# INTERNATIONAL JOURNAL OF ENGINEERING SCIENCES & MANAGEMENT
## FUZZY PREDICTION METHOD FOR STOCK PRICE BASED ON AN ENHANCED BOOSTING ALGORITHM

**Qiansheng Zhang[*1], Ziqi Wu[2], Qiting Chen[3] & Lushuo Wei[4]**
[*1,2,3&4]School of Finance, Guangdong University of Foreign Studies, Guangzhou 510006, China

## ABSTRACT

An enhanced boosting algorithm is proposed by employing triangular fuzzy number in the traditional boosting learning machine. Using the enhanced boosting algorithm can predict stock price by selecting some important predictable variables. And the stock prediction price can be modified by using fuzzy possibility expectation theory. Finally, the empirical analysis of stock price prediction for PING AN BANK of China by utilizing the enhanced boosting regression algorithm illustrates that MSE and MAPE prediction deviations decrease greatly, which improve the prediction performance of boosting algorithm without optimal arguments

**Keywords:** *Boosting Algorith Fuzzy Number Multi-criteria Regression Stock Price Prediction.*

## I.    INTRODUCTION

With the development of networks and the arrival of the big data era，it becomes more easier to predict stock price based on data mining and analysis in stock market and multi-criteria regression of various variables associated with t heestimated stock price.  As is well known, the OLS regression analysis is widely used in forecasting stock price. Th e regressed stock price Y can be evaluated as $Y = \beta_1 X_1 + \beta_2 X_2 + \ldots\ldots + \beta_n X_n + \varepsilon$, where $\beta = (\beta_1, \beta_2, \ldots\ldots, \beta_n)$ are t he regressionparameters, while the ε denotes the error term.  In the OLSregression, there may be correlations betwee n the  selected variables which increase the possibility of the multicollinearity problem leading to greater variance in regression estimation. Therefore, it should be serious to select the significant regression variable to estimate the stoc k price. Among them, the coefficient compression method is a commonly used one to select correct variable.Theesse nce of this method is an appropriate loss function,which is determined by$R(\beta) + \sum_{j=1}^{p} p_\lambda \left( \left| \beta_j \right| \right)$, where the penalty f unction $p_\lambda$ differs depending on different $\beta_j$ (j = 1,2,…,p) and different penalty parameter $\lambda$.  Recently, Hoerl[1] intr oduced theridge analysis to regression problem where the penalty function is constructed to compress the variable th rough constructing the penalty function $L_p (0 < p < 2)$. Tibshirani[2] also studied the Lasso regression where the pe nalty function is the least squares estimation penalty function $L_1$ (p=1).The construction of these two kinds of penalt y functions help to congress the coefficient,leading to regression with better fitting effect.

Ensemble Learning has attracted much attention in recent years, as it can improve the generalization ability of a system, and thus can decline the deviation in the process of single learner. In the field of ensemble learning, Breiman[3]introduced the classification tree and regression tree based on bagging alogorithm,and Friedman [4] introduced the boosting algorithm with gradient descent. Bu ̈hlmann[5-7] conducted the study on different models based on boosting algorithm. More recently, Breiman[8] studied on random forest based on multiple decision trees. Bagging, boosting, and random forest are the three most important ensemble learning algorithms in the construction of regression ensemble algorithms. Many researchers used ensemble learning in regression forecasting and achieved some good results. For example, Zhang and Ma [9] BP-Boosting regression algorithm used it to forecast the price of estate, while Chen and Jin[10] predict the price change of house estate with random forest. In addition, Wang[11] used boosting-ARMA forecasting algorithm to forecast the price of carbon. Thus ensemble learning can deal with large and complex data, and get higher precision data classification and regression results. It can be widely used in classification and regression problems.

