061)	(24.218)	(27.657)
$\mu_s$				-0.028	-0.027	-0.029	-0.027	-0.026	-0.026	-0.027	-0.026	-0.025
				(0.018)	(0.00)	(0.010)	(0.008)	(0.008)	(0.00)	(0.014)	(0.013)	(0.014)
$\sigma_s$				0.059	0.056	0.053	0.054	0.050	0.050	0.054	0.050	0.048
				(0.008)	(0.007)	(0.006)	(0.006)	(0.006)	(0.006)	(0.007)	(0.007)	(0.007)
ΓΓ	16722	16755	16755	16762	16783	16784	16769	16788	16788	16769	16788	16788

models
jump
Levy
One-factor
results (
estimation
parameter
In-sample
÷
Table '

This table reports the parameter estimation results for the one-factor Lévy jump models. Parameter estimates correspond to annual units and decimal notation for percentages. The estimation period is from January 2, 1987 to December 29, 2006. The estimation is performed using an extension of the maximum likelihood method proposed in Bates (2006). For each parameter, we report the maximum likelihood estimates and the standard errors in

Param	eters	SVYY model	ν.	SV	'DEXP mode	els	$\mathbf{S}$	VVG model:	S	IS	/YYD mode	s
3	0.500	1.000	1.043 (0.068)	0.500	1.000	1.018 (0.077)	0.500	1.000	1.197 (0.078)	0.500	1.000	1.001 (0.093)
$\mu_c$	0.061	0.071	0.077	0.053	0.064	0.067	0.050	0.062	0.074	0.059	0.070	0.072
	(0.029)	(0.025)	(0.027)	(0.028)	(0.021)	(0.021)	(0.020)	(0.019)	(0.024)	(0.027)	(0.026)	(0.027)
$\kappa_v$	3.775	2.110	1.840	3.531	2.049	2.021	3.507	2.040	1.395	3.705	2.129	2.162
	(0.674)	(0.655)	(0.677)	(0.651)	(0.563)	(0.763)	(0.628)	(0.572)	(0.702)	(0.658)	(0.673)	(0.801)
$\theta_v$	0.029	0.032	0.031	0.030	0.034	0.034	0.031	0.035	0.040	0.029	0.032	0.032
	(0.004)	(0.007)	(0.007)	(0.004)	(0.007)	(0.008)	(0.004)	(0.007)	(0.014)	(0.003)	(0.007)	(0.007)
$\sigma_v$	0.323	2.159	2.473	0.331	2.261	2.399	0.335	2.238	4.482	0.320	2.178	2.173
	(0.032)	(0.151)	(0.616)	(0.025)	(0.146)	(0.639)	(0.027)	(0.146)	(1.195)	(0.024)	(0.154)	(0.726)
$ ho_v$	-0.632	-0.711	-0.729	-0.640	-0.717	-0.717	-0.631	-0.711	-0.723	-0.632	-0.710	-0.713
	(0.054)	(0.048)	(0.042)	(0.041)	(0.037)	(0.039)	(0.044)	(0.041)	(0.046)	(0.042)	(0.048)	(0.046)
$f_{jump}$				0.289	0.301	0.297	0.338	0.339	0.365	0.998	1.000	0.997
• 2				(0.077)	(0.081)	(0.081)	(0.104)	(0.102)	(0.127)	(0.000)	(0.000)	(1.977)
$w_n$	0.330	0.348	0.351	0.832	0.848	0.848	0.832	0.839	0.855	0.319	0.348	0.345
	(0.126)	(0.114)	(0.118)	(0.106)	(0.106)	(0.107)	(0.110)	(0.109)	(0.110)	(0.074)	(0.126)	(0.637)
G	3.761	3.944	6.479	27.500	28.000	28.000	14.188	15.306	13.363	4.996	3.838	3.990
	(3.018)	(2.833)	(6.914)	(6.045)	(6.190)	(5.875)	(5.759)	(5.816)	(5.809)	(0.002)	(3.399)	(3.952)
$Y_n$	1.297	1.368	1.637							1.122	1.379	1.399
	(0.450)	(0.460)	(0.276)							(0.055)	(0.473)	(0.514)
$Y_p$	1.858	1.837	1.809							1.855	1.812	1.804
	(0.058)	(0.075)	(0.088)							(0.055)	(0.086)	(0.615)
ΓΓ	16779	16797	16797	16771	16790	16790	16772	16791	16791	16780	16797	16797

# Table 5: Model Confidence Set p-Values and Model Ranking Full Out-of-sample Period Using Predictive Likelihood

This table shows model confidence set results for the full out-of-sample period January 3, 2007 to December 31, 2014 using predictive likelihood as the ranking criteria. For details regarding notation, see Section 3.2 and Section 3.3. The first column indicates the number of the iterative elimination step for models running from i = 1 to total number of models  $(m_0 = 14)$ . The second column shows the *p*-values for the hypotheses  $H_{0,\mathcal{M}_i}$  and the third column presents the MCS *p*-value  $\hat{p}_{e_{\mathcal{M}_i}}$  for the model that is removed in the respective elimination step. The fourth column shows the model eliminated in each iterative step by the elimination rule and thereby presents the model ranking according the MCS criteria, with the worst model ranked at the top and the best model at the bottom of the table, respectively. For a given significance level  $\alpha$  any model for which holds  $\hat{p}_{e_{\mathcal{M}_i}} \geq \alpha$  is included in the MCS  $\hat{\mathcal{M}}^*_{1-\alpha}$ .

Elimination Rule	<i>p</i> -Value for $H_{0,\mathcal{M}_i}$	MCS $p\text{-Value}\ \hat{p}_{e_{\mathcal{M}_i}}$	Eliminated Model
$e_{\mathcal{M}_1}$	0.0111	0.0111	SVYYD-A
$e_{\mathcal{M}_2}$	0.0104	0.0111	SVYY-A
$e_{\mathcal{M}_3}$	0.0092	0.0111	SVDEXP-A
$e_{\mathcal{M}_4}$	0.0085	0.0111	SVSJ-A
$e_{\mathcal{M}_5}$	0.0083	0.0111	SVVG-A
$e_{\mathcal{M}_6}$	0.0081	0.0111	SVJ-A
$e_{\mathcal{M}_7}$	0.0070	0.0111	SVSJ-G
$e_{\mathcal{M}_8}$	0.0084	0.0111	SVYY-G
$e_{\mathcal{M}_9}$	0.0090	0.0111	SVYYD-G
$e_{\mathcal{M}_{10}}$	0.0100	0.0111	SVDEXP-G
$e_{\mathcal{M}_{11}}$	0.0200	0.0200	SVVG-G
$e_{\mathcal{M}_{12}}$	0.0681	0.0681	SVJ-G
$e_{\mathcal{M}_{13}}$	0.2308	0.2308	SV-A
$e_{\mathcal{M}_{14}}$	1.0000	1.0000	SV-G

