# Analysis of agricultural exports based on deep learning and text mining

**Jia-Lang Xu[1] · Ying-Lin Hsu[2]**

## Abstract

Agricultural exports are an important source of economic profit for many countries. Accurate predictions of a country's agricultural exports month on month are key to understanding a country's domestic use and export figures and facilitate advance planning of export, import, and domestic use figures and the resulting necessary adjustments of production and marketing. This study proposes a novel method for predicting the rise and fall of agricultural exports, called agricultural exports time series-long short-term memory (AETS-LSTM). The method applies Jieba word segmentation and Word2Vec to train word vectors and uses TF-IDF and word cloud to learn news-related keywords and finally obtain keyword vectors. This research explores whether the purchasing managers' index (PMI) of each industry can effectively use the AETS-LSTM model to predict the rise and fall of agricultural exports. Research results show that the inclusion of keyword vectors in the PMI values of the finance and insurance industries has a relative impact on the prediction of the rise and fall of agricultural exports, which can improve the prediction accuracy for the rise and fall of agricultural exports by 82.61%. The proposed method achieves improved prediction ability for the chemical/biological/medical, transportation equipment, wholesale, finance and insurance, food and textiles, basic materials, education/professional, science/technical, information/communications/broadcasting, transportation and storage, retail, and electrical and machinery equipment categories, while its performance for the electrical and optical categories shows improved prediction by combining keyword vectors, and its accuracy for the accommodation and food service, and construction and real estate industries remained unchanged. Therefore, the proposed method offers improved prediction capacity for agricultural exports month on month, allowing agribusiness operators and policy makers to evaluate and adjust domestic and foreign production and sales.

**Keywords** Long short-term memory (LSTM) · Purchasing managers' index · Principal component analysis · Non-manufacturing purchasing managers' index · Time series · Text mining

---

Extended author information available on the last page of the article

## 1 Introduction

Artificial intelligence-related technologies, resources, and infrastructure have gradually matured and can now be easily applied to various fields to solve multiple problems with good effect, including time series, image processing, audio signal processing, and natural language processing. While time series are widely used in various fields, long- and short-term memory models are preferred for deep learning analysis of time series. For example, Qin et al. [23] used time series to create a model used to detect abnormal behavior in controller area network (CAN) buses under tampering attacks. Tulensalo et al. [32] used local weather data to determine total local grid transmission losses. Shahid et al. [26] used time-series methods to predict the number of COVID-19 deaths and recovery cases in ten major countries. Tahvili et al. [30] proposed a natural language processing and data conversion approach that uses supervised learning methods to process and evaluate unbalanced data. Zhang et al. [40] applied text mining and natural language processing techniques to construction accident reports and used support vector machines (SVM), linear regression (LR), decision tree (DT), and other models plus an ensemble model to classify the causes of accidents.

Taiwan's early economic development was mainly based on agriculture, but with the transition to an industrial and technology-based economy, the importance of the agricultural sector gradually diminished. However, in order to guarantee food security and resource sustainability, the Taiwan government has begun to focus additional attention on the development of agriculture-related industries. What little arable land Taiwan has is not well-suited to large-scale agricultural operations, and farmers focus rather on cultivating agricultural products with high added value. This has led to steadily increasing exports, a trend the government hopes to encourage, allowing export figures to reflect actual income, in order to reduce the imbalance of domestic production and sales. The development of information technology allows for easy access to a wide range of information, and online news sources provide fast and convenient insight into what is happening at home or abroad. Wei et al. [33] suggested that purchasing manager indexes (PMI) can be effectively used to predict the prices of industrial stocks. Xu and Hsu [35] obtained good results using news related to climate change and oil prices to analyze and predict agricultural product prices. Chen and Gong [4] assessed the impact of global warming on the total factor productivity of agriculture, finding that global warming has an impact on agricultural products. Liu et al. (2020a, b) used multiple PMI as auxiliary variables to predict coal mining accidents. Su et al. [27] suggested there is a positive two-way causal relationship between the prices of agricultural products and oil. Sun and Li [28] suggested that the global financial crisis and common borders had a significant effect on China's trade profits on agricultural exports to ASEAN countries. This research investigates whether Taiwan's agricultural exports are impacted by international news on climate change, oil prices and other related matters, and changes in PMI for various industries. This research proposes an AETS-LSTM deep learning model, adjusting the characteristics and weight parameters of the learning target column, to successfully forecast future agricultural export trends.

The remainder of this research is arranged as follows: Sect. 2 reviews the literature on LSTM, text mining, and principal component analysis. Section 3 introduces the research architecture and methods used. Section 4 describes results and performance evaluation, and Sect. 5 presents conclusions.

## 2 Literature review

### 2.1 Long short-term memory

The core of long short-term memory (LSTM) consists of three control gates: input, forget, and output. The input gate uses the input value and the value in the newly generated memory cell in an activation function to determine whether the value must be added to the long-term memory neuron. The forget gate determines whether the current value is a new topic or data that is the opposite of the current value and determines whether the value needs to be filtered or kept in the memory. The output gate determines whether the current value needs to be added to the output. The activation function of the output valve is usually determined using the Sigmoid method. Finally, the activation function tanh is used to determine whether long-term memory should be added to the output. The value falls between [-1, 1], with -1 ordering the removal of long-term memory, while 1 means it should be retained. Following Hochreiter and Schmidhuber [8], Fig. 1 shows the LSTM architecture, followed by Eqs. (1) to (6).

