



# Research on sales information prediction system of e-commerce enterprises based on time series model

Jian Liu<sup>1</sup> · Chunlin Liu<sup>2</sup> · Lanping Zhang<sup>1</sup> · Yi Xu<sup>3</sup>

Received: 4 December 2018 / Revised: 19 December 2018 / Accepted: 3 January 2019 /  
Published online: 22 January 2019  
© Springer-Verlag GmbH Germany, part of Springer Nature 2019

## Abstract

Sales forecasting plays an important role in guiding the sales and marketing of e-commerce enterprises, and warehousing department planning warehouse location. At the same time, sales data can better reflect future sales trends. This paper establishes a sales forecasting and analysis model for commodities with common characteristics using their historical sales data through time series model, and forecasts the sales inventory of a certain kind of products from a quantitative point of view. In order to improve the predictive reliability, this paper introduces external observable data and qualitative analysis of historical data prediction model by using hidden Markov model to predict the characteristics of hidden values, so as to further improve the reliability of prediction model.

**Keywords** Information system · Time series · Sales forecasting · Prediction model · Hidden Markoff · Qualitative analysis



✉ liujian@vip.sina.com

Chunlin Liu  
liucl@nju.edu.cn

Lanping Zhang  
8000000315@czie.edu.cn

Yi Xu  
xuyi@tschina.com

- <sup>1</sup> School of Economics and Sports Management, Changzhou Vocational Institute of Engineering, Changzhou 213164, Jiangsu, China
- <sup>2</sup> Nanjing University Business School, Nanjing 210093, Jiangsu, China
- <sup>3</sup> Changzhou Tiansheng New Materials Co, Ltd, Changzhou 213000, Jiangsu, China

## 1 Introduction

In recent years, China has made rapid progress in the field of e-commerce, in which online transactions are the main form. With the improvement of the basic platform construction by Ali Group, Jingdong Mall and other companies, the development of e-commerce in our country has crossed the period of platform construction, and formally entered the stage of rapid development of small and medium-sized e-commerce enterprises relying on platform strength to compete and develop their own unique competitiveness. By introducing personalized products and services, small and medium-sized e-commerce enterprises can more meet the growing demand for online transactions, and achieve the segmentation and service optimization of e-commerce market. However, due to the limitation of technology and resources, small and medium-sized e-commerce enterprises have not paid enough attention to the large amount of transaction data resources they acquire (Arunraj and Ahrens 2015).

For small and medium-sized e-commerce enterprises, sales forecasting plays an important role in guiding sales and marketing of e-commerce enterprises, and warehouse department planning warehouse location. At the same time, sales data can better reflect future sales trends, improve operational efficiency, reduce operating costs, enhance user satisfaction, and ultimately achieve the goal of improving e-commerce enterprises (Michis 2015; Kechyn et al. 2018; Omar et al. 2016). The competitive advantage of industry in competition will increase company profits.

## 2 Quantitative prediction model based on historical data

Main purpose of time series analysis is to predict the future based on the existing historical data (Haviluddin et al. 2015). E-commerce product sales data is a typical time series data. Based on such time series, using the corresponding time series model, we can theoretically predict future sales through fitting and regression of historical data. However, the sequence of sales of different products needs to be treated differently.

The data studied in this paper come from the real data of e-commerce. Through data analysis, the total sales volume of e-commerce enterprises in the initial stage of development is relatively stable on the whole. Therefore, in this paper, the SARMA model is used in the following experiments.

The deterministic time series analysis model can divide the sequence into the following parts:

$$Y = f(T, S, e) \quad (1)$$

Among them, T is trend item, S is seasonal item and E is stochastic disturbance. Since the general trend of commodity sales data does not show linear growth after the stability of the store sales, but shows time-varying characteristics with the seasonal differences of commodity sales, the time-series additive model is adopted in the specific function model:

$$Y_t = \phi_1 Y_{t-1} + \phi_4 Y_{t-4} + \varepsilon_t \quad (2)$$

For the stationary seasonal time series model, in order to eliminate the interference of accidental factors, it is necessary to smooth the sequence. However, both the sliding average smoothing method and the exponential smoothing method are difficult to maintain the original seasonal trend. Therefore, it is necessary to transform the seasonal sequence into the ordinary sequence and then use the sliding smoothing or exponential smoothing method, as follows:

$$X_t = Y_t - \frac{\sum_{i=1}^n Y_{t-s_i}}{N} \tag{3}$$

$$\bar{X}_t = \frac{X_t + X_{t-1} + \dots + X_{t-M+1}}{M} \tag{4}$$

$$\bar{X}_t = \alpha X_t + (1 - \alpha)\alpha X_{t-1} + (1 - \alpha)^2 \alpha X_{t-2} + \dots \tag{5}$$