# Table 6: Model Confidence Set p-Values and Model RankingFirst Part and Second Part of Out-of-sample Using Predictive Likelihood

This table shows model confidence set results for the first part of the out-of-sample period January 3, 2007 to December 31, 2009 in the upper panel and the second part of the out-of-sample period January 4, 2010 to December 31, 2014 in the lower panel using predictive likelihood as the ranking criteria. For details regarding notation, see Section 3.2 and Section 3.3. The first column indicates the number of the iterative elimination step for models running from i = 1 to total number of models ( $m_0 = 14$ ). The second column shows the *p*-values for the hypotheses  $H_{0,\mathcal{M}_i}$  and the third column presents the MCS *p*-value  $\hat{p}_{e_{\mathcal{M}_i}}$  for the model that is removed in the respective elimination step. The fourth column shows the model ranking according the MCS criteria, with the worst model ranked at the top and the best model at the bottom of the table, respectively. For a given significance level  $\alpha$  any model for which holds  $\hat{p}_{e_{\mathcal{M}_i}} \geq \alpha$  is included in the MCS  $\widehat{\mathcal{M}}^*_{1-\alpha}$ .

Pa	nel A:	January	2007	to	December	2009	

Elimination Rule	<i>p</i> -Value for $H_{0,\mathcal{M}_i}$	MCS $p\text{-Value}\ \hat{p}_{e_{\mathcal{M}_i}}$	Eliminated Model
$e_{\mathcal{M}_1}$	0.1504	0.1504	SVYYD-A
$e_{\mathcal{M}_2}$	0.1499	0.1504	SVYY-A
$e_{\mathcal{M}_3}$	0.1504	0.1504	SVSJ-A
$e_{\mathcal{M}_4}$	0.1497	0.1504	SVDEXP-A
$e_{\mathcal{M}_5}$	0.1499	0.1504	SVVG-A
$e_{\mathcal{M}_6}$	0.1535	0.1535	SVJ-A
$e_{\mathcal{M}_7}$	0.1937	0.1937	SV-A
$e_{\mathcal{M}_8}$	0.5746	0.5746	SVSJ-G
$e_{\mathcal{M}_9}$	0.5490	0.5746	SVYY-G
$e_{\mathcal{M}_{10}}$	0.7044	0.7044	SVYYD-G
$e_{\mathcal{M}_{11}}$	0.7812	0.7812	SVDEXP-G
$e_{\mathcal{M}_{12}}$	0.7846	0.7846	SVVG-G
$e_{\mathcal{M}_{13}}$	0.9958	0.9958	SVJ-G
$e_{\mathcal{M}_{14}}$	1.0000	1.0000	SV-G

#### Panel B: January 2010 to December 2014

Elimination Rule	<i>p</i> -Value for $H_{0,\mathcal{M}_k}$	MCS $p$ -Value	Eliminated Model
$e_{\mathcal{M}_1}$	0.0033	0.0033	SVYYD-A
$e_{\mathcal{M}_2}$	0.0026	0.0033	SVYY-G
$e_{\mathcal{M}_3}$	0.0034	0.0034	SVDEXP-A
$e_{\mathcal{M}_4}$	0.0028	0.0034	SVVG-G
$e_{\mathcal{M}_5}$	0.0041	0.0041	SVVG-A
$e_{\mathcal{M}_6}$	0.0033	0.0041	SVSJ-G
$e_{\mathcal{M}_7}$	0.0061	0.0061	SVYYD-G
$e_{\mathcal{M}_8}$	0.0066	0.0066	SVSJ-A
$e_{\mathcal{M}_{9}}$	0.0082	0.0082	SVDEXP-G
$e_{\mathcal{M}_{10}}$	0.0095	0.0095	SVJ-A
$e_{\mathcal{M}_{11}}$	0.0211	0.0211	SVYY-A
$e_{\mathcal{M}_{12}}$	0.0029	0.0211	SVJ-G
$e_{\mathcal{M}_{13}}$	0.9490	0.9490	SV-G
$e_{\mathcal{M}_{14}}$	1.0000	1.0000	SV-A

# Table 7: Gneiting-Ranjan TestsFull Out-of-sample Dataset (No weighting).

This table reports the Gneiting and Ranjan (2011) test statistics  $t_n = \sqrt{n} \left(\overline{CRPS}_w^f - \overline{CRPS}_w^g\right) \hat{\sigma}_n^{-1}$  for several model pairs with CRPS denoting continuous ranked probability score and f and g denoting forecasting densities of the models to be tested against each other. The test statistic follows asymptotically a standard normal distribution. For details of calculation, see Section 3.3. The models in rows refer to forecasting density f and the models in columns to forecasting density g, respectively. A positive statistic therefore indicates that the model in the row is out-performed by the model in the column and vice versa.

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
SV-A (1)	_	2.39	-3.34	1.84	-3.76	1.31	-3.28	1.51	-3.74	1.34
SV-G (2)	-2.39	_	-3.09	-2.50	-3.25	-3.10	-3.14	-2.80	-3.26	-3.10
SVJ-A (3)	3.34	3.09	_	2.89	-1.45	2.50	-0.27	2.63	-1.55	2.51
SVJ-G(4)	-1.84	2.50	-2.89	_	-3.04	-2.99	-2.92	-2.03	-3.06	-2.92
SVYY-A (5)	3.76	3.25	1.45	3.04	_	2.73	2.50	2.83	-0.86	2.74
SVYY-G (6)	-1.31	3.10	-2.50	2.99	-2.73	_	-2.58	2.80	-2.75	2.14
SVVG-A (7)	3.28	3.14	0.27	2.92	-2.50	2.58	_	2.70	-3.04	2.59
SVVG-G (8)	-1.51	2.80	-2.63	2.03	-2.83	-2.80	-2.70	_	-2.86	-2.54
SVYYD-A $(9)$	3.74	3.26	1.55	3.06	0.86	2.75	3.04	2.86	_	2.76
SVYYD-G $(10)$	-1.34	3.10	-2.51	2.92	-2.74	-2.14	-2.59	2.54	-2.76	—

# Table 8: Gneiting-Ranjan TestsFull Out-of-sample Dataset Dataset (left tail weighting).

This table reports the Gneiting and Ranjan (2011) test statistics  $t_n = \sqrt{n} \left(\overline{CRPS}_w^f - \overline{CRPS}_w^g\right) \hat{\sigma}_n^{-1}$  for several model pairs with CRPS denoting continuous ranked probability score with weight function  $w(\alpha) = (1-\alpha)^2$  and f and g denoting forecasting densities of the models to be tested against each other. The test statistic follows asymptotically a standard normal distribution. For details of calculation see Section 3.3. The models in rows refer to forecasting density f and the models in columns to forecasting density g, respectively. A positive statistic therefore indicates that the model in the row is out-performed by the model in the column and vice versa.