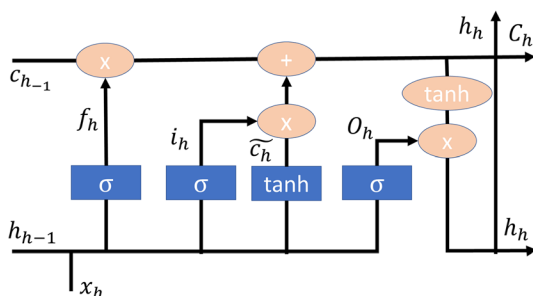$$f_h = \sigma(w_f \cdot \left[h_{h-1}, x_h\right] + b_f) \tag{1}$$

$$i_h = \sigma\left(w_i \cdot \left[h_{h-1}, x_h\right] + b_i\right) \tag{2}$$

$$\widetilde{c_h} = \tanh(w_c \cdot \left[h_{h-1}, , x_h\right] + b_c) \tag{3}$$

$$C_h = f_h * c_{h-1} + i_h * \widetilde{C_h} \tag{4}$$

$$O_h = \sigma\left(w_O\left[h_{h-1}, x_h\right] + b_0\right) \tag{5}$$

**Fig. 1** LSTM Architecture

$$h_h = O_h * \tanh(C_h) \tag{6}$$

where $f_h$ is the forget value, $i_h$ is the input value, $O_h$ is the output value, $\tilde{c}_h$ is the memory cell candidate, $h_{h-1}$ is the current output value, and $x_h$ is the input value. $w_i, w_c, w_o, w_f$ and $b_i, b_c, b_o, b_f$ are, respectively, the weight matrix and deviation vector. $C_h$ is a storage unit, and σ is the Sigmoid activation function.

Chen [2] proposed a nonlinear LSTM algorithm with good prediction accuracy for use in voltage prediction. Elsheikh et al. [6] proposed an LSTM model to predict the fresh water production of stepped solar distillers and conventional distillers. Kırbaş et al. [15] proposed using LSTM for predicting COVID-19 cases in Denmark, Belgium, Germany, France, the UK, Finland, Switzerland, and Turkey using performance indicators such as MSE, PNSR, RMSE, NRMSE, MAPE, and SMAPE to evaluate model accuracy, but with final results indicating that LSTM provided the best accuracy. Liu et al. (2020a, b proposed a model combining deep neural network (DNN) and long short-term memory (LSTM) to solve the problem of developing a system model based on given input and output data to predict sinter chemical composition. Miao et al. [20] proposed an LSTM framework for fog forecasting using hourly meteorological elements. They believed that the LSTM framework is more effective than traditional machine learning models in this application. Rahman et al. [24] proposed a novel diabetes classification model based on Conv-LSTM and used a grid search algorithm to perform hyperparameter optimization so that the applied model can find the best parameters. Tsantekidis et al. [31] proposed a combination of the ability of CNN to extract useful features and the ability of LSTM to analyze time series for evaluation. They demonstrated that their proposed model outperformed various compared LSTM and CNN models within the prediction range of the test. Yan et al. [37] proposed an ON-LSTM model to determine the remaining service life of gears and compared its performance with those of LSTM, GRU, DLSTM and DNN. The proposed ON-LSTM model achieved the best short-term and long-term prediction accuracy of the compared models.

## 2.2 Text mining

Very few hard and fast rules exist for written text, and this means that text mining techniques are unable to identify definable texts in either long or short passages, despite their clear comprehensibility in daily life. TF-IDF is a common statistical calculation method that can be divided into two parts: TF (term frequency) and IDF (inverse document frequency). TF is a calculation term, where the frequency of occurrence is described by Eq. (7), while IDF is a measure to calculate a word's general importance, as in Eq. (8). TF-IDF is used to filter common words, while retaining important words, as in Eq. (9).

$$TF_{i,j} = \frac{n_{ij}}{\Sigma_k n_{k,j}} \tag{7}$$

where $n_{ij}$ is the number of times a specific word or phrase appears in the news content, and $\Sigma_k n_{k,j}$ is the sum of all words or phrases in the news content.

$$IDF_i = \log\left(\frac{D}{d_i}\right) \tag{8}$$

where D is all news content, and $d_i$ is the number of words that appear in all news content. If a word is not contained in the news content, the denominator is 0, thus $t_i+1$ is generally used.

$$TFIDF_{i,j} = TF_{i,j} \times IDF_i \tag{9}$$

Therefore, TF-IDF is used to calculate the frequency of specific words or phrases in specific news content.

Word2Vec is a method of converting words into vectors which represent their meaning, because it is difficult to determine the relationships between words and phrases, as with synonyms, antonyms and corresponding words. This raises the importance of Word2Vec. Mikolov et al. [21] used large amounts of text data to represent the semantic meaning between words and phrases by means of their corresponding vectors. After embedding words in a space, words with similar meanings will have greater spatial proximity. The most common Word2Vec models are CBOW and Skip-gram. Skip-gram uses a given input word to predict the context, while CBOW uses given a context to predict the input word. Figures 2 and 3 show the CBOW and Skip-gram model architectures, respectively.