Finally,  $X_t$  is brought into the ARMA model to calculate, and the final prediction value plus the monthly mean is the prediction value of the monthly time series model.

### 3 Qualitative prediction model based on external single factor

By using time series prediction model and fitting different kinds of commodities, an interpretable prediction value can be given from the perspective of historical data. In fact, such time series forecasting values already contain such seasonal changes and regular promotional activities, so the forecasting results of time series forecasting model are more like black box test, and its forecasting results are somewhat unexplainable (Maciel et al. 2016). Therefore, the forecasting results based on time series forecasting model have certain limitations. Firstly, such forecasting values can not bring in factors with historical differences. For annual sales forecasting, such factors as lower temperature this year than last year and stronger promotion can not respond better in time series model. Come out. Secondly, there is no criterion to judge the predicted value of the model. There is no suitable criterion to judge whether the predicted value should be the maximum or the minimum. Therefore, the problem of historical differences will be solved in the future, and a reference standard for the upper and lower bounds of time series prediction value will be given. Hidden Markov prediction model is introduced here. Some statistical factors will be taken as observation variables, and sales change as implicit variables. Qualitative analysis of prediction results is made by quantitative method.

Before using the hidden Marco model, we first introduce the hypothesis that the Markoff chain must satisfy:

1. The probability distribution of  $t + 1$  system state is only related to the state of  $t$  time, and has nothing to do with the state before  $t$ :

$$P(x_{t+1} | x_t \dots x_1) = P(x_{t+1} | x_t) \tag{6}$$

2. The state transition from  $t$  to  $t + 1$  is independent of the value of  $t$ . The hidden Markov model parameters are as follows:

1.  $S = \{S_1, S_2, \dots, S_N\}$ : state set with  $N$  values.
2.  $V = \{V_1, V_2, \dots, V_M\}$ : observational sets with  $M$  values.
3.  $A = [a_{ij}]$ : state transition matrix.
4.  $B = b_j(k)$ : probability matrix of observations (confusion matrix).

$$b_j(k) = P(o_t = V_k | q_t = S_j) \quad (7)$$

$$j \in [1, N], t \in [1, T]$$

5.  $\pi = \{\pi_i\}$ : Initial probability distribution.

$$\pi = P(q_1 = S_i), i \in [1, N] \quad (8)$$

In this way, a Markov model can be marked as:

$$\lambda = (N, M, A, B, \pi)$$

Among them,  $q_t$  is the state value of  $t$  time,  $o_t$  is the observed value at  $t$  time.

Here we take temperature and volume changes as the two observation sequences. When the temperature of each month increases or decreases with respect to the same period last year as the observation sequence, the observation sequence is {growth, unchanged, decrease}. The proof of the transition probability of the growth change can be obtained by the same statistical method. Such as:

A: relative temperature rises next month; A: relative temperature rises this month.

B: relative temperature is unchanged next month; B: relative temperature remains unchanged this month.

C: relative temperature decreases next month; C: relative temperature decreases this month.

It can be obtained from the full probability formula:

$$(P(A), P(B), P(C)) = (P(a), P(b), P(c)) * A \quad (9)$$

$$A = \begin{bmatrix} P(A|a) & P(A|b) & P(A|c) \\ P(B|a) & P(B|b) & P(B|c) \\ P(C|a) & P(C|b) & P(C|c) \end{bmatrix} \quad (10)$$

Transfer matrix  $A$  can be obtained through statistical meteorological historical data. The relative temperature here is replaced by the average temperature and the mean temperature.

Vector change sequence is {increase, unchanged, decrease}. The relationship between temperature change and sales change can be obtained by statistics of sales change, that is, confusion matrix can also be obtained by statistical method.