In the real stock market, because the investors tend to hold shares with their subjective intention which would cause the future stock price to fluctuate and lead to deviation from the equilibrium, the future stock price has a fuzzy uncertainty, which means that there would be large deviation if we get to predict stock price with an accurate result. In 1965,Zadeh[12] introduced fuzzy set concept, which is widely used to describe the uncertainty of objects. After that, some scholars such as Zhang[13] and Li [14] applied the fuzzy set theory into financial forecasting and achieved some intresting results.In this paper, we will propose an enhanced boosting algorithm in regression

combined with fuzzy set theory to predict the stock price tend during certain period with the selection of non-optimal independent variables.

## II.    BOOSTING ALGORITHM

Boosting algorithm is a kind of ensemble learning which can improve the accuracy of weak classification algorithm. In this process, the total sample was divided, and then get the corresponding sub-sample with the original distribution. Through training, the machine would generate a series of weak learners. After reaching a certain times of circulation or achieving a metric, all weak learners are weighted to get a strong learner to predict the data.In the construction of boosting regression learning, it is significant to construct a reasonable loss function. This paper applied the commonly used ordinary least squares method, where the sample obeys Gauss distribution. Based on this hypothesis, the principle of regression optimization is to minimize the sum of squares of residuals.

We divide the squared loss by two so that the first derivative of the model is exactly the residual of the model and does not affect the solution of loss function to *F\*(x)*. Therefore, under n-dimensional samples, the loss function is

$$L(y, F(x)) = \sum_1^n (y_i - F(x_i))^2/2 \,,$$

where $(y_i - F(x_i))$ is residual, and $L(y, F(x))$ is the sum of the square of residuals.

Thus the result we would like to get is the estimated value of *F\*(x)* based on the loss function *L(y , F(x))*, and *F\*(x)* follows that:

$$F^*(x) = \arg\min_F \quad E_{y,x} L(y, F(x)) = \arg\min_F E_x[E_y(L(y, F(x)))|x],$$

where *F(x)* is a parametric function *F(x , P)*, and *P* is a limited parameter set $\{\beta_m, a_m\}_1^M$. In addition, *F(x , P)* is an additive model, whose general form is:

$$F(x , P) = \sum_{m=1}^M \beta_m h(x, a_m),$$

where $h(x, a_m)$is the single-variable function with independent variable x, and $a_m$ represents parameters including classificationvariable, classificationrule and so on, while $\beta_m$ means corresponding weight.

In practical solution, greedy stagewise is used for the piecewise solution of objective function $F_m(x) = F_{m-1}(x) + \beta_m h(x, a_m)$.

Since it is hard to obtain the base learner of the above function, we can find the parameter *a*of the base learner by constructing the unconstrained gradient. For any given approximation function$F_{m-1}(x)$, the function $\beta_m h(x, a_m)$ can be considered as the best greedy step of objective function. Therefore, it can be regarded as steepest descent under the constraint condition. By constructing an approximate unconstrained negative gradient $-g_m(x_i) = -[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)}]_{F(x)=F_{m-1}(x)}$, we can get the best steepest-descent direction $-g_m = \{-g_m(x_i)\}_1^N$ in the N-dimensional data space at $F_{m-1}(x)$. Based on $-g_m$, we can compute the parameter $a_m$ of the base learner $h(x, a_m)$ that can be extended to any *x* value

$$a_m = \arg\min_{a,\beta} \sum_{i=1}^N [-g_m(x_i) - \beta h(x_i, a)]^2.$$

After getting the parameter $a_m$ of the base learner $h(x, a_m)$, we can get the corresponding weight$\rho_m$ of base learner $h(x, a_m)$ based on loss function below.

$$\rho_m = \arg\min_\rho \sum_{i=1}^N L(y_i, F_{m-1} + \rho h(x_i, a_m)).$$

Estimated function is changed to$F_m(x) = F_{m-1}(x) + \rho_m h(x, a_m)$.

Finally initializing it $F_0(x) = \arg\min_\rho \sum_{i=1}^N L(y_i, \rho)$. When m reaches our target step, $F_m(x)$ is the target function *F\*(x)* we want to obtain.

