# SV MIXTURE, CLASSIFICATION USING EM ALGORITHM

## Ahmed Hachicha

*University of Sfax, Department of Economic Development, Faculty of Economics and Management of Sfax, Airport road Km 4.5- B.P. 3018, Sfax, Tunisia*

## Fatma Hachicha

*Department of Finance. Faculty of Economics and Management of Sfax  Airport road Km 4.5- B.P. 3018, Sfax, Tunisia*

## Afif Masmoudi

*Department of Mathematics. Faculty of Sciences of Sfax Soukra Road Km.3.5-  B.P.1171 Sfax, Tunisia*

## ABSTRACT

*The present paper presents a theoretical extension of our earlier work entitled"A comparative study of two models SV with MCMC algorithm" cited, Rev Quant Finan Acc (2012) 38:479-493 DOI 10.1007/s11156-011-0236-1 where we propose initially a mixture stochastic volatility model providing a tractable method for capturing certain market characteristics. To estimate the parameter of a mixture stochastic volatility model, we first use the Expectation-Maximisation (EM) algorithm. The second step is to adopt the clustering approach to classify the database into subsets (clusters).*

**Keywords:** Mixture stochastic volatitlity model, Expectation-Maximization algorithm, clustering approach.

**JEL classification:** C51, C53, E37

## INTRODUCTION

The stochastic volatility (SV) model proposed by Taylor (1982) has become more and more popular for clearing up the behavior of financial time series such as stock prices and exchange rates. It consists of couple of equations that describe how the returns depend on the volatility. Stochastic volatility models (SV) are designed to fit this time-varying behaviour in the conditional variance of returns see Harvey, Ruiz et al. (1994).

Danielson (1994) and Sandmann and Koopman (1998) used classical inference methods simulation to SV model while, Melino and Turnbull (1990) utilized a generalized method (GMM) approach.

Liesenfeld and Richard (2003) used a maximum likelihood approach based upon efficient importance sampling. Recently, particle method like Monte Carlo have been applied to the SV model Manabu Asai (2008) and Hachicha (2012), where they apply the Metroplis hasting algorithm to Markov chain.

A simple single-factor SV model appears to be sufficient to capture extreme value, but mixture model provides a tractable method for capturing certain market characteristics. Several research until now, have concentrated on analyzing of the SV models with a mixture of two normal distributions, Asai, McAleer et al. (2006), Asai (2009).

Ningning and Zhidong (2011) prove that the expectation maximization (EM) algorithm is an efficient iterative method for finding maximum likelihood computationally and works well for unrounded data in most cases. Wojciech (2012) has recently proposed an estimator for convex optimization algorithms with a multivariate binary distribution that incorporates the correlation between individual variables. In addition, the EM algorithm is used for estimation of many probabilistic models, i.e., Finite Mixture Models (FMM), Gaussian Mixture Models (GMM) (Bilmes 1998), McLachlan and Basford (1988).

In our case, we deal with the fixed parameter problem and adopt the EM algorithm to determine the distribution of the latent variables in the next expected step. Then, we adopt the clustering approach. This paper proceeds as follows, we first propose in section 1 a mixture of two stochastic volatility models; the ARSV-t and the SVOL model. Secondly, we use in section 2 the Expectation maximization algorithm (EM) in order to estimate the parameter. Finally, we apply in section 3 the clustering in order to classify the database according to the first or the second model.

## Mixture Model

In order to present the mixture model we first recall the $M_1$ and $M_2$.

$M_1$ is named the p-th order ARSV-t model, ARSV(p)-t, and is given by:

$$\begin{cases} Y_t = \sigma \xi \exp(V_t / 2) \\ V_t = \phi_1 V_{t-1} + .... + \phi_p V_{t-p} + \eta_{t-1} \end{cases} \qquad (M_1)$$

$$\xi_t = \frac{\varepsilon_t}{\sqrt{\kappa_t /(\nu - 2)}}, \quad \kappa_t \approx \chi^2(\nu)$$

where $\kappa_t$ is independent of $(\varepsilon_t, \eta_t)$, $Y_t$ is the stock market indices, and $V_t$ is the log-volatility which is assumed to follow a stationarity AR(p) process with a persistent parameter $|\phi| \prec 1$. By this specification, the conditional distribution, $\xi_t$, follows the standardised t-distribution with zero mean and unit variance. Since $\kappa_t$ is independent of $(\varepsilon_t, \eta_t)$, the correlation coefficient between $\xi_t$ and $\eta_t$ is also $\rho$.