Assuming that the sequence of vectors is {x, y, z}, x: growth, y: unchanged, z: decrease, the relationship between historical sales data and temperature change

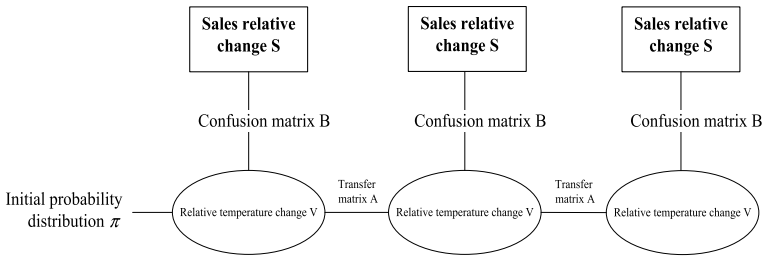


Fig. 1 Single factor hidden Markov model

can be obtained by statistics. The flow chart can be shown in Fig. 1. Confusion matrix is obtained.

$$B = \begin{bmatrix} P(x|A) & P(y|A) & P(z|A) \\ P(x|B) & P(y|B) & P(z|B) \\ P(x|C) & P(y|C) & P(z|C) \end{bmatrix} \tag{11}$$

$$(P(x), P(y), P(z)) = (P(A), P(B), P(C)) * B \tag{12}$$

Through the transfer matrix and confusion matrix, the probability distribution of sales volume change in the next month is finally obtained. According to the probability distribution, the marketing strategy can be adjusted appropriately by clearing the experience rules of salesmen.

To get the forecasting results of the two forecasting models, it is necessary to give the results of adjusting the marketing plan. In view of the single category commodities, the change percentage can be used to carry out marketing early warning. Here, the expert system can be constructed by using the experience of the marketers.

Firstly, the predicted value of time series prediction model is regarded as historical predicted value, which does not include external change factors. It can be concluded that in theory, if all historical conditions do not change, the predicted value will approach the true value (Bretschneider and Gorr 2016). However, in addition to the same historical factors as historical values, other observable and historical factors will also affect sales. At this time, qualitative boundary analysis of time series predictions is carried out by observing the predictive vectors of hidden Markov prediction model (Schneider and Gupta 2016).

#### 4 Qualitative prediction model based on external multiple factors

In the actual sales environment, the impact on sales volume is actually multiple, and the experimental theoretical value of single factor on sales volume is greater than the actual value (Fan et al. 2017). There are many factors related to sales, so the practical application value of qualitative prediction model with multiple external factors is higher than that of single factor prediction model (Wei et al.

2016; Vhatkar and Dias 2016). Therefore, the practical application value of qualitative prediction model with multiple external factors is higher than that of single factor prediction model. However, because there are certain correlations between various external factors, such as the number of colleges and universities and the economic development of the region, these two factors can not be regarded as two independent reasons affecting sales volume (Jiménez et al. 2017). At present, there are relatively few quantitative correlation analysis data for this type of problem, and there is no relatively feasible solution because of the large number of external data types. In order to study the prediction model of multiple external factors, this paper presents the prediction model of external factors under independent conditions.

For example, suppose that there are only two factors affecting sales, and the two factors are independent of each other. The change of either side will not affect the change of the other side, but directly affect sales. Under this premise, the first thing to be given is the weight ratio of the two factors to the sales volume. When the weight ratio is determined, because the two factors are independent of each other, the influence of the two factors on the sales volume data can be calculated independently by using single factor hidden Markov model, and finally the weight ratio pairs occupied by the two factors can be used to calculate the sales volume data independently. The final sales impact is weighted. The flow chart of multiple independent factors hidden Markov prediction model is shown in Fig. 2.