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
SV-A (1)	_	1.11	-1.83	1.01	-2.03	0.78	-1.77	0.90	-2.03	0.80
SV-G(2)	-1.11	_	-1.52	-0.64	-1.59	-1.07	-1.50	-0.84	-1.59	-1.04
SVJ-A (3)	1.83	1.52	_	1.56	-0.77	1.43	0.13	1.51	-0.81	1.44
SVJ-G(4)	-1.01	0.64	-1.56	_	-1.63	-1.47	-1.53	-0.96	-1.63	-1.43
SVYY-A (5)	2.03	1.59	0.77	1.63	_	1.53	1.74	1.60	-0.26	1.54
SVYY-G (6)	-0.78	1.07	-1.43	1.47	-1.53	_	-1.41	1.68	-1.53	1.25
SVVG-A (7)	1.77	1.50	-0.13	1.53	-1.74	1.41	_	1.49	-1.92	1.42
SVVG-G (8)	-0.90	0.84	-1.51	0.96	-1.60	-1.68	-1.49	_	-1.60	-1.52
SVYYD-A $(9)$	2.03	1.59	0.81	1.63	0.26	1.53	1.92	1.60	_	1.54
SVYYD-G $(10)$	-0.80	1.04	-1.44	1.43	-1.54	-1.25	-1.42	1.52	-1.54	_

# Table 9: Model Confidence Set p-Values and Model Ranking Full Out-of-sample Period Using CRPS

This table shows model confidence set results for the full out-of-sample period from January 3, 2007 to December 31, 2014 using continuous ranked probability score (CRPS) as the ranking criteria. For details of notation and calculation see Sections 3.2 and 3.3. The first column provides model specifications, the second column provides results for the non-weighted CRPS statistic. Columns 3 to 6 refer to the results for the weighted CRPS statistics. The weighting scheme "Center" applies more weight to the center of the predictive density when calculating CRPS and the weighting schemes "Tails", "Right Tail", and "Left Tail" work accordingly. For a given significance level  $\alpha$  models for which  $\hat{p}_{e_{\mathcal{M}_i}} \geq \alpha$  are included in the MCS  $\widehat{\mathcal{M}}^*_{1-\alpha}$ . We use \* (\*\*) to indicate that the model belongs to the 10% (25%) MCS.

Model Name	No Weight	Center	Tails	Right Tail	Left Tail
SV-A	0.0444	0.0286	0.0928	0.0156	$0.5016^{**}$
SV-G	$1.0000^{**}$	$1.0000^{**}$	$1.0000^{**}$	$1.0000^{**}$	$1.0000^{**}$
SVJ-A	0.0427	0.0286	0.0928	0.0156	$0.3975^{**}$
SVJ-G	0.0500	0.0393	0.0928	0.0156	$0.5527^{**}$
SVSJ-A	0.0427	0.0286	0.0928	0.0156	$0.3657^{**}$
SVSJ-G	0.0500	0.0393	0.0928	0.0156	$0.5016^{**}$
SVYY-A	0.0427	0.0286	0.0928	0.0156	$0.3368^{**}$
SVYY-G	0.0444	0.0291	0.0928	0.0156	$0.5016^{**}$
SVDEXP-A	0.0427	0.0286	0.0928	0.0156	$0.3657^{**}$
SVDEXP-G	0.0500	0.0393	0.0928	0.0156	$0.5527^{**}$
SVVG-A	0.0427	0.0286	0.0928	0.0156	$0.3691^{**}$
SVVG-G	0.0500	0.0393	0.0928	0.0156	$0.5527^{**}$
SVYYD-A	0.0427	0.0286	0.0928	0.0156	$0.3359^{**}$
SVYYD-G	0.0444	0.0291	0.0928	0.0156	$0.5016^{**}$

# Table 10: Model Confidence Set p-Values and Model RankingFirst Part and Second Part of Out-of-sample Using CRPS

This table provides model confidence set results for the first part of out-of-sample period from January 3, 2007 to December 31, 2009 in panel A and the second part of the out-of-sample period January 4, 2010 to December 31, 2014 in Panel B. The loss function is given by the continuous ranked probability score (CRPS). For details of notation and calculation see Sections 3.2 and 3.3. The first column provides model specifications, the second column provides results for the non-weighted CRPS statistic. Columns 3 to 6 refer to the results for the weighted CRPS statistics. The weighting scheme "Center" applies more weight to the center of the predictive density when calculating CRPS and the weighting schemes "Tails", "Right Tail", and "Left Tail" work accordingly. For a given significance level  $\alpha$  models for which  $\hat{p}_{e_{\mathcal{M}_i}} \geq \alpha$  are included in the MCS  $\widehat{\mathcal{M}}_{1-\alpha}^*$ . We use \* (\*\*) to indicate that the model belongs to the 10% (25%) MCS.