Jain et al. (2021a, b, c, d, ) a proposed a Cuckoo Search-eXtreme gradient boosting model and optimized the model to recommend airlines. They also (2021 b) proposed a sparse self-attentive network-based aspect-aware model that can effectively predict consumer recommendation decisions. In addition, they (2021 c) proposed a multi-label classification model for travel recommendation, and finally, they (2021 d) explored the applicability of consumer sentiment analysis in online reviews of machine learning models and explored the literature review. Quamer et al. [22] proposed a self-attentive convolutional neural network model that can effectively perform sentence matching and natural language inference. Yen et al. [36] used the text in online news and stock forums to conduct text exploration and predict future financial performance. Choi et al. [3] proposed using text mining to analyze social network texts to identify cyber bullying. Jung and Lee [13] used text mining of keywords and citation information in academic papers to examine the information value of research scopes and trends. Mosa [18] proposed the mining of large-scale social media data for recombination into a multi-objective optimization (MOO) task for abstract extraction using the gravity search algorithm (GSA) to optimize several expression targets to generate brief social media summaries. To explore resident sentiment toward an urban waste classification policy, Wu et al. [34] used text mining methods to collect and analyze public comments on Weibo. Zhong et al. [39] proposed a four-step modeling model: (1) using an implied Dirichlet distribution to identify dangerous topics, then (2) a convolutional neural network (CNN) algorithm for the automatic classification of such hazards, followed by (3) word co-occurrence networks (WCN) which determine the relationships between the hazards, and finally (4) a quantitative analysis of keywords through word cloud methods to create a visual overview
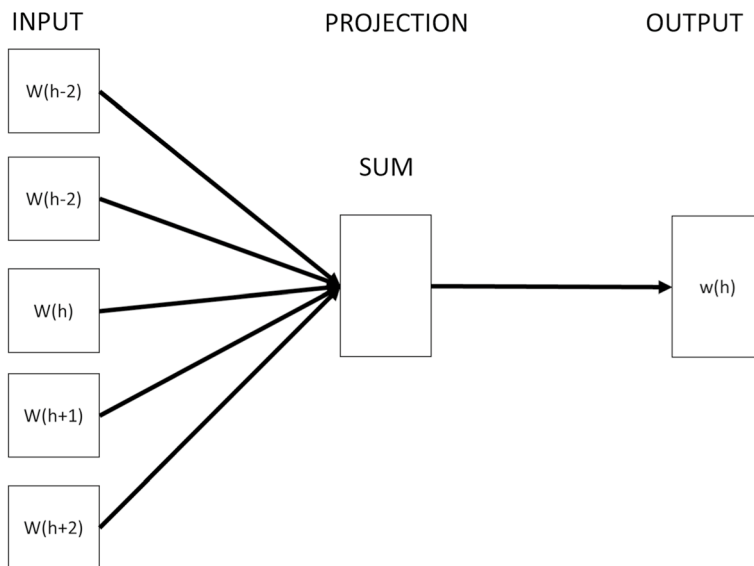
INPUT PROJECTION OUTPUT

W(h-2)

W(h-2)

SUM

W(h)

W(h+1)

W(h+2)

w(h)

**Fig. 2** CBOW Model

INPUT PROJECTION OUTPUT
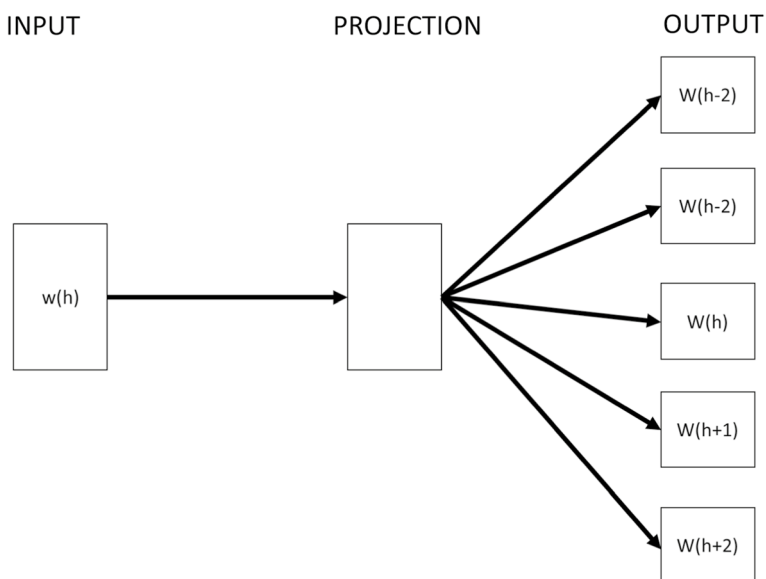
W(h-2)

W(h-2)

w(h)

W(h)

W(h+1)

W(h+2)

**Fig. 3** Skip-gram Model

of such hazards, thus providing managers with new knowledge and insights. Drury and Roche [5] applied text mining methods to a large number of papers and news reports cited in recent agricultural research papers, seeking to identify problems and potential applications.