The probability of change can be predicted by using the observed or predicted values of the current weather and the transfer matrix derived from statistics. The formulas are as follows:

$$V_{nm} = V_{n(m-1)} * A_n \tag{13}$$

$V_{nm}$  is the M probability distribution state sequence of N factors,  $A_n$  is the transfer matrix of the N factor. The probability distribution of external factors is transformed

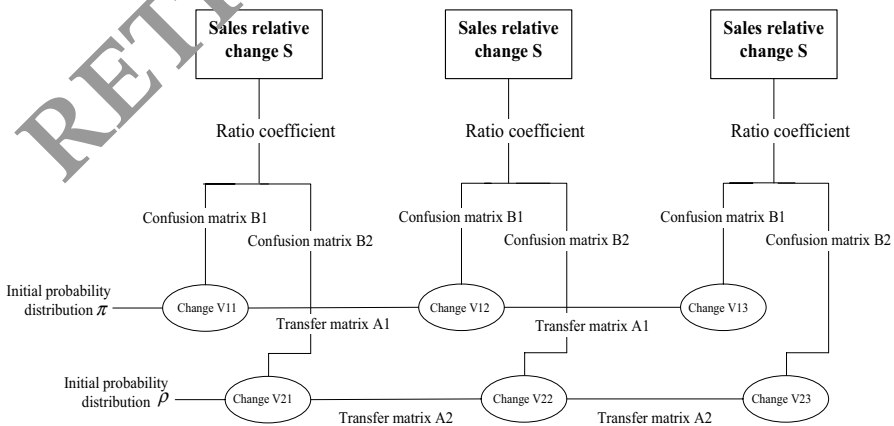


Fig. 2 Hidden Markov prediction model with multiple independent factors

into the transformation probability sequence of sales through confusion matrix. The formulas are as follows:

$$S_n = V_{nm} * B_n \tag{14}$$

Among them,  $B_n$  is the confusion matrix of the  $n$ th factor and  $S_n$  is the probability distribution sequence of the influence on sales under the influence of the  $n$ th factor. When external factors are independent of each other, the weight of external factors on sales is used to obtain the distribution probability of final sales forecast.

$$S = \sum_{i=1}^n \alpha_i * S_i$$

$$\sum_{i=1}^n \alpha_i = 1 \tag{15}$$

$\alpha_i$  B is the weight of I factor.

### 5 Experiment

This experiment firstly gives quantitative sales forecast by time series model, and then makes qualitative analysis by comparing actual sales value and the deviation between forecast value and actual value by using hidden Markov prediction model. The statistics are shown in Table 1.

The other is weather data from Beijing. The data are from China Weather Network. The final collection of tables is summarized in Table 2.

#### 5.1 Quantitative prediction

Using the time series prediction model, the sales data are taken into the model prediction results are shown in Table 3.

The predicted fitting curve is shown in Fig. 3.

Tab. 1 Sales statistics

Year	Month					
	1	2	3	4	5	6
2016	1146	950	827	543	123	12
2017	982	931	754	400	76	24
2018	1096	925	763	310	33	
Year	Month					
	7	8	9	10	11	12
2016	7	9	27	110	531	782
2017	17	3	34	128	476	824

**Table 2** Beijing monthly average monthly high temperature statistics

Date	1	2	3	4	5	6
2015	1	4	11	21	29	30
2016	0	4	12	28	28	28
2017	5	3	16	23	28	31
2018	5	7	14	22	28	30
Date	7	8	9	10	11	12
2015	31	30	26	21	9	-1
2016	32	32	26	19	12	6
2017	33	31	25	19	12	4
2018	31	32	26			

### 5.2 Qualitative prediction

From the results of using historical sales data to predict by time series prediction model, we can see that there are some deviations between the predicted value and the actual value. Some of these deviations are caused by similar random disturbances, and some are caused by external factors. Therefore, the study of bias can help company's better use predictive value.

#### 1. Single factor qualitative analysis

Based on the statistics of the monthly mean high temperature historical data in Beijing (see Table 2), the transfer matrix of temperature state is obtained as follows:

$$A = \begin{bmatrix} 4/13 & 5/13 & 4/13 \\ 3/7 & 1/7 & 3/7 \\ 5/12 & 3/13 & 4/12 \end{bmatrix}$$

Obfuscation matrix:

$$B = \begin{bmatrix} 2/7 & 4/7 & 1/7 \\ 2/5 & 0 & 3/5 \\ 2/5 & 2/5 & 1/5 \end{bmatrix}$$

The initial probability vectors (1, 0, 0) are constructed by using the temperature reduction in December 2017, and the temperature variation vectors in January and

**Table 3** Forecast results

Model	Jan-18	Feb-18	Mar-18
Forecast	1050	932	788
UCL	1116	1003	865
LCL	984	862	711



