Matematika, 2004, Jilid 20, bil. 1, hlm. 31–41 ©Jabatan Matematik, UTM.

Principal Component Analysis in Modelling Stock Market Returns

Kassim Haron & Maiyastri

Department of Mathematics, Universiti Putra Malaysia, 43400 UPM Serdang, Selangor, Malaysia

Abstract In this study, an alternative method to compare the performance of several GARCH models in fitting the KLCI daily rate of return series before and after the Asian financial crisis in 1997 using Principal Component Analysis (PCA) is sought. Comparison is then made with the results obtained from a known method based on the ranks of the Log Likelihood (Log L), Schwarzs Bayesian Criterion (SBC) and the Akaike Information Criterion (AIC) values. It is found that the best and the worst fit models identified by both methods are exactly the same for the two periods but some degree of disagreement, however, existed between the intermediate models. We also find that the proposed method has a clear edge over its rival because PCA uses actual values of the three criteria and hence the inability to exactly specify the relative position of each of the competing models as faced by the ranking method may be avoided. Another plus point is this method also enables models to be classified into several distinct groups ordered in such a way that each group is made up of models with nearly the same level of fitting ability. The two extreme classes of models are identified to represent the best and the worst groups respectively.

Keywords GARCH models, Returns, Fitting, Rank, Principal component

Abstrak Dalam kajian ini satu kaedah alternatif untuk membandingkan pencapaian beberapa model GARCH bagi penyuaian siri kadar pulangan harian KLCI sebelum dan selepas krisis kewangan Asia pada tahun 1997 menggunakan Analisis Komponen Prinsipal dicari. Perbandingan kemudiannya dibuat dengan keputusan yang diperoleh daripada kaedah yang diketahui berasaskan kepada pangkat nilai Log Likelihood (Log L), Schwarzs Bayesian Criterion (SBC) dan Akaike Information Criterion (AIC). Didapati bahawa model penyuaian terbaik dan terburuk yang dikenalpasti menggunakan kedua-dua kaedah adalah tepat sama bagi dua tempoh masa itu tetapi beberapa percanggahan, bagaimana pun, wujud di antara model pertengahan. Kami juga dapati kaedah yang dicadangkan mempunyai kelebihan yang ketara ke atas lawannya kerana PCA menggunakan nilai sebenar bagi ketiga-tiga kriteria dan oleh itu ketidakupayaan untuk menyatakan dengan tepat kedudukan secara relatif setiap model yang bersaing seperti yang dihadapi oleh kaedah pangkat dapat dielakkan. Kelebihan lain ialah kaedah ini juga dapat mengklasifikasikan model ke dalam beberapa kumpulan berbeza, disusun sedemikian rupa supaya setiap kumpulan terdiri daripada model dengan paras kebolehan penyuaian yang hampir sama. Dua kelas model ekstrim masing-masing dikenalpasti mewakili kumpulan terbaik dan terburuk.

Katakunci Model GARCH, Pulangan, Penyuaian, Pangkat, Komponen prinsipal

1 Introduction

Financial time series data such as stock return, inflation rates, foreign exchange rates have non-normal characteristics: leptokurtic and skew. In addition, they exhibit changes in variance over time. In such circumstances, the assumption of constant variance (homoscedasticity) is inappropriate. The variability in the financial data could very well be due to the volatility of the financial markets. The markets are known to be sensitive to factors such as rumours, political upheavals and changes in the government monetary and fiscal policies. Engle [4] introduced the Autoregressive Conditional Heteroscedasticity (ARCH) process to cope with the changing variance. Bollerslev [1] proposed a General ARCH (GARCH) model which has a more flexible lag structure because the error variance can be modelled by an Autoregressive Moving Average (ARMA) type of process. Such a model can be effective in removing the excess kurtosis. Nelson [8] proposed a class of exponential GARCH (EGARCH) model which can capture the asymmetry and skewness of the stock market return series. Several researchers such as Franses and Van Dijk [6], Choo [3] and Gokcan [7] had shown that models with a small lag like GARCH (1,1) is sufficient to cope with the changing variance. Nevertheless, due to the high volatility of the rate of returns of the KLCI, higher order lag models such as the GARCH (1,2), GARCH (2,1) and GARCH (2,2) will also be included in our study. In all, we shall compare the performance of eleven competing time series models for fitting the rate of returns data. The models are the ARCH (1), ARCH (2), GARCH (1,1), GARCH (1,2), GARCH (2,1), EGARCH (1,1), GARCH-NNG (1,1), SGARCH (1,1), IGARCH (1,1) and GARCH-M (1,1) and GARCH (2,2).