Model Name	No Weight	Center	Tails	Right Tail	Left Tail
SV-A	$0.2565^{**}$	$0.2746^{**}$	$0.1739^{*}$	0.0868	$0.3329^{**}$
SV-G	$1.0000^{**}$	$1.0000^{**}$	$1.0000^{**}$	$1.0000^{**}$	$0.9565^{**}$
SVJ-A	$0.1392^{*}$	$0.1437^{*}$	$0.1489^{*}$	0.0629	$0.3329^{**}$
SVJ-G	$0.4570^{**}$	$0.4730^{**}$	$0.2869^{**}$	$0.1175^{*}$	$1.0000^{**}$
SVSJ-A	$0.1676^{*}$	$0.1707^{*}$	$0.1489^{*}$	0.0629	$0.3329^{**}$
SVSJ-G	$0.4399^{**}$	$0.4730^{**}$	$0.2869^{**}$	$0.1175^{*}$	$0.8719^{**}$
SVYY-A	$0.1284^{*}$	$0.1437^{*}$	$0.1370^{*}$	0.0629	$0.3292^{**}$
SVYY-G	$0.2565^{**}$	$0.4730^{**}$	$0.2285^{*}$	0.0974	$0.7440^{**}$
SVDEXP-A	$0.1473^{*}$	$0.1464^{*}$	$0.1489^{*}$	0.0629	$0.3329^{**}$
SVDEXP-G	$0.4570^{**}$	$0.4730^{**}$	$0.2869^{**}$	$0.1175^{*}$	$0.9565^{**}$
SVVG-A	$0.1529^{*}$	$0.1437^{*}$	$0.1489^{*}$	0.0629	$0.3394^{**}$
SVVG-G	$0.4570^{**}$	$0.4730^{**}$	$0.2869^{**}$	0.0974	$0.9003^{**}$
SVYYD-A	$0.1301^{*}$	$0.1437^{*}$	$0.1370^{*}$	0.0629	$0.3329^{**}$
SVYYD-G	$0.2565^{**}$	$0.4730^{**}$	$0.2285^{*}$	0.0868	$0.8047^{**}$

Panel A: January 2007 to December 2009

Panel B: January 2010 to December 2014

Model Name	No Weight	Center	Tails	Right Tail	Left Tail
SV-A	0.2627**	0.1042*	$0.5905^{**}$	0.0788	1.0000**
SV-G	$1.0000^{**}$	$1.0000^{**}$	$1.0000^{**}$	$1.0000^{**}$	$0.8145^{**}$
SVJ-A	0.0889	0.0227	$0.1891^{*}$	0.0706	$0.5458^{**}$
SVJ-G	$0.1077^{*}$	0.0298	$0.1891^{*}$	0.0706	$0.5458^{**}$
SVSJ-A	0.0889	0.0273	$0.1490^{*}$	0.0706	$0.5189^{**}$
SVSJ-G	0.0889	0.0298	$0.1700^{*}$	0.0706	$0.5458^{**}$
SVYY-A	0.0889	0.0227	$0.1490^{*}$	0.0706	$0.4789^{**}$
SVYY-G	0.0889	0.0298	$0.1490^{*}$	0.0706	$0.5189^{**}$
SVDEXP-A	0.0889	0.0227	$0.1273^{*}$	0.0706	$0.3945^{**}$
SVDEXP-G	0.0889	0.0298	$0.1700^{*}$	0.0706	$0.5458^{**}$
SVVG-A	0.0889	0.0225	$0.1142^{*}$	0.0706	$0.4321^{**}$
SVVG-G	0.0889	0.0298	$0.1490^{*}$	0.0706	$0.5458^{**}$
SVYYD-A	0.0889	0.0221	$0.1078^{*}$	0.0698	$0.3710^{**}$
SVYYD-G	0.0889	0.0298	$0.1490^{*}$	0.0706	$0.5458^{**}$

### Table 11: Likelihood Ratio Tests (Sub-Sample Analysis). $H_0: \mu = \rho = 0 \text{ and } \sigma = 1$

This table reports the likelihood ratios (LR) for the likelihood ratio test proposed in Berkowitz (2001) as described in Section (5.4). Based on the estimated optimal parameters sets for time period 1987 to end of 2006, LRs are calculated for the full out-of-sample period from start of 2007 until end of 2014 and for each of the out-of-sample years separately. The LR is calculated as  $LR = -2[\mathcal{L}(0, 1, 0) - \mathcal{L}(\hat{\mu}, \hat{\sigma}^2, \hat{\rho})]$ , which serve as test statistics for the null hypothesis H0:  $(\mu = 0, \sigma^2 = 1, \rho = 0)$ , jointly testing the probability integral transforms for independence and mean and variance equal to (0, 1). The test statistic is distributed  $\chi^2(3)$  with critical values given by: 99% level –  $\chi^2(3) = 11.34$  (\*\*\*), 95% level:  $\chi^2(3) = 9.210$  (\*\*), and 90% level:  $\chi^2(3) = 6.635$  (\*).

Model	All	2007	2008	2009	2010	2011	2012	2013	2014
SV-A	25.18***	5.36	34.61***	5.11	3.12	7.39*	0.79	6.70*	1.01
SV-G	24.83***	$7.40^{*}$	$13.65^{***}$	3.01	5.93	$7.29^{*}$	0.27	$6.85^{*}$	0.77
SVJ-A	24.24***	6.45	33.17***	2.38	2.48	6.60	0.91	6.09	0.45
SVJ-G	$23.35^{***}$	8.08*	$13.89^{***}$	2.05	4.33	6.03	0.14	6.25	1.24
SVSJ-A	22.37***	6.39	26.01***	2.22	2.38	$6.89^{*}$	0.75	6.13	0.42
SVSJ-G	22.31***	8.13*	$12.65^{***}$	1.65	4.07	6.24	0.16	6.47	1.50
SVYY-A	27.02***	$7.09^{*}$	$30.83^{***}$	2.66	2.80	8.99*	0.78	5.94	0.45
SVYY-G	$25.60^{***}$	9.00*	13.73***	1.76	4.56	$7.84^{*}$	0.15	6.37	1.75
SVDEXP-A	22.73***	6.51	$26.43^{***}$	2.17	2.36	$7.27^{*}$	0.80	6.09	0.43
SVDEXP-G	22.87***	$8.26^{*}$	12.80***	1.68	4.11	6.51	0.16	6.51	1.55
SVVG-A	22.67***	6.56	$26.09^{***}$	2.14	2.39	$7.29^{*}$	0.83	6.19	0.44
SVVG-G	22.77***	8.29*	12.86***	1.68	4.09	6.57	0.16	6.55	1.51
SVYYD-A	$26.80^{***}$	$7.05^{*}$	30.76***	2.54	2.75	$8.97^{*}$	0.77	5.93	0.44
SVYYD-G	25.55***	8.96*	13.69***	1.78	4.60	7.78*	0.15	6.38	1.74

	,
liffusions).	
ıp c	
jum	
factor .	į
(Multi-	
results	. 8
ion	
estimat	•
ter	
ame	¢
para	,
mple	•
[n-sa]	
12: ]	
Table	

This table reports the parameter estimation results for the two-factor jump-diffusion models. The estimation period is from 2 January 1987 to 29 December 2006. The estimation is performed using an extension of the maximum likelihood method proposed in Bates (2006). For each parameter, we report the maximum likelihood estimates and the standard errors in parenthesis. Log-likelihood values for each model are given in the last row.