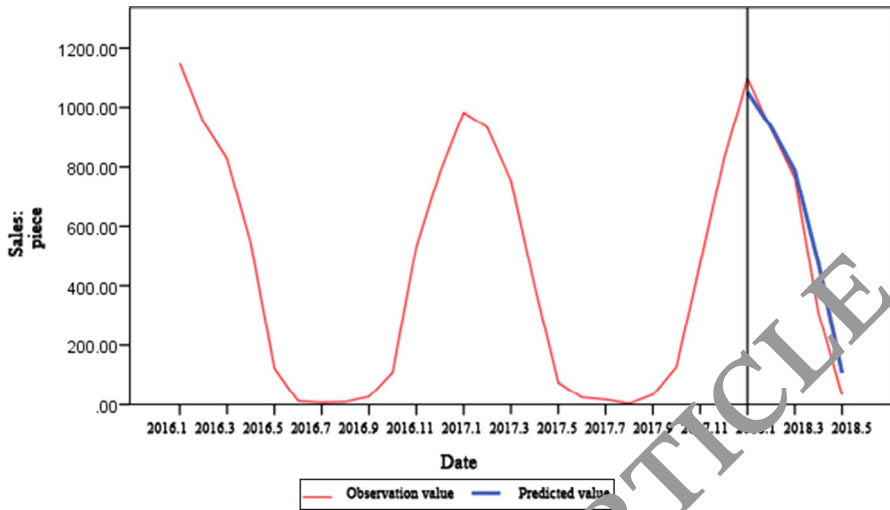


Fig. 3 Sales forecast of seasonal time series model

February 2018 are predicted by using the transfer matrix as follows: (0.31, 0.38, 0.31), (0.39, 0.2, 0.41). From the temperature vectors, we can see that under the condition of temperature decreasing in December last year, the probability of temperature decreasing in January is 0.31, the probability of invariance is 0.38, and the probability of increase is 0.31. By multiplying the temperature vector with the confusion matrix, the sales change vector is obtained as follows: (0.36, 0.3, 0.34), (0.36, 0.39, 0.25). The significance of this sequence is: on the premise that the temperature in December was higher than that in the previous year, the probability of decreasing, unchanged and increasing sales of this kind of commodity in January was (0.36, 0.3, 0.34), and the probability of decreasing, unchanged and increasing sales of this kind of commodity in February was (0.36, 0.39, 0.25). According to the calculation, the risk of sales decline in February is relatively small. We should keep appropriate inventory and prepare for possible sales growth while maintaining the expected sales in February.

## 2. Multi factor qualitative analysis

From the sales curve of the experimental commodity data, it can be seen that the sales change has a strong correlation with the temperature change. Therefore, the temperature data are used as the external observation data to conduct qualitative analysis of the sales forecast value. However, there are many sales data of commodities are not sensitive to seasonal temperature. With the promotion and sale of merchants, the sales volume of commodities shows a steady growth pattern. This paper takes the backpack sales of a brand as an example, and its sales curve is shown in Fig. 4.