Franses and Van Dijk [6] and Choo [3] chose the best models for fitting time series data by comparing the ranks of the values of three goodness-of-fit statistics namely the Log Likelihood (Log L), Schwarz's Bayesian Criterion (SBC) and the Akaike Information Criterion (AIC) respectively. We notice that such a method has one weakness. The use of an ordinal scale type of measurement, that is the rank, instead of the actual values in calculating the criteria would cause some loss of information. Therefore, in this paper an alternative method capable of overcoming such a problem is sought and used to compare the performance of potential models for fitting the return series for both periods. We proposed the use of Principal Component Analysis (PCA) procedures to produce a new set of variables called the principal components formed by a linear combination of the three statistics. By looking at the component values together with the PCA plots, one would be able to classify the relative performance of the competing models. Finally, the results from the two methods are studied and compared.

2 Data Description

The data used in this study is the daily rate of returns of the KLCI (Kuala Lumpur Composite Index) registered from week 1 of January 1989 to week 4 of December 2000. In the fourth quarter of 1997, financial crisis which hit the Asian region had badly hurt the performance of most of the stock markets including the KLSE (Kuala Lumpur Stock Exchange). From the plot in Figure 1, we can see that starting from September 1997, the rate of returns of the KLCI were noticeably volatile.

For this reason, we shall divide the data into two periods:

Period I : From January 1989 to September 1997

Period II : From October 1997 to December 2000.



Figure 1: Daily rate of returns of the KLCI fro January 1989 to Decemver 2000

Some descriptive statistics for the daily return of the KLCI are presented in Table1.

Period Ν Mean Standard Variance Skewness Kurtosis $(x10^{-4})$ Deviation (x10⁻⁴) 2128 4.24 0.012005 1.44 -0.5471 11.13897 Ι Π 770 0.027812 7.74 0.3862 18.20129 -1

Table 1: Summary statistics of the rate of daily returns of the KLCI

As seen from Table 1, the distribution of the rate of daily returns in Period I is negatively skewed and leptokurtic. However, for Period II, that is after the financial crisis, the standard deviation of the data is about twice as large as that in Period I. This result, according to Gokcan [7] indicates the rate of returns in Period II is more volatile than in Period I. As a result, the data have a positive skew and leptokurtic distribution.

3 The Models

The daily rate of returns r_i of the KLCI are calculated using the following formula:

$$r_i = \log\left(\frac{I_t}{I_{t-1}}\right), \ t = 1, 2, \dots, T$$

where I_t denotes the reading on the composite index at the close of t^{th} trading day. As noted earlier, the rate of daily returns of the KLCI displays a changing variance over time. There are many ways to describe the changes in variance and one of them is by considering the Autoregressive Conditional Heteroscedasticity (ARCH) model.

The ARCH regression model for the series r_t can be written as $\phi_m(B)r_t = \mu + \varepsilon_t$, for the model with intercept and $\phi_m(B)r_t = \varepsilon_t$, for the non-intercept model, with

$$\phi_m(B) = 1 - \phi_1 B - \phi_2 B^2 - \dots - \phi_m B^n$$

where B is the backward shift operator defined by $B^k X_t = X_{t-k}$. The parameter μ reflects a constant term (intercept) which in practice is typically estimated to be close or equal to zero. The order m is usually 0 or small, indicating that there are usually no opportunities to forecast r_t from its own past. In other words, there is never an autoregressive process in r_t .