Paran	neters	SV models			SVJ models			SVSJ models		S	VSJJ model	
Х	0.500	1.000	1.175 (0.035)	0.500	1.000	$1.176 \\ (0.041)$	0.500	1.000	1.057 (0.040)	0.500	1.000	1.127 (0.075)
$\mu_c$	0.045	0.038	0.060	0.038	0.047	0.059	0.046	0.048	0.051	0.046	0.048	0.053
	(0.002)	(0.027)	(0.026)	(0.021)	(0.028)	(0.026)	(0.024)	(0.022)	(0.024)	(0.027)	(0.027)	(0.021)
$\kappa_v$	26.928	16.078	11.402	8.017	9.052	8.400	8.544	9.366	9.287	8.451	9.467	9.285
	(2.793)	(1.820)	(1.521)	(1.419)	(1.486)	(1.588)	(1.486)	(1.466)	(1.539)	(1.391)	(1.508)	(1.547)
$\sigma_v$	0.633	4.186	6.670	0.349	2.856	5.501	0.348	2.802	3.488	0.345	2.829	4.547
	(0.028)	(0.152)	(0.724)	(0.027)	(0.198)	(0.715)	(0.027)	(0.202)	(0.586)	(0.027)	(0.205)	(1.307)
$\kappa_m$	1.408	1.022	0.946	0.574	0.534	0.516	0.602	0.593	0.586	0.628	0.601	0.570
	(0.171)	(0.117)	(0.125)	(0.221)	(0.182)	(0.197)	(0.210)	(0.161)	(0.188)	(0.170)	(0.205)	(0.202)
$\theta_m$	0.018	0.016	0.014	0.019	0.017	0.016	0.018	0.016	0.016	0.018	0.016	0.016
	(0.002)	(0.002)	(0.002)	(0.004)	(0.003)	(0.003)	(0.003)	(0.003)	(0.003)	(0.003)	(0.003)	(0.003)
$\sigma_m$	0.228	1.727	3.765	0.194	1.794	4.156	0.190	1.821	2.358	0.196	1.821	3.210
	(0.028)	(0.247)	(0.798)	(0.052)	(0.413)	(1.431)	(0.045)	(0.348)	(0.609)	(0.034)	(0.426)	(1.350)
$ ho_v$	-0.651	-0.832	-0.857	-0.744	-0.861	-0.906	-0.747	-0.865	-0.880	-0.747	-0.864	-0.891
	(0.034)	(0.023)	(0.022)	(0.036)	(0.027)	(0.026)	(0.036)	(0.027)	(0.027)	(0.037)	(0.027)	(0.027)
$\lambda_{c}$				0.976	0.852	1.385				0.012	0.011	0.048
				(0.285)	(0.299)	(0.461)				(0.440)	(0.010)	(0.303)
$\lambda_v$							56.381	51.751	52.239	57.074	50.258	49.514
							(17.275)	(16.977)	(16.886)	(25.771)	(17.426)	(20.483)
$\mu_s$				-0.028	-0.024	-0.008	-0.027	-0.026	-0.025	-0.026	-0.027	-0.025
				(0.012)	(0.020)	(0.007)	(0.012)	(0.001)	(0.011)	(0.004)	(0.002)	(0.002)
$\sigma_s$				0.062	0.062	0.041	0.059	0.054	0.053	0.060	0.054	0.053
				(0.002)	(0.007)	(0.002)	(0.007)	(0.007)	(0.007)	(0.007)	(0.007)	(0.008)
$\Gamma\Gamma$	16722	16781	16789	16771	16812	16816	16777	16816	16818	16777	16816	16820

### Table 13: Model Confidence Set *p*-Values and Model Ranking Multi-factor Models for Full Out-of-sample Period using Predictive Likelihood

This table shows model confidence set results for the full out-of-sample period January 3, 2007 to December 31, 2014 using predictive likelihood as the ranking criteria. For details of notation and calculation see Section 3.2 and Section 3.3. Multi-factor models are tested against SV-A, SV-G, SVSJJ-A, SVSJJ-G, SVYYD-A, and SVYYD-G. The first column indicates the number of the iterative elimination step for models running from i = 1 to total number of models  $(m_0 = 17)$ . The second column shows the *p*-values for the hypotheses  $H_{0,\mathcal{M}_i}$  and the third column presents the MCS *p*-value  $\hat{p}_{e_{\mathcal{M}_i}}$  for the model that is removed in the respective elimination step. The fourth column shows the model eliminated in each iterative step by the elimination rule and thereby presents the model ranking according to the MCS criteria, with the worst model ranked at the top and the best model at the bottom of the table respectively. For a given significance level  $\alpha$  any model for which holds  $\hat{p}_{e_{\mathcal{M}_i}} \geq \alpha$  is included in the MCS  $\hat{\mathcal{M}}_{1-\alpha}^*$ .

Elimination Rule	<i>p</i> -Value for $H_{0,\mathcal{M}_k}$	MCS $p$ -Value	Eliminated Model
$e_{\mathcal{M}_1}$	0.0205	0.0205	MF-SVSJJ-A
$e_{\mathcal{M}_2}$	0.0242	0.0242	MF-SVSJ-A
$e_{\mathcal{M}_3}$	0.0302	0.0302	SVYYD-A
$e_{\mathcal{M}_4}$	0.0301	0.0302	MF-SVJ-A
$e_{\mathcal{M}_5}$	0.0410	0.0410	SVSJJ-A
$e_{\mathcal{M}_6}$	0.0439	0.0439	MF-SV-A
$e_{\mathcal{M}_7}$	0.1958	0.1958	SVYYD-G
$e_{\mathcal{M}_8}$	0.1955	0.1958	SVSJJ-G
$e_{\mathcal{M}_9}$	0.1958	0.1958	MF-SVSJJ-C
$e_{\mathcal{M}_{10}}$	0.2057	0.2057	MF-SVJ-C
$e_{\mathcal{M}_{11}}$	0.2261	0.2261	MF-SVSJ-G
$e_{\mathcal{M}_{12}}$	0.2280	0.2280	MF-SVJ-G
$e_{\mathcal{M}_{13}}$	0.2353	0.2353	MF-SV-C
$e_{\mathcal{M}_{14}}$	0.4827	0.4827	SV-A
$e_{\mathcal{M}_{15}}$	0.2439	0.4827	MF-SVSJJ-G
$e_{\mathcal{M}_{16}}$	0.9147	0.9147	SV-G
$e_{\mathcal{M}_{17}}$	1.0000	1.0000	MF-SV-G