## III. THE ENHANCED BOOSTING ALGORITHM BASED ON TRIANGULAR FUZZY NUMBER

In the practical application of boosting algorithm to predict stock price, investors are difficult to accurately obtain the optimal variable required for training, and it would result in a big error between the results under the condition with non-optimal variables and the actual equilibrium price. Meanwhile, in real market, the investor always holds the subjective will for each stock, so that the future stock price will change with the will of the people. So the final transaction price will still deviate from the equilibrium price. In 1965, Zadeh proposed the use of membership degree to describe the fuzzy phenomenon with unclear extension in *Fuzzy Sets*. As a result, we quote the fuzzy theory introduced by Zadeh to describe the fuzziness of the estimated stock price in the previous section resulting from non-optimal independent variables and subjective will of investors in order to correct the model.

In this paper, we apply triangular fuzzy number $\tilde{A} = (\alpha, l, u)$ to correct the forecasts, where a is the mid-value when degree of membership equals 1, while $l,u$ are the upper bound and lower bound of which membership degree equals 0, respectively. Obviously, the constructed triangular fuzzy numbers satisfy the following properties:
1. $\alpha(\tilde{A})$ and $\alpha \in (0, 1]$ must be a closed interval (i.e. must be convex).
2. The underlying sets of $\tilde{A}$ must be bounded
3. The membership functions need to satisfy the following equations:

$$\mu(x) = \begin{cases} \dfrac{x}{\alpha - l} - \dfrac{l}{\alpha - l}, & x \in [\,l,\alpha\,]: \\[2mm] \dfrac{x}{\alpha - u} - \dfrac{u}{\alpha - u}, & x \in [\,\alpha,u\,]: \end{cases}$$

$$0, \quad \text{otherwise}:$$

Where $\alpha$ is the mid-value of the fuzzy number, $l$ is the lower bound of the fuzzy number, and $u$ is the upper bound of fuzzy number.

Let $x = \begin{pmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \cdots \\ x_{31} & x_{32} & x_{33} \\ & \vdots & \end{pmatrix}$ be the sample data, where $(x_{1j} x_{2j} x_{3j} \cdots)'$ $(j = 1,2,3 \dots)$ is a sample set of

independent variables for predicting interval segments and $(x_{i1} x_{i2} x_{i3} \cdots)$ $(i = 1,2,3 \dots)$ is a collection of all independent variables in a unit interval of time. By the target function $F_m(x)$ under the target step m which we gained in the above section, we can get the predicted value of the unit time in the forecast interval $F_m\left(\{x_{ij}\}_{j=1}^M\right)$ $(i = 1,2,3 \dots)$. During the unit time i(i=a), $x_a = \{x_{aj}\}_{j=1}^M$, let the predicted value $F_m(x_a)$ be the mid-value of the triangular fuzzy number $\alpha$, let the visible lowest price in a unit of time be the lower bound of the fuzzy number $l$, and let the visible highest price in a unit of time be the upper bound of the fuzzy number $u$.

Therefore, in the unit time i = a, the forecasting value of $F^*(x)$ is changed to:
$$F_m'(x_a) = \tilde{A} = (F_m(x_a), l, u) \quad (1).$$

As it is hard to witness the change of price directly through the triangular fuzzy number, we let the mean value $E(\tilde{A})$ as the final predicted value in order to judge the change of stock price：
1. let F be all fuzzy set defined in the real domain R；
2. let $\tilde{A} \in F$, $[\tilde{A}]_\gamma = [\underline{a}(\gamma), \overline{a}(\gamma)]$ $(\gamma \in [0,1])$ be the level set with the degree of membership $\gamma$, then $\underline{a}(\gamma) = l + \gamma(F_m(x) - l), \overline{a}(\gamma) = u - \gamma(u - F_m(x))$.