The conditional distribution of the series of disturbances which follows the ARCH process can be written as

$$\varepsilon|_{\Phi_{\tau}} \sim N(0, h_t) \tag{1}$$

where Φ_{τ} denotes all available information at time $\tau < t$. The conditional variance h_t is

$$h_t = \omega + \sum_{i=1}^q \alpha_i \varepsilon_{t-1}^2$$
(2)
$$\varepsilon_t = \sqrt{h_t} e_t, \quad e_t \sim N(0, 1).$$

Bollerslev [1] introduced a Generalized ARCH(p,q) or GARCH(p,q) model where the conditional variance h_t is given by

$$h_{t} = \omega + \sum_{i=1}^{q} \alpha_{i} \varepsilon_{t-1}^{2} + \sum_{j=1}^{p} \beta_{j} h_{t-j}, \ p \ge q, q > 0 \ \text{and} \ \omega > 0, \alpha_{i} > 0, \beta_{j} \ge 0$$
(3)

If the parameters are constrained such that

$$\sum_{i=1}^{q} \alpha_i + \sum_{j=1}^{p} \beta_j < 1$$

we have a weakly stationary GARCH(p,q) or SGARCH(p,q) process since the mean, variance and autocovariance are finite and constant over time. If

$$\sum_{i=1}^q \alpha_i + \sum_{j=1}^p \beta_j = 1$$

we then have the Integrated GARCH(p, q) or IGARCH(p, q) model.

Nelson [11] proposed a class of exponential GARCH or EGARCH models. In this model h_t is defined by

$$\ln(h_t) = \omega + \sum_{i=1}^{q} \alpha_i g(e_{t-i}) + \sum_{j=1}^{p} \beta_j \ln(h_{t-i})$$

where

$$g(e_t) = \theta e_t + \gamma |e_t| - \gamma E |e_t|.$$

The coefficient of the second term in $g(e_t)$ is set to be $1(\gamma = 1)$ in this formulation. Unlike the linear GARCH model there are no restrictions on the parameters to ensure non-negativity of

the conditional variances. The EGARCH model allows good news (positive return shocks) and bad news (negative return shocks) to have a different impact on volatility whereas the linear GARCH model does not. If $\theta = 0$, a positive return shock has the same effect on volatility with the negative return shock of the same amount. If $\theta < 0$, a positive return shock actually reduces volatility and if $\theta > 0$, a positive return shock increases volatility. The conditional variance (h_t) follows equation (3) and we write the model as EGARCH(p, q).

In the GARCH-in-Mean or GARCH-M model, the GARCH effects appear in the mean of the process, given by $\varepsilon_t = \sqrt{h_t} e_t$ where $e_t \sim N(0, 1)$ and $r_t = \mu + \delta \sqrt{h_t} + \varepsilon_t$ for the model with intercept and $r_t = \delta \sqrt{h_t} + \varepsilon_t$ for the non-intercept model. Engle and Mustafa [5] reported there is a significant test statistics for ARCH model specially for stock returns. For the model GARCH(p, q) specification, Bollerslev [2] suggested to adopt low orders for the lag lengths p and q. Typical examples are the GARCH(1,1), GARCH(1,2) and GARCH (2,1).

4 Research Method

The parameters of the models considered in this study are estimated using the maximum likelihood method. The likelihood function for the ARCH and GARCH models can be written as follows:

$$L = -\frac{T}{2}\log(2\pi) - \frac{1}{2}\log\sum_{t=1}^{T}\log(h_t) - \frac{1}{2}\sum_{t=1}^{T}\frac{\varepsilon_t^2}{h_t}$$

where

$$T=$$
 total number of daily rate of returns
$$\label{eq:expansion} \varepsilon_t = r_t$$

and h_t is the conditional variance. However, when estimating GARCH-M(p,q) model, we take $\varepsilon_t = r_t - \delta \sqrt{h_t}$.

The procedures for performing the ranking and the proposed PCA methods are as follows.

(a) The Ranking Method

In this method, the values of the three goodness-of-fit statistics, namely the Log Likelihood (Log L), Schwarzs Bayesian Criterion (SBC) and the Akaikes Information Criterion (AIC) are first calculated for each rival models. The values of AIC and SBC are computed as follows;

$$AIC = -2\ln(L) + 2k$$
$$SBC = -2\ln(L) + \ln(T)k$$

where k is the number of free parameters and T is the number of residuals that can be computed for the time series. The ranks of the models based on the calculated values are then determined for each statistic. Finally, the ranks of the models are averaged across the statistics and the models with the smallest and the largest averages are taken to be the best and the worst fit models respectively.