# Table 14: Model Confidence Set p-Values and Model RankingMulti-factor Models for Full Out-of-sample Period Using CRPS

This table shows model confidence set results for the full out-of-sample period January 3, 2007 to December 31, 2014 using continuous ranked probability score (CRPS) as the ranking criteria. For details of notation and calculation see Sections 3.2 and Section 3.3. Multi-factor models are tested against SV-A, SV-G, SVSJJ-A, SVSJJ-G, SVYYD-A, and SVYYD-G. First column gives models tested listed from least complex model at the top to most complex model at the bottom of table. Column 2 gives results of the non-weighted CRPS version. Columns 3 to 6 refer to the results of the weighted CRPS versions. The weighting scheme "Center" puts more weight on the center of the predictive density when calculating CRPS and the weighting schemes "Tails", "Right Tail", and "Left Tail" work accordingly. For a given significance level  $\alpha$  any model for which holds  $\hat{p}_{e_{\mathcal{M}_i}} \geq \alpha$  is included in the MCS  $\widehat{\mathcal{M}}_{1-\alpha}^*$ . One \* indicates the model belongs to the 10% MCS and two \*\* indicate model belongs to the 25% MCS.

Model Name	No Weight	Center	Tails	Right Tail	Left Tail
MF-SV-A	0.0221	0.1337*	0.0477	0.0400	0.2438*
MF-SV-G	$0.6180^{**}$	$0.6506^{**}$	$0.6586^{**}$	$0.2857^{**}$	$0.5211^{**}$
MF-SV-C	$0.6180^{**}$	$0.6506^{**}$	$0.6586^{**}$	$1.0000^{**}$	$0.4998^{**}$
MF-SVJ-A	0.0196	0.0150	$0.1073^{*}$	0.0124	$0.4129^{**}$
MF-SVJ-G	$0.4500^{**}$	$0.4851^{**}$	$0.5362^{**}$	$0.2842^{**}$	$0.5211^{**}$
MF-SVJ-C	$0.6180^{**}$	$0.5607^{**}$	$0.6586^{**}$	$0.2857^{**}$	$0.4998^{**}$
MF-SVSJ-A	$0.1077^{*}$	$0.1160^{*}$	0.0646	0.0909	$0.3662^{**}$
MF-SVSJ-G	$0.5195^{**}$	$0.4851^{**}$	$0.5919^{**}$	$0.2842^{**}$	$0.5211^{**}$
MF-SVSJ-C	$0.6180^{**}$	$0.6506^{**}$	$0.6586^{**}$	$0.2857^{**}$	$0.5211^{**}$
MF-SVSJJ-A	0.0221	0.0186	0.0513	0.0208	$0.2872^{**}$
MF-SVSJJ-G	$0.5195^{**}$	$0.5607^{**}$	$0.5919^{**}$	$0.2857^{**}$	$0.5211^{**}$
MF-SVSJJ-C	$0.6180^{**}$	$0.6506^{**}$	$0.6586^{**}$	$0.2857^{**}$	$0.5211^{**}$
SV-A	$0.3884^{**}$	$0.4017^{**}$	$0.4911^{**}$	$0.1110^{*}$	$0.5211^{**}$
SV-G	$1.0000^{**}$	$1.0000^{**}$	$1.0000^{**}$	$0.2857^{**}$	$1.0000^{**}$
SVSJ-A	$0.1532^{*}$	0.0424	$0.3063^{**}$	0.0118	$0.5211^{**}$
SVSJ-G	$0.5195^{**}$	$0.5607^{**}$	$0.5919^{**}$	$0.1110^{*}$	$0.5211^{**}$
SVYYD-A	0.0508	0.0146	$0.2290^{*}$	0.0106	$0.4998^{**}$
SVYYD-G	$0.4500^{**}$	$0.4851^{**}$	$0.5362^{**}$	0.0909	$0.5211^{**}$

# Table 15: In-sample parameter estimation results (Discrete-time GARCH Models).

This table reports the parameter estimation results for discrete-time GARCH models. The estimation period is from January 2, 1987 to December 29, 2006. The estimation is performed using maximum likelihood method. For each parameter, we report the maximum likelihood estimates and the standard errors in parenthesis. Log-likelihood values for each model are given in the last row. For exact model definitions see Section (6.2).

	GJR-N	MF-GJR-N	GJR-N-J	MF-GJR- N-J	GJR-t	MF-GJR-t
$\mu$	0.0318	0.0393	0.0285	0.0287	0.0467	0.0472
	(0.0116)	(0.0113)	(0.0124)	(0.0132)	(0.0106)	(0.0105)
$ar{q}$	1.0918		1.0396		0.9571	
	(0.1093)		(0.1412)		(0.1896)	
$lpha_h$	0.0136	0.0005	0.0196	0.0147	0.0178	0.0001
	(0.0058)	(0.0194)	(0.0067)	(0.0094)	(0.0089)	(0.0269)
$\beta_h$	0.9040	0.6100	0.8943	0.9212	0.9124	0.7552
	(0.0041)	(0.0494)	(0.0075)	(0.0174)	(0.0076)	(0.0705)
$\gamma_h$	0.1308	0.2122	0.1211	0.0900	0.1110	0.1233
	(0.0067)	(0.0212)	(0.0112)	(0.0232)	(0.0132)	(0.0373)
$q_q$		1.0810		0.5825		1.4592
		(0.1773)		(0.0581)		(0.8374)
$\alpha_q$		0.0250		0.0008		0.0332
		(0.0076)		(0.0155)		(0.0144)
$\beta_q$		0.9519		0.8714		0.9487
		(0.0046)		(0.3938)		(0.0076)
$\gamma_q$		0.0289		-0.0002		0.0266
		(0.0123)		(0.0288)		(0.0244)
$\lambda_b$			0.0114	0.0065		
			(0.0035)	(0.0024)		
$\mu_r$			-1.6804	-2.6934		
			(0.8805)	(1.7811)		
$\sigma_r$			2.8360	3.6652		
			(0.2697)	(0.4840)		
$\eta$					6.8810	6.8291
					(0.5240)	(0.5641)
LL	16623	16663	16741	16748	16774	16785

### Table 16: Model Confidence Set *p*-Values and Model Ranking Discrete-time GARCH Models for Full Out-of-sample Period using Predictive Likelihood

This table shows model confidence set results for the full out-of-sample period January 3, 2007 to December 31, 2014 using predictive likelihood as the ranking criteria. For details of notation and calculation see Section 3.2 and Section 3.3. Discrete-time models are tested against SV-A, SV-G, SVSJ-A, SVSJ-G, SVYYD-A, and SVYYD-G. The first column indicates the number of the iterative elimination steps for models running from i = 1 to total number of models ( $m_0 = 12$ ). Second column shows the *p*-values for the hypotheses  $H_{0,\mathcal{M}_i}$  and third column presents MCS *p*-Value  $\hat{p}_{e_{\mathcal{M}_i}}$  for the model that is going to be eliminated in the respective elimination step. Fourth column shows the model eliminated in each iterative step by the elimination rule and thereby presents the model ranking according the MCS criteria, with the worst model ranked at the top and the best model at the bottom of the table, respectively. For a given significance level  $\alpha$  any model for which holds  $\hat{p}_{e_{\mathcal{M}_i}} \geq \alpha$  is included in the MCS  $\hat{\mathcal{M}}_{1-\alpha}^*$ .