The mean of a fuzzy number is denoted as
$$E(\tilde{A}) = \int_0^1 \gamma\left(\underline{a}(\gamma) + \overline{a}(\gamma)\right) d\gamma \quad (2).$$

After getting the mean of the fuzzy number, we set MSE and MAPE as standard to judge the fitting effect and observe the effect after correction. The idiographic expression are
$$MSE = \frac{1}{n}\sum_1^n \left(\hat{Y}_i - Y_i\right)^2,$$

where $\hat{Y}_i$ is the predicted value, $Y_i$ is the actual value. The smaller MSE gets, the better the fitting effect.

$$MAPE = \frac{1}{n}\sum_1^n \frac{|\hat{Y}_i - Y_i|}{Y_i} \times 100\%,$$

where $\hat{Y}_i$ is the predicted value, $Y_i$ is the actual value. The smaller MAPE gets, the better the fitting effect.

## IV.     PREDICTION OF STOCK PRICE AND ITS EMPIRICAL ANALYSIS

This article selects monthly data from 2000 to 2016 of PING AN BANK of China in which the data during 2000~2015 are the learning samples and the data in 2016 will be the prediction interval. For bank share, we choose some factors that strongly related with banking industry, including opening price, CPI index, money supply, value of real estate investment and value of foreign trade. Among them, CPI index includes food consumer price CPI, clothing consumer price CPI, living consumer price CPI, family articles and services consumer price CPI, transportation and communication consumer price index, culture and entertainment consumer price index and medical consumer price index. Money supply covers M1 and M2. Value of real estate investment includes residence, office building and commercial building investment. Value of foreign trade includes volume of imports and exports. One-to-one matching between each independent variable and parameter during prediction process is listed asin the following Table 1.

*Table 1.Parameter corresponding each selected variable*

| Independent variable | Parameter |
|---|---|
| opening price | open |
| food consumer price CPI | food cpi |
| clothing consumer price CPI | cloth cpi |
| living consumer price CPI | live cpi |
| family articles and services consumer price CPI | family cpi |
| transportation and communication consumer price index | traffic cpi |
| culture and entertainment consumer price index | entertainment cpi |
| medical consumer price index | medical cpi |
| M1 money supply | M1 |
| M2 money supply | M2 |
| residence investment | house |
| office building investment | office |
| commercial building investment | commercial |
| volume of imports | import |
| volume of exports | export |

Data source：Data of PING AN BANK comes from Zhongtou securities software， the rest data are from official website of Statistics Department of Republic of China.

We select gbm expansion package form R language to do regression analysis. As matter of the author's experience, learning rate should be set up between 0.01-0.001, while iterations should be 3000-100000. Therefore, the initialized optional loss function has a Gaussian distribution whose iteration is 5000, learning rate is 0.005 and depth of decision tree is 4. After learning based on the function and setting above and confirm the best iteration, there comes out the Ordinary Least Squares–Iterations curve as the following figure 1.
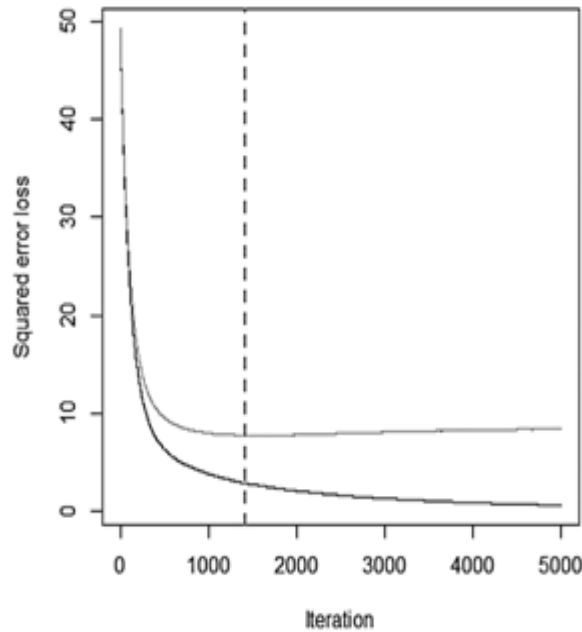
*Figure 1.  Ordinary Least Squares–Iterations curve*

In this curve, we can observe that ordinary loss square comes to its least when iteration was around 2000. And we finally determine that the best one was 1931 by using internal function.