(b) The PCA method

The principal component $_{n}\mathbf{Y}_{q} = (\mathbf{y}_{1}, \mathbf{y}_{2}, \dots, \mathbf{y}_{q})$ is obtained by multiplying the data matrix $_{n}\mathbf{X}_{p}$ with $_{p}\mathbf{A}_{q}$, (Johnson and Wichern [10]), where $_{p}\mathbf{A}_{q}$ is a matrix resulted from Singular Value Decomposition (SVD) of matrix $_{n}\mathbf{X}_{p}$, that is

$$_{n}\mathbf{X}_{p} = {}_{n}\mathbf{U}_{q}\mathbf{L}_{q}\mathbf{A}_{p}^{\prime}$$

where

 $_{q}\mathbf{L}_{q} =$ diagonal matrix $q \times q$ non zero characteristics root $\mathbf{X}'\mathbf{X}$

 $_{p}\mathbf{A}_{q} = (\mathbf{a}_{1}, \mathbf{a}_{2}, \dots, \mathbf{a}_{q})$ in which $\mathbf{a}_{j}, j = 1, \dots, q$ are vector characteristics of $\mathbf{X}'\mathbf{X}$.

The PCA enables a new set of variables called the principal components or components formed by a linear combination of the old variables to be produced. In this case, the components formed are a linear combination of the Log L, SBC and AIC statistics. The components of PCA are ordered by their importance in such a way that the first component contains more information than the second, the second component is better than the third and so on. The classification of models is obtained from the PCA plots.

5 Results and Discussion

(a) Period I

Table 2 shows the parameter estimates and the values of *t*-ratio. All parameter estimates, with the exception of β_1 for GARCH (2,2) and θ for EGARCH (1,1) are significant at 5% level.

Results of the goodness-of-fit test are presented in Table 3.

Since the proportion of variation for the first component in the PCA is found to be 99.8%, only the first component would be used to determine model performances. Figure 2 shows the resulting PCA plot together with the component values.

The results presented in Table 3 and Figure 2 clearly suggest that GARCH(1,2) is the best fit whilst ARCH(1) is the worst fit models for both methods. However, the performances of the intermediate models displayed some degree of disagreement between the two methods. Figure 2 also shows that the models GARCH(1,2), GARCH(2,1) and GARCH(2,2) have almost the same level of performances and may therefore form a group consisting of models which are better than the others. They are followed by GARCH-M(1,1), GARCH(1,1), GARCH-NNG(1,1) and SGARCH(1,1) which formed the second group. The worst models which belong to the ARCH family, namely the ARCH(2) and ARCH(1), formed groups which are clearly separated from the others. In other words, the GARCHswith lag 2 formed the best group followed by the GARCHs with lag 1, with the exception of EGARCH(1,1) and IGARCH(1,1). Also note that, if we were to use the ranking method for choosing the models, we would not be able to identify which models exhibit about the same level of performances and which are clearly different.

(b) Period II

Table 4 shows the parameter estimates and the values of *t*-ratio. All parameter estimates, with the exception of β_1 for GARCH(2,2), β_2 for GARCH (2,1) and δ for GARCH-M(1,1) are significant at 5% level. Results of the goodness-of-fit test are presented in Table 4.

Model	Ø	t-ratio	α_1	t-ratio	ρ_1	t-ratio	α_2	t-ratio
	X10 ⁻⁵							
ARCH(1)	9.1203	65.562	0.4201	14.547	-	-	-	-
ARCH(2)	7.8234	61.592	0.291	9.029	-	-	0.190	9.099
GARCH(1,1)	1.6686	8.872	0.181	10.212	0.706	26.801	-	-
GARCH(1,2)	1.1268	6.146	0.279	8.371	0.788	25.941	-0.139	3.994
GARCH(2,1)	1.6130	8.530	0.248	10.359	0.030	2.076	-	-
EGARCH(1,1)	-1.252	-5.852	0.336	7.736	0.858	35.870	-	-
G-NNG(1,1)	1.6685	8.870	0.181	10.211	0.706	26.799	-	-
SGARCH(1,1)	1.6688	8.876	0.181	10.220	0.707	26.804	-	-
IGARCH(1,1)	0.9753	8.745	0.258	12.399	0.742	35.737	-	-
G-M(1,1)	1.6582	8.937	0.176	10.107	0.711	27.320	-	-
GARCH(2,2)	1.8143	7.820	0.195	8.368	0.000	0.0000	0.041	1.940
Model	β_2	t-ratio	δ	t-ratio	θ	t-ra	tio	
GARCH(2,1)	0.627	20.564						
GARCH(2,2)	0.640	14.951						
EGARCH(1,1)	-	-	-	-	-0.1	064 -1.2	268	
GARCH-M(1,1)) -	-	0.047	2.231			-	