## 2.3 Principal component analysis

Principal component analysis (PCA) is a widely used unsupervised learning linear transformation technique that allows for original data to be converted into different modes of expression and also allows for data processing. PCA converts high-dimensional data into lower-dimensional data, thus reducing calculation time and memory space requirements to facilitate storage and analysis. For example, if the data have $q$ internally correlated continuous variables, $x_1, x_2, \ldots, x_q$, there must be independent variables, using Eq. (1) for linear transformation, while in Eq. (11) the vector $v_i$ is the eigenvector of array $A$, and $\lambda_i$ is the eigenvalue corresponding to the eigenvector $v_i, i = 1, 2, \ldots, q$, and in Eq. (12), S explains the cumulative proportion of variance as the main component:

$$
\begin{aligned}
y_1 &= v_{11}x_1 + v_{12}x_2 + \ldots + v_{1q}x_q \\
y_2 &= v_{21}x_1 + v_{22}x_2 + \ldots + v_{2q}x_q \\
&\vdots \\
y_q &= v_{q1}x_1 + v_{q2}x_2 + \ldots + v_{qq}x_q
\end{aligned}
\tag{10}
$$

$$
\lambda v = A v
\tag{11}
$$

$$
S = \frac{\lambda_1 + \lambda_2 + \ldots + \lambda_r}{\lambda_1 + \lambda_2 + \lambda_3 + \ldots + \lambda_q} * 100\%
\tag{12}
$$

where the selection of r is 90% above S.

Boubchir and Aourag [1] proposed a new multivariate technique with PCA and PLS to effectively test the statistical influence of the stability of inverse perovskites and perovskites. Mahmoudi et al. al. (2020) used PCA to classify the spread of COVID-19 in France, Germany, Iran, Italy, Spain, the UK and the USA. Raj et al. [25] applied PCA for the simplified, economical and sensitive classification of vesicles and bronchus. García-Gil et al. [7] proposed a new method using the open-source cluster computing framework Apache Spark platform and main component analysis to reduce data volume, finding that high-dimensional data will affect the algorithm's calculation time. Jolliffe and Cadima [14] suggested that PCA can be used to reduce data dimensionality, thus increasing data interpretability and minimizing information loss. New data obtained are learned from the data set, making PCA an adaptive data analysis technique.

## 3 Research design and process

The experimental environment of this research includes two levels of hardware and software. The operating system used is Windows 10 64 bits, running on an i7 CPU, with 24G RAM and a GeForce GTX1650Ti graphics card. Python 3.6 is used for development.

This research uses web-crawlers to obtain international news, the MI of each industry, and the export index of agricultural products. The collected data are pre-processed, merged, and then used to train the LSTM model. The overall process is shown in Fig. 4.

### 3.1 Collecting International News 2014—2019

This research uses Web crawlers to collect all news reports on the international news section of the ETtoday Web site, using HTML tags to filter all international news content using search terms including agriculture, petroleum, climate and other related topics in the period January 1, 2014, to December 31, 2019.

### 3.2 Using Jieba for international news segmentation

To clarify and smooth the overall calculation and structure, the international news content data collected for this research were processed using Jieba word segmentation to facilitate subsequent data annotation.

### 3.3 Using TF-IDF find keywords

Following Jieba segmentation, the TF-IDF method is used to identify the top 10 key words for each international news article.

### 3.4 Using Word2Vec convert words into word vectors

Following Jieba segmentation, the Word2Vec method is used to vectorize words, allowing for the calculation of word similarity. The tested feature vector dimensions are 100, 200, and 300, and the tested and interpretable word vectors are all consistent. Therefore, the Word2Vec feature vector is set at 100 dimensions, and CBOW is used for training.

### 3.5 Word cloud analysis by date

Keyword terms were categorized by month, and then, the ten most critical words were determined by means of word cloud analysis, as shown in Fig. 5.

### 3.6 Deriving word vectors from word cloud

Identified key words were then processed using the trained Word2Vec model to obtain the word vectors of ten key words.
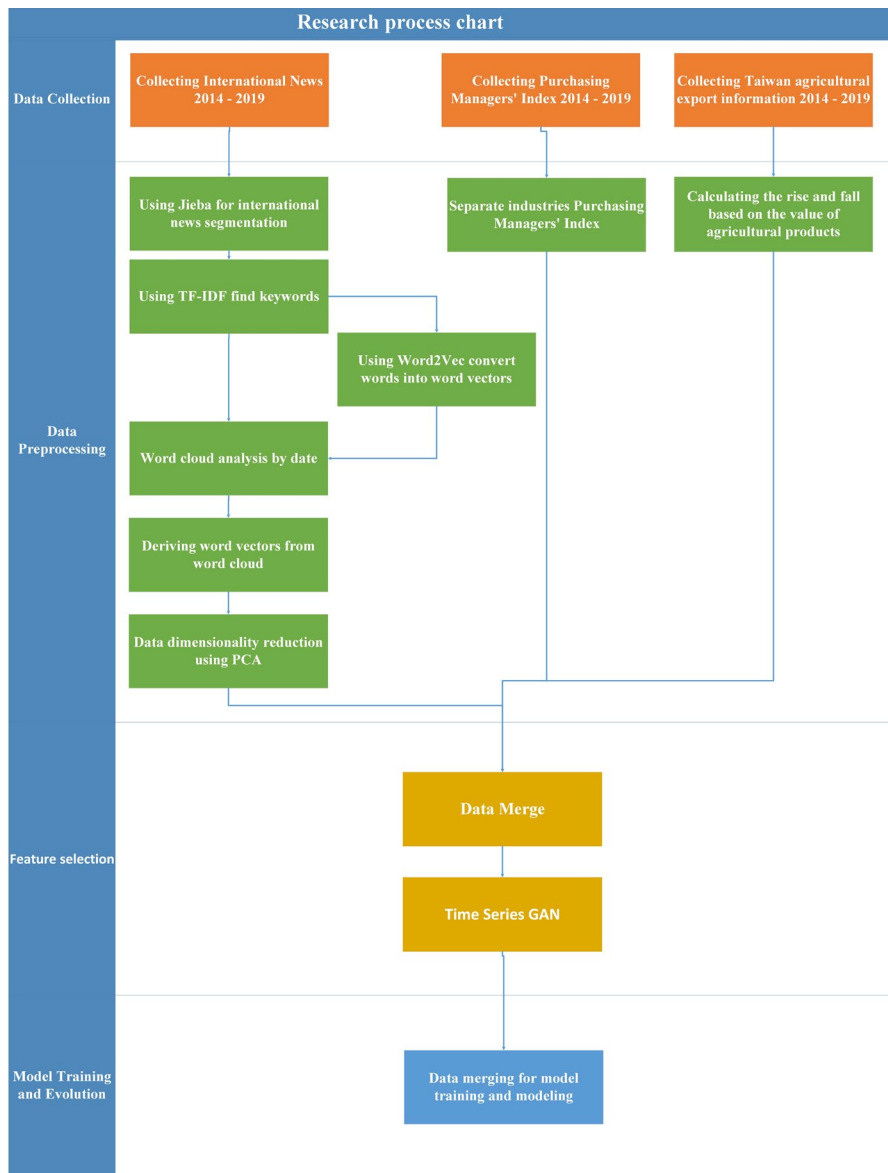
**Fig. 4** Research process chart