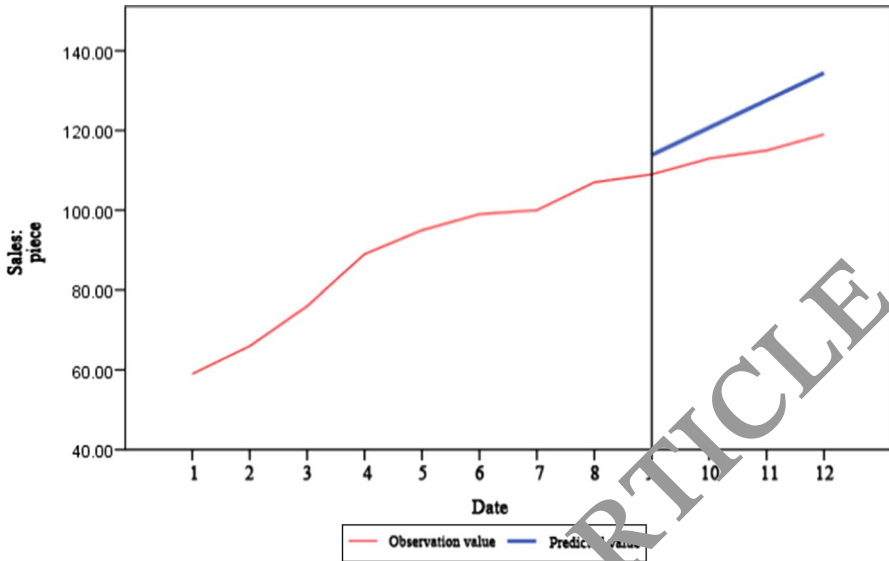


Fig. 4 Sales forecast for non seasonal time series

Table 4 CPI, weather and sales volume table

	Jan-17	Feb-17	Mar-17	Apr-17	May-17	Jun-17
CPI	103.3	101.8	102.1	101.5	102.1	102.2
Weather	5	6	16	23	28	31
Sales	59	60	76	89	95	99
	Jul-17	Aug-17	Sep-17	Oct-17	Nov-17	Dec-17
CPI	101.5	101.5	101.1	100.5	100.7	100.8
Weather	33	31	25	19	12	4
Sales	100	104	110	123	103	117

For multi-factor qualitative analysis, this experiment uses Beijing average high temperature historical data and Beijing CPI data as shown in Table 4. As two opposing factors, the correlation coefficients between Beijing average monthly maximum temperature in 2017 and Beijing CPI index in 2017 are calculated:  $cor(CPI, weather) = 0.0035$ .

Correlation coefficient between CPI and weather in sales volume is:

$$cor(CPI, sales) = -0.8091481$$

$$cor(weather, sales) = 0.3900728$$

Through correlation calculation, we can see that although the correlation between the two external factors is very weak, it has different degrees of correlation effect on sales. Temperature transfer matrix  $A_1$ :

**Table 5** Prediction probability distribution based on transfer matrix

	September observation	October forecast	November forecast	December forecast
Monthly average high temperature change probability	(1, 0, 0)	(0.31, 0.38, 0.31)	(0.38, 0.24, 0.36)	(0.38, 0.28, 0.34)
CPI change probability	(1, 0, 0)	(0.45, 0.55, 0)	(0.75, 0.25, 0)	(0.58, 0.42, 0)

**Table 6** Relationship between temperature change and sales volume change

Unit: times	Temperature reduction	Constant temperature	Temperature rise
Sales reduction	15	5	7
Sales volume is unchanged	5	6	4
Increase in sales	8	7	23

$$A_1 = \begin{bmatrix} 4/13 & 5/13 & 4/13 \\ 3/7 & 1/7 & 2/7 \\ 5/12 & 3/13 & 4/12 \end{bmatrix}$$

Because the domestic CPI data is always growing, the CPI is processed differently, and then the transfer matrix is obtained according to statistics  $A_2$ :

$$A_2 = \begin{bmatrix} 8/18 & 10/18 & 0 \\ 1 & 0 & 0 \\ 5/13 & 8/13 & 0 \end{bmatrix}$$

Through the time series forecasting model, the sales volume in the 3 months after simulated forecast is shown in Fig. 4.

In September 2017, the monthly mean maximum temperature in Beijing was lower than that in September 2016, so the initial probability was (1, 0, 0). In September 2017, the CPI index was lower than that in August. Therefore, the initial probability of CPI was (1, 0, 0). The predicted results by confusion transfer matrix were shown in Table 5.

Confusion matrix of the impact of temperature change on sales volume is shown through historical data (as shown in Table 6):

$$B_1 = \begin{bmatrix} 15/27 & 5/27 & 7/27 \\ 5/15 & 6/15 & 4/15 \\ 8/38 & 7/38 & 23/38 \end{bmatrix}$$

As a single factor, the probability of temperature change on sales change is (Table 7):

**Table 7** Temperatures-HMM forecast result

	October forecast	November forecast	December forecast
Monthly average high temperature change probability	(0.36, 0.27, 0.37)	(0.38, 0.24, 0.38)	(0.38, 0.25, 0.37)

**Table 8** relationship between CPI and sales volume changes

	CPI growth reduction	CPI growth remains unchanged	CPI growth is bigger
Sales decrease	14	9	4
Sales remain unchanged	4	7	3
Sales increase	5	12	22

**Table 9** CPI-HMM forecast results

	October forecast	November forecast	December forecast
Probability distribution of CPI impact on sales volume	(0.39, 0.43, 0.18)	(0.43, 0.38, 0.16)	(0.42, 0.4, 0.18)