If $\phi \approx N(0,1)$, then;

$$\phi_1 = \frac{\left(\sum_{t=1}^{T} V_t V_{t-1}\right) - \phi_2 \left(\sum_{t=1}^{T} V_{t-1} V_{t-2}\right) + \overline{\phi_1}}{\left(\sum_{t=1}^{T} V_{t-1}^2\right) - 1}$$

Where,

$$\phi_2 = \frac{\left(\sum_{t=2}^{T} V_t V_{t-2}\right) - \phi_1 \left(\sum_{t=2}^{T} V_{t-1} V_{t-2}\right) + \overline{\phi_2}}{\left(\sum_{t=2}^{T} V_{t-2}^2\right) - 1}$$

The conditional posterior distribution of the volatility of $M_1$ for $\theta_1 = (\phi_1, \phi_2 .... \theta_p)$ is expressible as:

$$f_1(V \mid \theta_1, V) \propto e^{\left(\frac{1}{2*\sigma^2}\left(\sum_{t=1}^{T} Y_t^2 e^{-V_t}\right) - \frac{1}{2}\sum_{t=1}^{T}(V_t - \phi_1 V_{t-1} - \phi_2 V_{t-2})^2 - \frac{1}{2}\sum_{t=1}^{T}(V_{t+1} - \phi_1 V_t - \phi_2 V_{t-1})^2\right)}$$

The second model $M_2$ is named SVOL model of Jacquier, Polson et al. (1994) with normally distributed conditional errors and it is presented by:

$$\begin{cases} Y_t = \sqrt{V_t}\,\varepsilon_t^s \\ \log(V_t) = \alpha + \delta \log(V_{t-1}) + \sigma_v \varepsilon_t^v \end{cases} \qquad (M_2)$$

$$(\varepsilon_t^s, \varepsilon_t^v) \approx N(0, I_2)$$

The conditional posterior distribution of the volatility of $M_2$ for $\theta_2 = (\alpha, \delta, \sigma_v)$ is given by:

$$f_2(Y/\theta_2, V) \propto \frac{1}{V_t^{0.5}} \exp\left(\frac{-Y_t^2}{2V_t}\right) \frac{1}{V_t} \exp\left(-\frac{(\log V_t - \mu_t)^2}{2\sigma^2}\right)$$

Next, we combine the two models recall the $M_1$ and $M_2$. On one mixture model denoted by M, the density of the mixture model can be represented by the following formula:

$$f(y_t/\theta, M) = \pi_1 f_1(y/\theta_1, V, M_1) + \pi_2 f_2(y/\theta_2, V, M_2)$$

Which is equivalent to:

$$f(y_t/\theta, p) = \pi_1 f_1(y_t/\theta_1) + \pi_2 f_2(y_t/\theta_2)$$

Where $\theta$ is parameter estimation ($\theta_1$, $\theta_2$), with $\theta_1 = (\phi_1, \phi_2, ....\phi_p)$ and $\theta_2 = (\alpha, \delta, \sigma_v)$ and $y_t$ is the stock market indices. $\pi_1$ and $\pi_2$ denote respectively the proportion of belonging to the model $M_1$ and $M_2$ which satisfies the condition $\pi_1 + \pi_2 = 1$.

The formula of bayes applied for the mixture model gives that:

$$f(\theta, V/Y) \propto f(Y/V, \theta) P(\theta)$$

$$\propto \left[\pi_1 f_1(Y/\theta, V, M_1) + \pi_2 f_2(Y/\theta, V, M_2)\right] f_v(V/\theta) f_o(\theta)$$

Where,

$f_v$ is the density of the volatility model and $f_0$ is the prior density of $\theta$.

**Expectation- Maximization Algorithm**

In order to maximize $f(\theta, V)$, we use the EM (Expectation-Maximization) algorithm. However, using the EM technique, we can find locally MAP (maximum-a-posteriori) parameter estimates for the mixture model.

Let $Y = (y_1, y_2, ......y_T)$ be a sample of T independent observation from a mixture of two multivariate normal distributions of dimension d, and let $Z = (z_1, z_2, ......z_T)$ be the latent variables that determine the component from which observation originates.

We consider that:

$$y_i /(z_i = 1) \approx f_1(y_i / \phi_1, \phi_2, .... \phi_p)$$

$$y_i /(z_i = 2) \approx f_2(y_i / \alpha, \delta, \sigma_v)$$

Where,

$$P(z_i = 1) = \pi_1 \quad \text{and} \quad P(z_i = 2) = \pi_2 = 1 - \pi_1$$

Where the likelihood function is:

$$L(\theta, Y, Z) = P(Y, Z / \theta) = \prod_{i=1}^{n} \sum_{j=1}^{2} \mathbf{1}_{(z_i = j)} \pi_j f_j(y_i / \theta_i)$$

Where **1** is an indicator function and $f$ is the Probability Density Function (PDF).

Camila, Rignaldo et al. (2012) develop a local influence analysis for measurement error models inspired by the EM algorithm proposed by (Lachos, Vilca-Labra et al. 2010).

The EM algorithm is an efficient iterative process to compute the maximum likelihood (ML) by iteratively applying the following two steps: the E-step, and the M-step. In the E-step, the absent data are estimated given the observed data. This is achieved by the conditional expectation to evaluate the posterior assignment probabilities. In the M-step, the likelihood function is maximized under the hypothesis that is used instead of the actual missing data. The EM has become a fashionable instrument in statistical estimation problems concerning uncompleted data which can be posed in a similar form like mixture estimation.