**Fig. 5** Word Cloud

### 3.7 Data dimensionality reduction using PCA

Because of the excessive dimensionality of the generated word vectors, PCA is used to reduce to dimensions to 55, with an interpretability rate of 90%. Therefore, in this research, word dimensions should be reduced from 1000 to 55.

### 3.8 Collecting purchasing managers' index 2014—2019

The manufacturing and non-manufacturing purchasing managers' indexes for 2014 to 2019 are obtained from the Business Indicators Database Web site.

### 3.9 Separate industries purchasing managers' index

Manufacturing and non-manufacturing indexes were obtained for various industries including chemical/biological/medical, transportation equipment, accommodation and food service, wholesale, finance and insurance, food and textiles, basic materials, education/professional, scientific/technical, information/communications/broadcasting, transportation and storage, retail, electrical and machinery equipment, electronic and optical industry, and construction and real estate.

#### 3.9.1 Collecting Taiwan agricultural export information 2014—2019

Public data for the total export value of agricultural products from 2014 to 2019 were obtained from Taiwan's Council of Agriculture.

#### 3.9.2 Calculating the rise and fall based on the value of agricultural products

This research obtained the fluctuation data of the total export value by subtracting the difference between the present and previous months.

#### 3.9.3 Data Merging

This research merged the parameters for the various data including agricultural export fluctuations, the PMI of each industry, the outlook for each industry in the coming six months, and the 55-dimensional word PCA.

### 3.9.4 Time-series GAN

The data from 2014 to 2017 are insufficient for effective training. Following Yoon et al. [38], time-series GAN was used to generate real samples through various real and synthetic time-series data, thereby generating sufficient data for training.

### 3.9.5 Data merging for model training and modeling

The input parameters of the proposed AETS-LSTM model are as follows: fluctuations in agricultural exports, 55-dimensional words after PCA, the PMI of each industry, and the outlook for each industry for the next six months. The output parameter is fluctuations in agricultural exports over the following month. The training sample used is data from 2014 to 2017 following application of the time-series GAN algorithm. The test sample is data from 2018 to 2019. Srivastava et al. [29] noted that Dropout is independent of each neuron and each iteration in the hidden layer and can improve the problem of overfitting. Setting position and size will also be a key factor. For example, too high a value will cause neurons to be completely covered and the model will not learn the training characteristics. Too low a value may lead to model overfitting in training, leading to Dropout. Zeng et al. [41] noted that L1 or L2 regularized sparse representation methods have been applied in different fields, where a representation based on L1 regularization is sparser, while a representation based on L2 is simpler and faster. Table 1 shows the parameter settings used for training and modeling in this research.

## 4 Experimental results and performance evaluation

First, this research uses the proposed AETS-LSTM model to predict the rise and fall of agricultural exports in the PMI of each industry. Figure 6 shows that the chemical/biological/medical, accommodation and food service, financial and insurance, basic materials, education/professional, science/technical, information/communications/broadcasting, transportation and storage, and retail