Table 2: Estimation results of the daily rate of returns for Period I

Table 3: Performance by ranking the average rank of the goodness-of-fit statistics values for Period I

Model	LOGL	SBC	AIC	Rank	Rank	Rank	Avg.	Rank
				LOGL	SBC	AIC	Rank	Avg.
ARCH(1)	6538.139	-13061.0	-13072.3	11	11	11	11	11
ARCH(2)	6584.430	-13145.9	-13162.9	10	10	10	10	10
GARCH(1,1)	6620.209	-13217.4	-13234.4	6	3	5	5	4
GARCH(1,2)	6626.179	-13221.7	-13244.4	2	1	2	2	1
GARCH(2,1)	6625.998	-13221.3	-13244.0	3	2	3	3	2.5
EGARCH(1,1)	6615.444	-13200.2	-13222.9	8	8	8	8	8
G-NNG(1,1)	6620.211	-13217.4	-13234.4	5	4	6	5	5
SGARCH(1,1)	6620.206	-13217.4	-13234.4	7	5	7	6	7
IGARCH(1,1)	6602.353	-13189.4	-13200.7	9	9	9	9	9
G-M(1,1)	6622.647	-13214.6	-13237.3	4	7	4	5	6
GARCH(2,2)	6627.753	-13217.2	-13245.5	1	6	1	3	2.5



Figure 2: PCA plot for period I

Table 4: Estimation result of the	daily rate of returns	for Period II
-----------------------------------	-----------------------	---------------

Model	ø	t-ratio	α_1	t-ratio	β_1	t-ratio	α_{2}	t-ratio
	X10-5							
ARCH(1)	39.2	27.537	0.554	11.861	-	-	-	-
ARCH(2)	21.4	14.656	0.185	3.692	-	-	0.66	15.67
GARCH(1,1)	1.883	5.638	0.164	6.903	0.815	37.550	-	-
GARCH(1,2)	2.396	5.907	0.113	2.535	0.777	29.398	0.085	1.743
GARCH(2,1)	1.766	3.230	0.152	3.117	0.921	2.996	-	-
EGARCH(1,1)	-0.203	-2.495	0.271	5.615	0.972	90.101	-	-
G-NNG(1,1)	1.895	5.679	0.160	6.821	0.816	37.858	-	-
SGARCH(1,1)	1.523	5.160	0.176	8.784	0.821	40.656	-	-
IGARCH(1,1)	1.659	5.134	0.189	9.207	0.811	39.477	-	-
G-M(1,1)	1.870	5.598	0.165	6.853	0.815	36.919	-	-
GARCH(2,2)	4.783	6.744	0.091	2.979	0.041	0.675	0.258	7.196
	в	t votio	5	t votio	۵	t natio		
Model	ρ_2	1-rano	0	1-14110	0	t-ratio		
GARCH(2,1)	-0.093	-0.356	-	-	-	-		
GARCH(2,2)	0.553	8.623	-	-	-	-		
EGARCH(1,1)	-	-	-	-	-0.154	-2.193		
GARCH-M(1,1))-	-	0.026	0.687	-	-		

Since the proportion of variation for the first component in the PCA is found to be 99.6%, only the first component would be used to determine model performances. Figure 3 shows the resulting PCA plot together with the component values.

The results presented in Table 5 and Figure 3 clearly suggest that SGARCH(1,1) is the best fit whilst ARCH(1) is the worst fit model for both methods. However, the performances of the intermediate models also displayed some degree of disagreement between the two methods Figure 3 shows that the GARCH family models are separated from the ARCH in such a way that the GARCHs are better than the ARCHs. Generally, the GARCH family can be divided into three groups with the best consists of SGARCH(1,1) and IGARCH(1,1) followed by the second group which is made up of the GARCH(1,1), GARCH-NNG(1,1) and EGARCH(1,1). The rest of the models form the third group.