Step E: use the current $\theta^{(t)}$, to estimate component membership of each unlabeled document, i.e.,

the probability that each mixture component generated each document.

$$Q(\theta / \theta^{(t)}) = E_{Z/Y, \theta}\left[\log L(\theta, Y, Z)\right]$$

Step M: re-estimate $\theta^{(t)}$. Use maximum a posteriori parameter estimation to find:

$$\theta^{(t+1)} = \arg\max Q(\theta / \theta^{(t)})$$

Let;

$$\tau_{j,i}^{(t)} = P(z_i = j / y_i, \theta^{(t)}) = \frac{\pi_j^t f(y_i, \theta_j^{(t)})}{\pi_1^{(t)} f_1(y) + \pi_2^{(t)} f_2(y)}$$

Where;
$$\pi_j^{(t)} = \frac{1}{n} \sum_{i=1}^{n} \tau_{j,i}^{(t)}$$

## The Clustering

The classification is used in most cases for classifying documents, news articles and web pages (Lewis and Gale 1994), Joachims (1998), Craven, DiPasquo et al. (1998).

After applying the EM algorithm, we obtain the estimation parameter. In order to classify the data in model $M_1$ or $M_2$, we opt for the hierarchical algorithm method and we compare the quantity

$\hat{\pi}_1 \hat{f}_1$ and $\hat{\pi}_2 \hat{f}_2$

Si $\hat{\pi}_1 \hat{f}_1(y_{t+1}) \succ \hat{\pi}_2 \hat{f}_2(y_{t+1})$ alors $y_{t+1} \approx \langle M_1 \rangle$

$$\text{Sinon } y_{t+1} \approx \langle M_2 \rangle$$

## CONCLUSION

Through this paper, we have presented a theoretical approach to estimate the parameters of the mixture model by applying Expectation-Maximization (EM) algorithm as a first step. Then, we have classified data according to stochastic volatility models to apply the technique of EM to resolve problems of system identification.

## REFERENCES

Asai, M. (2008). "Autoregressive stochastic volatility models with heavy-tailed distributions: A comparison with multifactors volatility models." Journal of Empirical Finance 15: 332-341.

Asai, M. (2009). "Bayesian analysis of stochastic volatility models with mixture-of-normal distribution." Mathematics and computers in simulation 79: 2579-2596.

Asai, M., M. McAleer, et al. (2006). "Multivariate stochastic volatility: A review." Econometric review 25: 145-175.

Bilmes, A. J. (1998). "Gentle tutorial on the EM algorithm and its application to parameter estimation for Gaussian mixture and hidden Markov models, Tech." International Computer Science Institute (ICSI). TR-97-021, USA.

Camila, B. Z., R. C. Rignaldo, et al. (2012). "On diagnostics in multivariate measurement error models under asymmetric heavy-tailed distributions." Statistical Papers 53: 665-683.

Craven, M., D. DiPasquo, et al. (1998). Learning to extract symbolic knowledge from the World Wide Web. . Proceedings of the Fifteenth National Conference on Artificial Intellligence (AAAI- 98): 792-799.

Danielson, J. (1994). "Stochastic volatility in asset prices: estimation with simulated maximum likelihood." Journal of Econometrics 61: 375-400.

Hachicha, e. a. (2012). "A comparative study of two models SV with MCMC algorithm Review." Journal of Quantitative Finance and Accounting 38: 479-493.

Harvey, A. C., E. Ruiz, et al. (1994). "Multivariate stochastic variance models." Review of Economic Studies 61: 247-264.

Jacquier, E., N. G. Polson, et al. (1994). "Bayesian analysis of stochastic volatility models." Journal of Business & Economic Statistics 12: 371-417.

Joachims, T. (1998). Text categorization with Support Vector Machines: Learning with many relevant features. In Machine Learning: ECML-98, Tenth European Conference on Machine Learning: 137-142.

Lachos, V. H., F. E. Vilca-Labra, et al. (2010). "Robust multivariatemeasurement error models based on scale mixtures of the skew-normal distribution." Statistics 44: 541–556.

Lewis, D. D. and W. A. Gale (1994). A sequential algorithm for training text classifiers. SIGIR '94: Proceedings of the Seventeenth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval: 3-12.

Liesenfeld, R. and J. F. Richard (2003). "Univariate and multivariate stochastic volatility models: estimation and diagnostics." Journal of Empirical Finance 10: 505-531.

McLachlan, G. J. and K. E. Basford (1988). Mixture Models: Inference and Applications to Clustering. New York, M. Dekker. .

Melino, A. and S. M. Turnbull (1990). "Pricing foreign currency options with stochastic volatility." Journal of Econometrics 45: 239-265.

Ningning, Z. and B. Zhidong (2011). Analysis of rounded data in mixture normal model. Statistical  Papers.

Sandmann, G. and S. Koopman (1998). "Estimation of stochastic volatility models via monte carlo maximum likelihood." Journal of Econometrics 87: 271-301.

Taylor, S. J. (1982). Financial returns modelled by the product of two stochastic processes — a study of the daily sugar prices 1961–75.

Wojciech, G. (2012). Maximum likelihood estimation for ordered expectations of correlated binary variables. Statistical Papers.