**Table 1** Model parameters

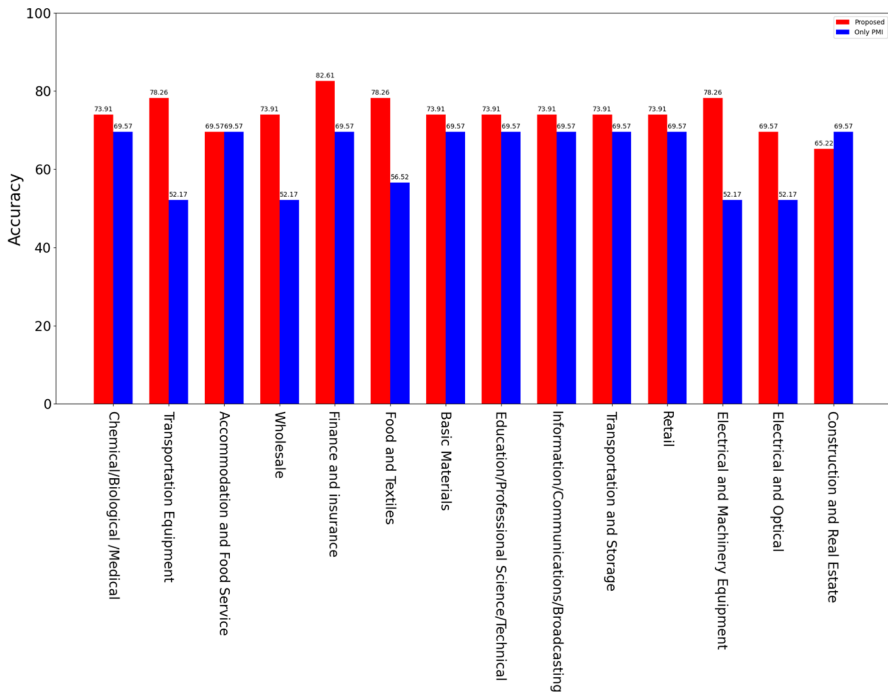| Parameter Name | Parameter Value |
|---|---|
| LSTM Layer | 2 Layers |
| Output Layer | 1 Layer |
| Activation Function | LSTM Layer: tanh |
| | Output Layer: softmax |
| Parameter Settings | Loss Function: categorical_crossentropy |
| | Optimizer: Adam |
| Dropout | 0.2 |
| Regularizer | 0.0001 |

**Fig. 6** Comparison of the prediction effect of combining keyword vectors and only PMI

industries produce better predictions of the rise and fall of agricultural product exports than those obtained using other industries. Next, this research combines the PMI of each industry and the keyword vectors into the proposed of AETS-LSTM model to compare and evaluate the rise and fall of agricultural exports. The indicators for predicting the rise and fall of exports use precision, recall, f-score, sensitivity, specificity, and accuracy, and the result analysis is shown in Table 2 and Fig. 6. With the exception of the accommodation and food service and the construction and real estate industries, all other industries performed well using AETS-LSTM models. Among them, the finance and insurance industries can show that more than 80% of the results of all performance evaluations, while their prediction accuracy improved from 69.57% of the original PMI to 82.61%, which is better than other industries.

Finally, this research compares the prediction results obtained using the proposed AETS-LSTM model and the neural network and SVM model for the four industries with the best prediction results, namely the finance and insurance, transportation equipment, food and textile, and electrical and machinery equipment industries. Figure 7 shows that the proposed AETS-LSTM model achieves excellent prediction results for the rise and fall of agricultural exports.

**Table 2** Comparison of the predictive performance indicators of various industries

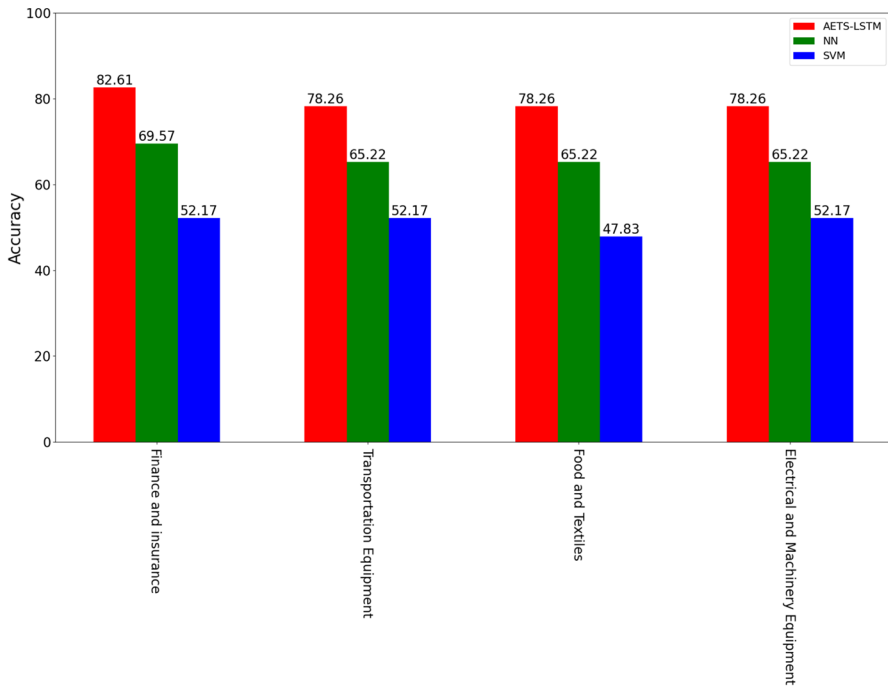| Industry | Precision (%) | Recall (%) | f-score (%) | Sensitivity (%) | Specificity (%) | Accuracy (%) |
|---|---|---|---|---|---|---|
| Chemical/Biological/Medical | 77.78 | 63.64 | 70.00 | 63.64 | 83.33 | 73.91 |
| Transportation Equipment | 87.50 | 63.64 | 73.68 | 63.64 | 91.67 | 78.26 |
| Accommodation and Food Service | 70.00 | 63.64 | 66.67 | 63.64 | 75.00 | 69.57 |
| Wholesale | 77.78 | 63.64 | 70.00 | 63.64 | 83.33 | 73.91 |
| Finance and insurance | 81.82 | 81.82 | 81.82 | 81.82 | 83.33 | 82.61 |
| Food and Textiles | 80.00 | 72.73 | 76.19 | 72.73 | 83.33 | 78.26 |
| Basic Materials | 72.73 | 72.73 | 72.73 | 72.73 | 75.00 | 73.91 |
| Education/Professional Science/Technical | 77.78 | 63.64 | 70.00 | 63.64 | 83.33 | 73.91 |
| Information/Communications/Broadcasting | 77.78 | 63.64 | 70.00 | 63.64 | 83.33 | 73.91 |
| Transportation and Storage | 77.78 | 63.64 | 70.00 | 63.64 | 83.33 | 73.91 |
| Retail | 77.78 | 63.64 | 70.00 | 63.64 | 83.33 | 73.91 |
| Electrical and Machinery Equipment | 87.50 | 63.64 | 73.68 | 63.64 | 91.67 | 78.26 |
| Electrical and Optical | 70.00 | 63.64 | 66.67 | 63.64 | 75.00 | 69.57 |
| Construction and Real Estate | 66.67 | 54.55 | 60.00 | 54.55 | 75.00 | 65.22 |