Results from the two periods revealed that the ranking method has one weakness. By using measurements in ordinal scale in calculating the values of the three criteria, one tends to lose some information regarding the relative position of the models concerned. This is obvious in Table 3 where two models, GARCH(2,1) and GARCH(2,2), have a tied rank of 2.5 whereas the PCA method firmly singled out the former as the second while the latter as the third best models.

6 Conclusions

This study managed to come out with an alternative method for selecting the best model from a set of competing GARCH models for fitting the stock market return series. The PCA method identified exactly the same best and worst fit models as the ranking method for the two periods. However, as a whole, the models occupying the intermediate positions differ in the two methods. The proposed method is seen to be superior and should be preferred because PCA uses actual values of the three goodness-of-fit statistics and hence the inability to exactly specify the relative position of each of the competing models as faced by the ranking method may be avoided. Another advantage is this method also enables models to be classified into several distinct groups ordered in such a way that each group is made up of models with about the same level of fitting ability. The two extreme classes of models are identified to represent the best and the worst groups respectively.

References

- T.Bollerslev, Generalized autoregressive conditional heteroscedasticity, Journal of Econometrics, 31, (1986), 307–27.
- [2] T. Bollerslev, Y.R. Chou & F.K. Kroner, ARCH modeling in finance, Journal of Econometrics, 52, (1992), 5–59
- [3] W. C. Choo, et al, *Performance of GARCH Models in Forecasting Stock Market Volatility*, Journal of Forecasting, 18, (1999), 333–334
- [4] R. F. Engle, Autoregressive conditional heteroscedasticity with estimates of variance of United Kingdom Inflation, Econometrica, 50,(4) (1982), 987–1008.
- [5] R. F. Engle & C. Mustafa, Implied GARCH models from options prices, Journal of Econometrics, 52, (1992), 289–311



Figure 3: PCA plot for period II

Model	LOGL	SBC	AIC	Rank	Rank	Rank	Mean	Rank
				LOGL	SBC	AIC	Rank	Model
ARCH(1)	1778.081	-3542.87	-3552.16	11	11	11	11.00	11
ARCH(2)	1855.175	-3690.41	-3704.35	10	10	10	10.00	10
GARCH(1,1)	1871.577	-3723.22	-3737.15	7	4	4	5.00	4
GARCH(1,2)	1872.342	-3718.10	-3736.68	4	7	6	5.67	5
GARCH(2,1)	1871.741	-3716.90	-3735.48	6	9	9	8.00	9
EGARCH(1,1)	1875.836	-3725.09	-3743.67	2	3	3	2.67	2
G-NNG(1,1)	1871.534	-3723.13	-3737.07	8	5	5	6.00	6
SGARCH(1,1)	1875.733	-3731.53	-3745.47	3	1	1	1.67	1
IGARCH(1,1)	1870.289	-3727.29	-3736.58	9	2	7	6.00	7
G-M(1,1)	1871.869	-3717.15	-3735.74	5	8	8	7.00	8
GARCH(2,2)	1876.881	-3720.53	-3743.76	1	6	2	3.00	3

Table 5: Performance by ranking the average rank of the goodness-of-fitstatistics values for Period II

- [6] P. H. Franses & R. Van Dijk, Forecasting stock market volatility using (non-linear) Garch models, Journal of Forecasting, 15, (1996), 229–35.
- S. Gokcan, Forecasting volatility of emerging stock markets: Linear versus Non-Linear GARCH models, Journal of Forecasting, 19, (2000), 499–504.
- [8] J. D. Hamilton, Time Series Analysis, Princeton, New Jersey, 1994.
- [9] A. C. Harvey, Time Series Model, Harvester, New York, 1993
- [10] R. A. Johnson & D. W. Wichern, Applied Multivariate Statistical Analysis, Prentice Hall, New Jersey, 1988.
- [11] D. B. Nelson, Conditional heteroscedasticity in asset returns: A new approach, Econometrica 52, (2) (1991) 347–70.