At the same time, we analyse importance of each independent variable by using summary function from which we obtain the bar chart as following figure 2.
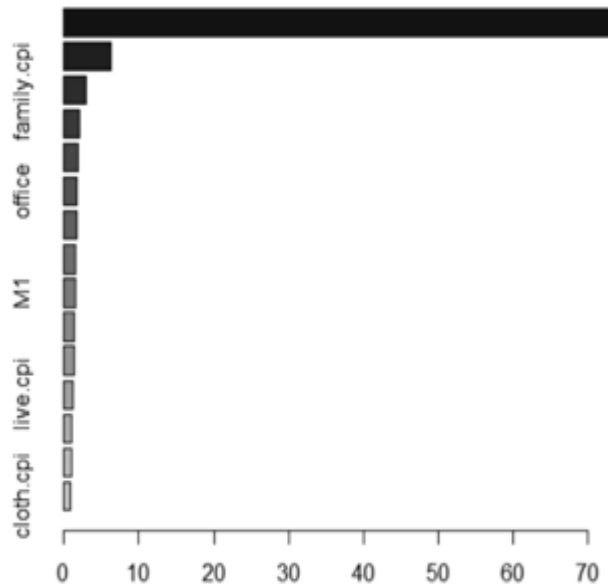


*Figure 2.  Correlation distribution of each independent variable*

*Table 2.  Degree of correlation of variables          Unit：%*

| Independent Variables | Degree of correlation |
|---|---|
| opening price | 74.2865822 |
| culture and entertainment consumer price index | 6.2714176 |
| family articles and services consumer price index | 2.8312574 |
| residence investment | 2.0003960 |
| volume of exports | 1.8144464 |
| commercial building investment | 1.7203278 |
| food consumer price index | 1.5988569 |
| volume of imports | 1.5533530 |
| M1 money supply | 1.5147095 |
| M2 money supply | 1.3591534 |
| commercial building investment | 1.2808006 |
| living consumer price index | 1.1032630 |
| transportation and communication consumer price index | 0.9920670 |
| medical consumer price index | 0.9479887 |
| clothing consumer price index | 0.7253804 |

From Table 2 we can find that the opening and the closing price have strong relativity which up to 74.29%. For other variable, all of them have weak relativity. For instance, in CPI index, culture and entertainment consumer price index and family articles and services consumer price index have stronger relativity while other consumer price index are weak related. For real estate investment, residence investment has stronger relativity than any other factor. Moreover, volume of export has higher correlation than volume of import. M1 and M2 also present a similar result as above.

Utilizing the above trained model, the prediction of testing data in 2016 can produce the outcomes as the following Table 3.

*Table 3.Stock price prediction based on boosting algorithm*

| Date | Actual Value | Prediction Value($\hat{Y}_i$) | $(\hat{Y}_i - Y_i)^2$ |
|---|---|---|---|
| 2016.01 | 10 | 12.482829 | 6.16444 |
| 2016.02 | 9.56 | 10.907009 | 1.814433 |
| 2016.03 | 10.64 | 10.849048 | 0.043701 |
| 2016.04 | 10.57 | 11.563007 | 0.986063 |
| 2016.05 | 10.55 | 10.799841 | 0.062421 |
| 2016.06 | 8.7 | 10.768961 | 4.2806 |
| 2016.07 | 9.2 | 9.414621 | 0.046062 |
| 2016.08 | 9.49 | 9.693153 | 0.041271 |
| 2016.09 | 9.07 | 10.331103 | 1.590381 |
| 2016.10 | 9.15 | 9.592796 | 0.196068 |
| 2016.11 | 9.55 | 9.380959 | 0.028575 |
| 2016.12 | 9.1 | 10.286423 | 1.4076 |
| MSE | 1.388468 | MAPE | 9.54% |

From Table 3 one can easy find that there are also huge deviation led to a bigger MSE, most of prediction values are approach to actual values. Fundamental reason should be that we have chosen imperfect variables related weakly with target variables, leading to unsatisfied prediction result.