Elimination Rule	<i>p</i> -Value for $H_{0,\mathcal{M}_i}$	MCS $p\text{-Value}\ \hat{p}_{e_{\mathcal{M}_i}}$	Eliminated Model
$e_{\mathcal{M}_1}$	0.0000	0.0000	MF-GJR-N
$e_{\mathcal{M}_2}$	0.0001	0.0001	MF-GJR-N-J
$e_{\mathcal{M}_3}$	0.0001	0.0001	GJR-N
$e_{\mathcal{M}_4}$	0.0012	0.0012	GJR-N-J
$e_{\mathcal{M}_5}$	0.0038	0.0038	SVYYD-A
$e_{\mathcal{M}_6}$	0.0023	0.0038	SVSJ-A
$e_{\mathcal{M}_7}$	0.0016	0.0038	MF-GJR-t
$e_{\mathcal{M}_8}$	0.0027	0.0038	SVYYD-G
$e_{\mathcal{M}_9}$	0.0055	0.0055	GJR-t
$e_{\mathcal{M}_{10}}$	0.0320	0.0320	SVSJ-G
$e_{\mathcal{M}_{11}}$	0.2308	0.2308	SV-A
$e_{\mathcal{M}_{12}}$	1.0000	1.0000	SV-G

### Table 17: Model Confidence Set *p*-Values and Model Ranking Discrete-time Models for Full Out-of-sample Period Using CRPS

This table shows model confidence set results for the full out-of-sample period January 3, 2007 to December 31, 2014 using continuous ranked probability score (CRPS) as the ranking criteria. For details of notation and calculation see Sections 3.2 and Section 3.3. The first column provides model specifications, the second column provides results for the non-weighted CRPS statistic. Columns 3 to 6 refer to the results for the weighted CRPS statistics. The weighting scheme "Center" applies more weight to the center of the predictive density when calculating CRPS and the weighting schemes "Tails", "Right Tail", and "Left Tail" work accordingly. For a given significance level  $\alpha$  models for which  $\hat{p}_{e_{\mathcal{M}_i}} \geq \alpha$  are included in the MCS  $\widehat{\mathcal{M}}^*_{1-\alpha}$ . We use \* (\*\*) to indicate that the model belongs to the 10% (25%) MCS.

Model Name	No Weight	Center	Tails	Right Tail	Left Tail
GJR-N	0.0090	0.0045	0.2849**	0.0022	$0.5315^{**}$
MF-GJR-N	0.0044	0.0016	0.0220	0.0022	$0.1777^{*}$
GJR-N-J	0.0090	0.0045	0.0555	0.0022	$0.3470^{**}$
MF-GJR-N-J	0.0090	0.0045	0.0220	0.0022	$0.2795^{**}$
GJR-t	0.0634	0.0088	$0.2849^{**}$	0.0022	$1.0000^{**}$
MF-GJR-t	0.0090	0.0088	0.0220	0.0022	$0.5315^{**}$
SV-A	0.0090	0.0088	0.0220	0.0022	$0.3470^{**}$
SV-G	$1.0000^{**}$	$1.0000^{**}$	$1.0000^{**}$	$1.0000^{**}$	$0.6755^{**}$
SVSJ-A	0.0090	0.0045	0.0220	0.0022	$0.2795^{**}$
SVSJ-G	0.0090	0.0167	0.0555	0.0056	$0.5315^{**}$
SVYYD-A	0.0090	0.0045	0.0220	0.0022	$0.2795^{**}$
SVYYD-G	0.0090	0.0088	0.0220	0.0022	$0.5315^{**}$

#### Table 18: VaR Specification Tests 1% (Full sample)

This table shows model confidence set results for the full out-of-sample period January 3, 2007 to December 31, 2014 using the asymmetric VaR loss function proposed by González-Rivera *et al.* (2004). For details of notation and calculation see Section 3.2 and Section 3.3. A subset of the most relevant models representing each of the model classes analyzed in this paper are tested against each other. The first column indicates the number of the iterative elimination step for models running from i = 1 to total number of models ( $m_0 = 18$ ). Second column shows the *p*-values for the hypotheses  $H_{0,\mathcal{M}_i}$  and third column presents MCS *p*-Value  $\hat{p}_{e_{\mathcal{M}_i}}$  for the model that is going to be eliminated in the respective elimination step. Fourth column shows the model eliminated in each iterative step by the elimination rule and thereby presents the model ranking according the MCS criteria, with the worst model ranked at the top and the best model at the bottom of the table, respectively. For a given significance level  $\alpha$ any model for which holds  $\hat{p}_{e_{\mathcal{M}_i}} \geq \alpha$  is included in the MCS  $\hat{\mathcal{M}}_{1-\alpha}^*$ .