The confusion matrix of the impact of CPI changes on sales volume (as shown in Table 8) is:

$$B_2 = \begin{bmatrix} 14/27 & 9/27 & 4/27 \\ 4/14 & 7/14 & 3/14 \\ 5/39 & 12/39 & 22/39 \end{bmatrix}$$

As a single factor, the impact of CPI change on sales volume is (Table 9):

Through correlation analysis, we can see that there are significant differences between the two factors from the perspective of relevance. Therefore, the weight of the impact on the final sales is also different. According to the results of correlation analysis, the empirical weight distribution is given. The proportion of CPI influence is 0.8, and the proportion of temperature change is 0.2. The results of weighted calculation are shown in Table 10.

From the probability distribution of sales volume change after weighting, we can see that the probability of sales volume decreasing or keeping unchanged in 10, 11 and 12 months is significantly greater than the probability of sales volume rising. Although there are some discrepancies between the qualitative prediction results and the actual results in October, from the perspective of all 3 months, the qualitative prediction results are in line with the qualitative prediction results.

**Table 10** Multivariate qualitative analysis results

	October forecast	November forecast	December forecast
CPI sales change probability distribution	(0.39, 0.43, 0.18)	(0.46, 0.38, 0.16)	(0.42, 0.4, 0.18)
Temperature distribution change probability distribution	(0.36, 0.27, 0.37)	(0.38, 0.24, 0.38)	(0.38, 0.25, 0.37)
Weighted posterior probability distribution	(0.384, 0.398, 0.218)	(0.444, 0.352, 0.204)	(0.412, 0.370, 0.218)

## 6 Summary

In this paper, we first use time series model to predict historical sales data. Experiments show that the time series model has a good effect on sales forecasting, but considering that the time series model can not introduce external variables, the hidden Markov model is used to introduce external variables into the forecasting model. In the part of hidden Markov prediction, this paper analyses the possible influence of external factors on the prediction value from two perspectives: single factor and multi-factor. Finally, the practicability of the hidden Markov model in qualitative prediction is verified through experiments.

**Acknowledgements** Supported by the National Natural Science Foundation of China (Grant No: 71572075).

## References

- Arunraj NS, Ahrens D (2015) A hybrid seasonal autoregressive integrated moving average and quantile regression for daily food sales forecasting. *Int J Prod Econ* 170:321–335
- Bretschneider S, Gornow W (2016) Economic, organizational, and political influences on biases in forecasting state sales tax receipts. *Int J Forecast* 7(4):457–466
- Fan ZP, Che YJ, Chen ZY (2017) Product sales forecasting using online reviews and historical sales data: a method combining the Bass model and sentiment analysis. *J Bus Res* 74:90–100
- Haviluddin H, Alfred R, Obit JH et al (2015) A performance comparison of statistical and machine learning techniques in learning time series data. *Adv Sci Lett* 21(10):3037–3041
- Jiménez F, Sánchez G, García JM et al (2017) Multi-objective evolutionary feature selection for online sales forecasting. *Neurocomputing* 234(C):75–92
- Kechyn G, Yu L, Zang Y et al (2018) Sales forecasting using WaveNet within the framework of the Kaggle competition. *Int J Forecast* 6(2):332
- Maciel L, Ballini R, Gomide F (2016) Evolving granular analytics for interval time series forecasting. *Granul Comput* 1(4):1–12
- Michis AA (2015) A wavelet smoothing method to improve conditional sales forecasting. *J Oper Res Soc* 66(5):832–844
- Omar H, Hoang VH, Liu DR (2016) A hybrid neural network model for sales forecasting based on ARIMA and search popularity of article titles. *Comput Intell Neurosci* 2016(4):9656453
- Schneider MJ, Gupta S (2016) Forecasting sales of new and existing products using consumer reviews: a random projections approach. *Int J Forecast* 32(2):243–256
- Vhatkar S, Dias J (2016) Oral-care goods sales forecasting using artificial neural network model. *Procedia Comput Sci* 79:238–243

Wei J, Zhu J, Huang C et al (2016) A novel prediction model for sales forecasting based on grey system. In: IEEE international conference on service-oriented computing and applications. IEEE, pp 23–24

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

RETRACTED ARTICLE