Another reason can be found during observing marginal utility of each variable (Please refer attachment for more details.) Taking medical consumer price index (medical cpi) and commercial building investment (Commercial) as example. We chose a value section between 97.8-98.7 for medical cpi, but the chart of variable marginal utility shows that it cause little effect when it is lower than 99.7 which means the change of medical cpi contributed nothing to the prediction process. And Commercial investment was very similar with medical cpi that valued between 1202.73 and 1867.71 which had no effect to prediction, for most data of prediction were set higher than 1300. This kind of problem also exits in other variables, making the change of variable cannot lead to different results.

Therefore we quoted Zadeh's Fuzzy Forecast to analyst, trying to reduce the influence of the problems introduced by imperfect variables. Then we built a triangular fuzzy number $\tilde{A}$ as formula (1) which equipped by the following parameters in Table 4. For the extend part, we select the lowest closing price in a month as the lower bound, the highest opening price in the same month as upper bound. Then we get the prediction values of stock price through boosting algorithm as vertex value whose fuzzy membership degree is 1.

*Table 4 . Triangular fuzzy number of each month's stock price*

| Date | Medium Value($F_m(x)$) | Upper bound($l$) | Lower bound($u$) |
|---|---|---|---|
| 2016.01 | 12.482829 | 7.87 | 9.9 |
| 2016.02 | 10.907009 | 7.72 | 8.46 |
| 2016.03 | 10.849048 | 7.86 | 8.99 |
| 2016.04 | 11.563007 | 8.51 | 9.01 |
| 2016.05 | 10.799841 | 8.35 | 8.83 |
| 2016.06 | 10.768961 | 8.47 | 8.74 |
| 2016.07 | 9.414621 | 8.67 | 9.24 |
| 2016.08 | 9.693153 | 8.93 | 9.8 |
| 2016.09 | 10.331103 | 9.01 | 9.52 |
| 2016.10 | 9.592796 | 9.03 | 9.3 |
| 2016.11 | 9.380959 | 9.01 | 9.78 |

To modify the predicted stock price, we add fuzzy number–mean value function $\hat{Y}_i' = E(\tilde{A})(2)$ into analyst as a final tool. So we can obtain the modified prediction value of stock price as displayed in the following Table 5.

*Table 5 . Improved prediction result of stock price based on enhanced boosting algorithm*

| Date | Actual Value | Modified Value($\hat{Y}_i'$) | $(\hat{Y}_i' - Y_i)^2$ |
|---|---|---|---|
| 2016.01 | 10 | 11.2836 | 1.647629 |
| 2016.02 | 9.56 | 9.968 | 0.166464 |
| 2016.03 | 10.64 | 10.041 | 0.358801 |
| 2016.04 | 10.57 | 10.6287 | 0.003446 |
| 2016.05 | 10.55 | 10.0632 | 0.236974 |
| 2016.06 | 8.7 | 10.0476 | 1.816026 |

| 2016.07 | 9.2 | 9.2614 | 0.00377 |
|---------|----------|--------|----------|
| 2016.08 | 9.49 | 9.5838 | 0.008798 |
| 2016.09 | 9.07 | 9.9757 | 0.820292 |
| 2016.10 | 9.15 | 9.4502 | 0.09012 |
| 2016.11 | 9.55 | 9.3856 | 0.027027 |
| 2016.12 | 9.1 | 9.9893 | 0.790854 |
| MSE | 0.497517 | MAPE | 5.82% |

As is shown in Table 5, values of MSE and MAPE have decreased a lot by the modifying prediction algorithm. For some specific dates, before modifying the value of 2016.01, 2016.02, 2016.06 present a high variance, which decrease rapidly after modifying from 6.16, 1.84, 4.28 to 1.65, 0.17, 1.81. In conclusion, the stock price prediction is performed much betterwith lower deviation and variance by using triangular fuzzy number and the enhanced boosting algorithm.