Elimination Rule	<i>p</i> -Value for $H_{0,\mathcal{M}_i}$	MCS $p\text{-Value}\ \hat{p}_{e_{\mathcal{M}_i}}$	Eliminated Model
$e_{\mathcal{M}_1}$	0.1195	0.1195	MF-SV-A
$e_{\mathcal{M}_2}$	0.1217	0.1217	MF-GJR-N-J
$e_{\mathcal{M}_3}$	0.1629	0.1629	MF-SVJ-A
$e_{\mathcal{M}_4}$	0.1474	0.1629	MF-GJR-N
$e_{\mathcal{M}_5}$	0.1500	0.1629	MF-SVJ-G
$e_{\mathcal{M}_6}$	0.1415	0.1629	SVYYD-A
$e_{\mathcal{M}_7}$	0.1267	0.1629	MF-SV-G
$e_{\mathcal{M}_8}$	0.1227	0.1629	MF-SVSJJ-A
$e_{\mathcal{M}_9}$	0.0993	0.1629	SVYYD-G
$e_{\mathcal{M}_{10}}$	0.1271	0.1629	SVSJ-G
$e_{\mathcal{M}_{11}}$	0.1954	0.1954	SVSJ-A
$e_{\mathcal{M}_{12}}$	0.1942	0.1954	MF-SVSJJ-G
$e_{\mathcal{M}_{13}}$	0.2171	0.2171	GJR-N
$e_{\mathcal{M}_{14}}$	0.3186	0.3186	SV-A
$e_{\mathcal{M}_{15}}$	0.1507	0.3186	SV-G
$e_{\mathcal{M}_{16}}$	0.1836	0.3186	GJR-N-J
$e_{\mathcal{M}_{17}}$	0.2529	0.3186	MF-GJR-t
$e_{\mathcal{M}_{18}}$	1.0000	1.0000	GJR-t

are also reported.						
Parameter	$\mu$	$\kappa$	$\theta$	$\sigma_v$	$ ho_v$	$\gamma$
simulated	0.050	3.500	0.025	0.400	-0.650	0.500
estimated	0.041	3.796	0.025	0.399	-0.663	
RMSE	0.026	0.681	0.004	0.026	0.040	
standard	0.003	0.068	0.000	0.003	0.004	
error						
simulated	0.050	3.500	0.025	2.530	-0.650	1.000
estimated	0.043	4.381	0.023	2.404	-0.683	
RMSE	0.033	0.988	0.003	0.178	0.051	
standard	0.003	0.099	0.000	0.018	0.005	
error						
simulated	0.050	3.500	0.025	1.006	-0.650	0.750
estimated	0.042	3.850	0.024	1.051	-0.678	0.761
RMSE	0.032	0.830	0.004	0.308	0.046	0.073
standard	0.003	0.083	0.000	0.031	0.005	0.007
error						
simulated	0.050	3.500	0.025	3.658	-0.650	1.100
estimated	0.044	4.780	0.022	3.252	-0.683	1.070
RMSE	0.033	1.139	0.003	1.086	0.053	0.082
standard	0.003	0.114	0.000	0.109	0.005	0.008
error						

# Table 19: Simulation study: SV model.This table reports the parameter estimation results from a Monte Carlo study where 100 sample paths

with 4000 daily returns are simulated from the true model with parameters shown as *simulated*. The simulation is performed using an Euler discretization with 100 time steps per day. The average estimated parameter of these simulated paths are reported in line *estimated*. RMSE and standard errors (*std error*)

Table 20: Simulati	on study:	$\mathbf{SV}$	model.
--------------------	-----------	---------------	--------

This table reports the parameter estimation results from a Monte Carlo study where 100 sample paths with 4000 daily returns are simulated from the true model with parameters shown as *simulated*. The simulation is performed using an Euler discretization with 100 time steps per day. The average estimated parameter of these simulated paths are reported in line *estimated*. RMSE and standard errors (*std error*) are also reported.

Paramete	er $\mu$	$\kappa$	$\theta$	$\sigma_v$	$ ho_v$	$\lambda_c$	$\lambda_v$	$\mu_s$	$\sigma_s$	$\gamma$
sim	0.050	3.500	0.025	0.400	-0.650	0.500	30.000	-0.020	0.050	0.500
$\mathbf{est}$	0.046	3.765	0.025	0.400	-0.657	0.546	34.957	-0.024	0.043	
RMSE	0.026	0.722	0.004	0.028	0.044	0.439	26.084	0.024	0.014	
std err	0.003	0.072	0.000	0.003	0.004	0.044	2.608	0.002	0.001	
$\sin$	0.050	3.500	0.025	2.530	-0.650	0.500	30.000	-0.020	0.050	1.000
$\mathbf{est}$	0.051	4.440	0.023	2.415	-0.676	0.417	44.148	-0.024	0.044	
RMSE	0.029	1.049	0.003	0.218	0.057	0.382	35.525	0.020	0.014	
std err	0.003	0.105	0.000	0.022	0.006	0.038	3.552	0.002	0.001	
$\sin$	0.050	3.500	0.025	1.006	-0.650	0.500	30.000	-0.020	0.050	0.750
$\mathbf{est}$	0.050	3.862	0.024	1.052	-0.672	0.509	37.628	-0.023	0.044	0.761
RMSE	0.027	0.840	0.004	0.325	0.050	0.470	29.295	0.020	0.013	0.072
std err	0.003	0.084	0.000	0.032	0.005	0.047	2.929	0.002	0.001	0.007
$\sin$	0.050	3.500	0.025	3.658	-0.650	0.500	30.000	-0.020	0.050	1.100
$\mathbf{est}$	0.052	4.853	0.022	3.334	-0.680	0.424	42.383	-0.023	0.044	1.073
RMSE	0.032	1.343	0.003	1.232	0.072	0.429	32.784	0.025	0.015	0.091
std err	0.003	0.134	0.000	0.123	0.007	0.043	3.278	0.002	0.001	0.009



Figure 1: In-sample Sequential Likelihood Ratios.

These graphs show S&P 500 index returns in the upper part of the graphs and sequential likelihood ratios in the lower part of the graphs for the in-sample time period from January 2, 1987 to December 29, 2006. The left graph shows results for single factor jump-diffusion models and the right graph for the Lévy-jump models, respectively. Sequential likelihoods are calculated as a byproduct of the filtering procedure proposed by Bates (2006). All sequential likelihood ratios are calculated relative to the benchmark model SV-A.



Figure 2: Out-of-sample Sequential Likelihood Ratios.

These graphs show S&P 500 index returns in the upper part of the graphs and sequential likelihood ratios in the lower part of the graphs for the out-of-sample time period from January 3, 2007 to December 31, 2014. The left graph shows results for the estimated single factor jump-diffusion models and the right graph for the Lévy-jump models, respectively. Sequential likelihoods are calculated as a byproduct of the employed estimation procedure proposed by Bates (2006). All sequential likelihood ratios are calculated relative to the benchmark model SV-A.



## Figure 3: In-sample and Out-of-sample Sequential Likelihood Ratios for multi-factor jump-diffusion models.

The graphs show S&P 500 index returns in the upper part of the graphs and sequential likelihood ratios for the estimated multi-factor jump-diffusion models in the lower part of the graphs. The left graph shows the sequential likelihood ratios for the in-sample time period from January 2, 1987 until December 29, 2006 and the right graph for the Lévy-jump models for the out-of-sample time period January 3, 2007 until December 31, 2014.