**Fig. 7** Comparison of algorithm accuracy for the top four industries

## 5 Conclusion

This research proposes a new AETS-LSTM model for effectively predicting agricultural export trends. Agricultural export trends are affected by many factors, such as news content, climate change, and the PMIs of various industries. Few studies have examined the impact of news content and PMI trends on agriculture. Experimental results show that PMI values for the finance and insurance industries combined with keyword vectors have a relative impact the prediction of the rise and fall of agricultural exports and can improve the prediction accuracy for the rise and fall of agricultural exports by 82.61%. The prediction accuracy for chemical/biological/medical, transportation equipment, wholesale, finance and insurance, food and textiles, basic materials, education/professional science/technical, information/communications/broadcasting, transportation and storage, retail, electrical and machinery equipment, and the electrical and optical industries can be improved by combining keyword vectors, while the prediction accuracy for the accommodation and food service and construction and real estate industries remained unchanged. Therefore, this research can enhance the understanding of agribusiness operators and policy makers with regard to the rise and fall of agricultural exports month on month and allow them to better evaluate and adjust domestic and foreign production and sales.

Taiwan's agricultural product export data are presented monthly, restricting the number of data points available, and relevant privacy laws restrict access to data

beyond publicly available information, such as public news reports and Open Data resources. These data availability restrictions place limitations on feature prediction accuracy. The current study focuses exclusively on agricultural exports, but future work could expand on the current results by applying the model structure to other types of exports.

**Declaration**

**Conflict of interest** The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

# References

1. Boubchir M, Aourag H (2020) Materials genome project: The application of principal component analysis to the formability of perovskites and inverse perovskites. Computational Condensed Matter 24:e00495
2. Chen Y (2020) Voltages prediction algorithm based on LSTM recurrent neural network. Optik 220:164869
3. Choi YJ, Jeon BJ, Kim HW (2020) Identification of key cyberbullies: A text mining and social network analysis approach. Telemat Inform 56:101504
4. Chen S, Gong B (2020) Response and adaptation of agriculture to climate change: Evidence from China. J Develop Econom 148:102557
5. Drury B, Roche M (2019) A survey of the applications of text mining for agriculture. Comput Electron Agri 163:104864
6. Elsheikh AH, Katekar VP, Muskens OL, Deshmukh SS, Elaziz MA, Dabour SM (2020) Utilization of LSTM neural network for water production forecasting of a stepped solar still with a corrugated absorber plate. Process Saf Environ Prot 148:273–282
7. García-Gil D, Ramírez-Gallego S, García S, Herrera F (2018) Principal components analysis random discretization ensemble for big data. Knowl-Based Syst 150:166–174
8. Hochreiter S, Schmidhuber J (1997) Long short-term memory. Neural Comput 9(8):1735–1780
9. Jain PK, Yekun EA, Pamula R, Srivastava G (2021) Consumer recommendation prediction in online reviews using Cuckoo optimized machine learning models. Comput Electr Eng 95:107397
10. Jain PK, Quamer W, Pamula R, Saravanan V (2021) SpSAN: sparse self-attentive network-based aspect-aware model for sentiment analysis. J Ambient Intell Humanized Comp. https://doi.org/10.1007/s12652-021-03436-x
11. Jain PK, Pamula R, Yekun EA (2021) A multi-label ensemble predicting model to service recommendation from social media contents. J Supercomput. https://doi.org/10.1007/s11227-021-04087-7
12. Jain PK, Pamula R, Srivastava G (2021) A systematic literature review on machine learning applications for consumer sentiment analysis using online reviews. Comp Sci Rev 41:100413
13. Jung H, Lee BG (2020) Research trends in text mining: Semantic network and main path analysis of selected journals. Expert Syst Appl 162:113851
14. Jolliffe IT, Cadima J (2016) Principal component analysis: a review and recent developments. Phil Trans Royal Soc A: Math, Phys Eng Sci 374(2065):20150202
15. Kırbaş I, Sözen A, Tuncer AD, Kazancıoğlu FS (2020) Comparative analysis and forecasting of COVID-19 cases in various European countries with ARIMA, NARNN and LSTM approaches. Chaos Solitons Fractals 138:110015
16. Liu Q, Liu J, Gao J, Wang J, Han J (2020) An empirical study of early warning model on the number of coal mine accidents in China. Safety Sci 123:104559
17. Liu S, Liu X, Lyu Q, Li F (2020) Comprehensive system based on a DNN and LSTM for predicting sinter composition. Appl Soft Comp 95:106574
18. Mosa MA (2020) A novel hybrid particle swarm optimization and gravitational search algorithm for multi-objective optimization of text mining. Appl Soft Comp 90:106189
19. Mahmoudi MR, Heydari MH, Qasem SN, Mosavi A, Band SS (2020) Principal component analysis to study the relations between the spread rates of COVID-19 in high risks countries. Alexandria Eng J. 60:467