## V.    CONCLUSION

Compared with single Decision Tree, boosting algorithm not only work better at processing lost data but also not so sensitive to noisy data. On the other hand, boosting algorithm work without shortcomings of single Decision Tree which makes it can fit a complex non-linear relationship by using different return tree as weak returning device, making the accuracy higher. It can avoid overfitting by controlling iteration, so that speed of calculation can be improved. At the same time, compared with bagging algorithm and random forest algorithm, boosting algorithm making weak returning device more balanced because it takes weight of each returning tree into consideration. But it is still a problem that the iteration of boosting algorithm can only calculate by using cross check instead of setting the iteration directly. The problem also took place in practical experiment, in which different cross check will lead to different iteration even based on the same sample. In practical returning, there are still some large deviation due to the overfitting led by some noisy data, the volume of data and so on. But it is also related with variables. Therefore, bringing triangle fuzzy number into analyst can in some way decrease the possibility of choosing independent variables, reducing the value of MSE and MAPE and improving the final result at the same time. But even triangle fuzzy number, there are still some weaknesses. For example, we can't confirm the extend part of triangle fuzzy number accurately, so it is hard to make good use of historical data for prediction. In the future, we will continue doing the prediction research on this issue in order to improve the predicting capability of stock price based on the improved boosting algorithm.

## VI.    ACKNOWLEDGMENT

## REFERENCES

[1] Hoerl AE, 1962, *Application of ridge analysis to regression problems, Chemical Engineering Progress,* 1958, 54–59
[2] Tibshirani, R. (1996). *Regression shrinkage and selection via the lasso. J. Royal. Statist. Soc B.,* Vol. 58, No. 1, pages 267-288.
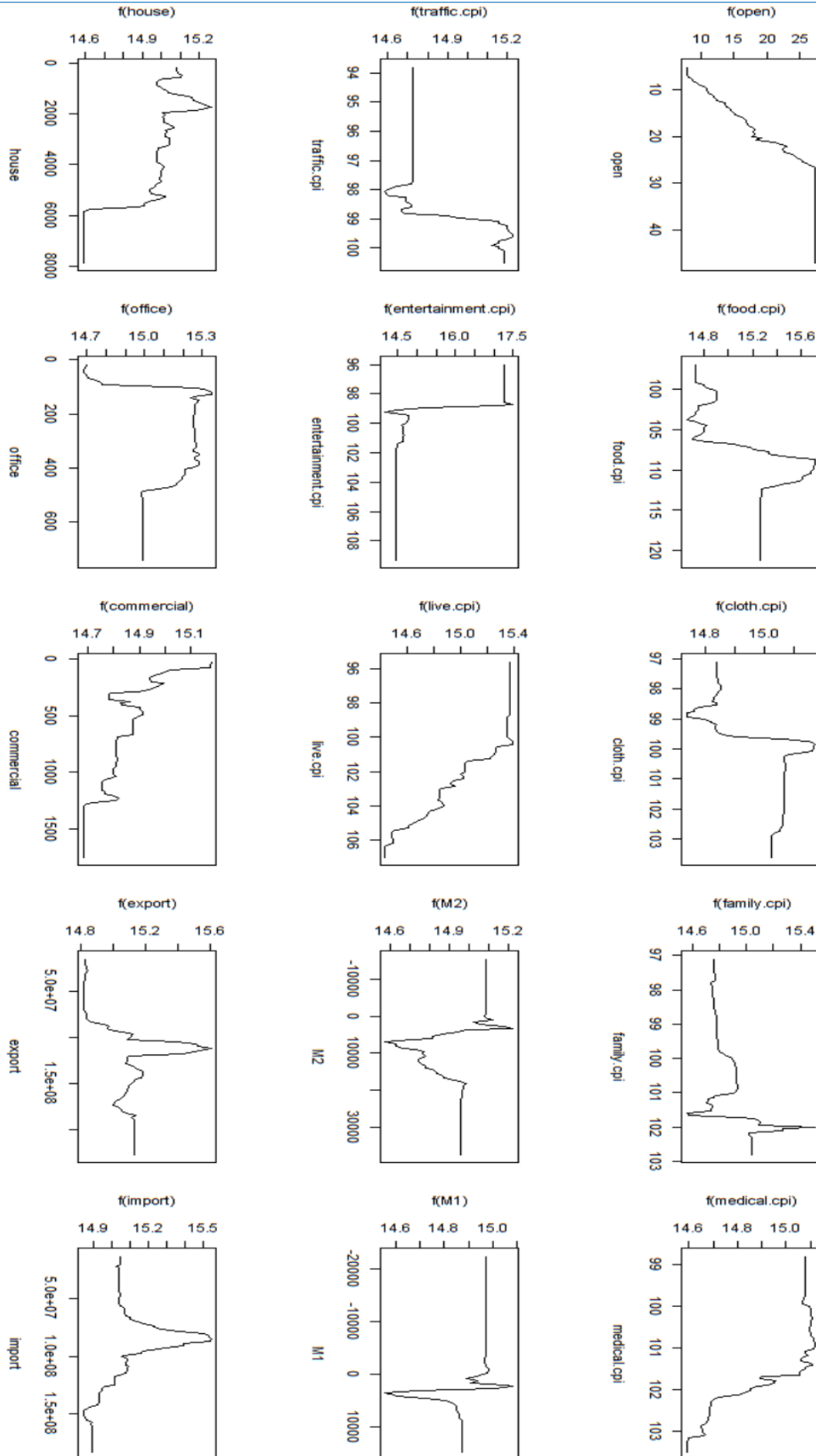[3] Breiman L. *Bagging predictors[J]. Machine Learning,* 1996, 24(2): 123-140.
[4] Friedman J H. *Greedy function approximation: a gradient Boosting machine[J]. The Annals of Statistics,*2001, 29(5): 1189-1232.
[5] Bu̇hlmann P, Yu B. *Boosting with the L2 loss: regression and classification[J]. Journal of the American Statistical Association,* 2003, 98(462): 324-339.

[6] *Bu˙hlmann P. Boosting for high-dimensional linear models[J]. The Annals of Statistics, 2006, 34(2): 559-583*

[7] *Bu˙hlmann P, Hothorn T. Boosting algorithms: regularization, prediction and model fitting[J]. Statistical Science, 2007, 22(4): 477-505.*

[8] *Breiman L. Random forests[J]. Machine Learning, 2001, 45(1): 5-32.*

[9] *Zhang YanZhou, Ma QiuXiang. The Commodity Residential House Price Prediction Based on BP-Boosting. [J]. HENAN SCIENCE, 2014(12):2588-2592.*

[10]*Chen ShiPeng, Jin ShengPing. The House Price Prediction Based on Radom Forest. [J]. Technology Innovation and Application, 2016(4):52-52.*

[11]*Wang Na. Forecasting of Carbon Price Based on Boosting-ARMA Model [J].Statistics & Information Forum, 2017, 32(3).*

[12]*Zadeh L A. Fuzzy sets[J]. Information and Control, 1965, 8(3): 338-353.*

[13]*Zhang W G, Zhang X L, Xiao W L. Portfolio selection under possibilistic mean-variance utility and a SMO algorithm[J]. European Journal of Operational Research, 2009, 197(2): 693-700.*

[14]*Li J, Xu J P. A novel portfolio selection model in a hybrid uncertain environment [J]. Omega, 2009, 37(2): 439-449.*

## Appendix

Marginal Utility of independent variables

**Prediction of Price Trend**



Price Trend