20. Miao KC, Han TT, Yao YQ, Lu H, Chen P, Wang B, Zhang J (2020) Application of LSTM for short term fog forecasting based on meteorological elements. Neurocomputing 408:285–291
21. Mikolov T, Chen K, Corrado G, Dean J (2013) Efficient estimation of word representations in vector space arXiv preprintarXiv:1301.3781
22. Quamer W, Jain PK, Rai A, Saravanan V, Pamula R, Kumar C (2021) SACNN: self-attentive convolutional neural network model for natural language inference. ACM Trans. Asian Low-Resour Lang Inf Process 20(3):1–16
23. Qin H, Yan M, Ji H (2020) Application of controller area network (CAN) bus anomaly detection based on time series prediction. Veh Commun 27:100291
24. Rahman M, Islam D, Mukti RJ, Saha I (2020) A deep learning approach based on convolutional LSTM for detecting diabetes. Comput Biol Chem 88:107329
25. Raj V, Renjini A, Swapna MS, Sreejyothi S, Sankararaman S (2020) Nonlinear time series and principal component analyses: Potential diagnostic tools for COVID-19 auscultation. Chaos, Solitons Fractals 140:110246
26. Shahid, F., Zameer, A., & Muneeb, M. (2020). Predictions for COVID-19 with deep learning models of LSTM, GRU and Bi-LSTM. *Chaos Solitons Fractals*, *140*, 110212.
27. Su CW, Wang XQ, Tao R, Oana-Ramona L (2019) Do oil prices drive agricultural commodity prices? Further evidence in a global bio-energy context. Energy 172:691–701
28. Sun Z, Li X (2018) The trade margins of Chinese agricultural exports to ASEAN and their determinants. J Integr Agric 17:2356–2367
29. Srivastava N, Hinton G, Krizhevsky A, Sutskever I, Salakhutdinov R (2014) Dropout: a simple way to prevent neural networks from overfitting. J Mach Learning Res 15(1):1929–1958
30. Tahvili S, Hatvani L, Ramentol E, Pimentel R, Afzal W, Herrera F (2020) A novel methodology to classify test cases using natural language processing and imbalanced learning. Eng Appl Arti Intell 95:103878
31. Tsantekidis A, Passalis N, Tefas A, Kanniainen J, Gabbouj M, Iosifidis A (2020) Using deep learning for price prediction by exploiting stationary limit order book features. Appl Soft Comput 93:106401
32. Tulensalo J, Seppänen J, Ilin A (2020) An LSTM model for power grid loss prediction. Electric Power Syst Res 189:106823
33. Wei Y, Bai L, Yang K, Wei G (2021) Are industry-level indicators more helpful to forecast industrial stock volatility? Evidence from Chinese manufacturing purchasing managers index. J Forecast 40(1):17–39
34. Wu Z, Zhang Y, Chen Q, Wang H (2020) Attitude of Chinese public towards municipal solid waste sorting policy: a text mining study. Sci Total Environ 756:142674
35. Xu JL, Hsu YL (2021) The Impact of News Sentiment Indicators on Agricultural Product Prices. Comput Econ https://doi.org/10.1007/s10614-021-10189-4
36. Yen MF, Huang YP, Yu LC, Chen YL (2021) A Two-Dimensional Sentiment Analysis of Online Public Opinion and Future Financial Performance of Publicly Listed Companies. Comp Econ https://doi.org/10.1007/s10614-021-10111-y
37. Yan H, Qin Y, Xiang S, Wang Y, Chen H (2020) Long-term gear life prediction based on ordered neurons LSTM neural networks. Measurement 165:108205
38. Yoon J, Jarrett D, van der Schaar M (2019) Time-series generative adversarial networks. In: Advances in neural information processing systems, pp 5508–5518
39. Zhong B, Pan X, Love PED, Sun J, Tao C (2020) Hazard analysis: A deep learning and text mining framework for accident prevention. Adv Eng Inform 46:101152
40. Zhang F, Fleyeh H, Wang X, Lu M (2019) Construction site accident analysis using text mining and natural language processing techniques. Autom Constr 99:238–248
41. Zeng S, Gou J, Deng L (2017) An antinoise sparse representation method for robust face recognition via joint l1 and l2 regularization. Expert Syst Appl 82:1–9

## Authors and Affiliations

**Jia-Lang Xu[1] · Ying-Lin Hsu[2]**

✉  Ying-Lin Hsu
    ylhsu@nchu.edu.tw

    Jia-Lang Xu
    jlxu.academy@gmail.com

[1]   Doctoral Program in Data Science and Industrial Analytics, National Chung Hsing University,
      Taichung City 402, Taiwan

[2]   Department of Applied Mathematics and Institute of Statistics, National Chung Hsing
      University, Taichung City 402, Taiwan