

Overview of the NTCIR-14 FinNum Task: Fine-Grained Numeral Understanding in Financial Social Media Data

Chung-Chi Chen¹[0000-0003-3680-9277], Hen-Hsen Huang^{2,4}[0000-0001-9169-3081],
Hiroya Takamura³[0000-0002-3244-8294], and Hsin-Hsi
Chen^{1,4}[0000-0001-9757-9423]

¹ Department of Computer Science and Information Engineering, National Taiwan University, Taiwan

² Department of Computer Science, National Chengchi University, Taiwan

³ Artificial Intelligence Research Center, National Institute of Advanced Industrial Science and Technology, Japan

⁴ Most Joint Research Center for AI Technology and All Vista Healthcare, Taiwan
cjchen@nlg.csie.ntu.edu.tw, hhuang@nccu.edu.tw,
takamura.hiroya@aist.go.jp, hhchen@ntu.edu.tw

Abstract. Numeral is the crucial part of financial documents. In order to understand the detail of opinions in financial documents, we should not only analyze the text, but also need to assay the numeric information in depth. Because of the informal writing style, analyzing social media data is more challenging than analyzing news and official documents. In this paper, we give an overview of the results of a shared task called FinNum in NTCIR-14 for fine-grained numeral understanding in financial social media data, i.e., to identify the category of a given numeral in a tweet. This task attracted 13 participants to register, received 16 submissions from 9 participants, and, finally, accepted 6 papers from participants.

Keywords: numeral understanding · financial social media · numeral corpus.

1 Introduction and Motivation

When analyzing a financial instrument, investors always focus on two sides, fundamental and technical. Investors using fundamental analysis attempt to evaluate the intrinsic value of the financial instrument. For the security of company, they may focus on the numerals in financial statements. For the treasury bond, they may evaluate the price depending on US Fed Funds Target Rate. Those who use technical analysis may employ the technical indicator like moving average (MA), relative strength index (RSI), and so on. No matter which analysis method investors use, numeral plays an important role, and provides much pivotal information in financial data.

Numeral contains much important information in financial domain. For example, investors may use price-earnings ratio (P/E ratio) to evaluate the value

2 CC. Chen et al.

of security of certain company, where both P/E ratio and the value of security are numeric information. For the purpose of understanding the fine-grained numeric information in social media data, we adopt the taxonomy for numerals [4], and classify numerals into 7 categories and further extend several categories into subcategories. Especially, the most important category, Monetary, is extended into 8 subcategories. (T1) is an instance that contains several numerals, and the categories of the numerals are dissimilar. For example, 8 is the numeral about quantity, 17.99 is about stop loss price, 200 is the parameter of technical indicator, and 1 is the price of stock. In such a short sentence, there are 4 kinds of numerals. That shows the importance of numerals in financial narrative.

(T1) 8 breakouts: \$CHMT (stop: \$17.99), \$FLO (200-day MA), \$OMX (gap), \$SIRO (gap). One sub-\$1 stock. Modest selection on attempted swing low

In the development of the FinNum shared task, we adopt the fine-grained numeral taxonomy [4], and provide more instances for the meaning of each (sub)category in Section 2. The task settings and the details of corpus, called FinNum 2.0, are described in Section 3. The methods and the experimental results of participants are shown and discussed in Sections 4 and 5. Finally, we conclude the remarks and give some future directions in Section 6.

2 Taxonomy

To analyze the opinion of an investor, understanding the implications of numerals is one of the important challenges. Numerals in financial tweets are classified into 7 categories based on experts' experiences. Four of them are further extended to several subcategories. The most important category for financial tweets is Monetary, which is divided into 8 subcategories [4]. The details are elaborated in the subsequent sections.

2.1 Monetary

All numerals about monetary will be annotated as Monetary category. For example, 110.20 in (T2) quoting the price of Facebook's security is one of the cases of Monetary category. We classify the cases in Monetary categories into the following 8 subcategories: "money", "quote", "change", "buy price", "sell price", "forecast", "stop loss" and "support or resistance".

(T2) \$FB (110.20) is starting to show some relative strength and signs of potential B/O on the daily.

The ideas to distinguish these subcategories are: (1) "money", "quote" and "change" narrate the status, but not contain the opinion; (2) the other subcategories present the opinions of a tweet poster. The numerals about money like

loss \$3 billion will be classified into “money” subcategory. 110.20 in (T2) is an example for “quote”.

The numeral describing the change of the price or money will be seen as “change”. For example, \$AAPL -\$3 today is the description of the change of the price of \$AAPL.

To capture the “buying” and “selling” prices of an individual investor can help us understand the performance of the writer. Based on the performance information, we can give different weights for the opinions of each investor. 137.89 in (T3) is an instance for “buying” subcategory. 36.50 in (T4) is an example for “selling” aspect.

(T3) \$SPY Long 1/2 position 137.89

(T4) \$KOG Took a small position- hopefully a better outcome than getting kneecapped by \$BEXP selling itself dirt cheap at 36.50

Some investors may “forecast” the price of the instruments depending on their analysis results. The numeral about the prediction of monetary will be classified into “forecast” subcategory. 14.35 in (T5) is an example for “forecast” subcategory. This kind of opinions can be considered as the summarization of the analysis results, which provide not only the market sentiment and the sentiment degree information, but also the exactly price level [5]. “stop loss” price is the price level that investors may close their positions. 17.99 in (T1) is one of the examples.

(T5) \$CIEN, CIEN seems to have broken out of a major horizontal resistance. Targets \$14.35.

Subcategory “support or resistance” price captures the prediction of price movement. Some investors think that when the price reaches the resistance price, it will fall down. On the other hand, they think that when the price reaches the support price, it will rebound. This subcategory could help us indicate the boundary of price movement. 46 in (T6) is an instance for “support or resistance” subcategory.

(T6) \$CTRP, \$46 Breakout Should be Confirmed with Wm%R Stochastic Up

2.2 Percentage

There are many numerals about ratio in financial documents. For example, there are a lot of accounting ratios like P/E ratio, current ratio, and so on. All numerals about percentage information will be classified into Percentage category, and further extended into two subcategories, “absolute” and “relative”. Subcategory “absolute” indicates the proportion of a certain amount, and “relative” is about the change relative to original amount. 167.1 in (T7) is an example of “absolute”, and 1.64, -2.7, -2.5, and -1.6 are the examples of “relative”.

(T7) no trades today...currently 167.1% net long..ended the day down 1.64% due to \$CASY (-2.7%), \$NKE (-2.5%), \$SRCL (-1.6%) and \$JJSF (-1.6%)

4 CC. Chen et al.

2.3 Option

Option is widely discussed in financial social media. There are two numerals in Option category, “maturity date” and “exercise price”. Capturing both information can also help us evaluate the performance of investors as “target” price in Monetary category. (T8) shows the case for “maturity date”, and (T9) shows the case for “exercise price” (\$111).

(T8) looks like a big feb 18-22 \$put spread on \$cree. (T9) Bought \$FB \$111 calls for \$0.62.

2.4 Indicator

Some investors use technical indicator to analyze the price movement. In order to match the analysis result with price, we need to get the parameter they used. (T10) is an example, which shows the necessary of identifying the parameter of technical indicator.

(T10) \$AAPL hit my short term target of the 100 SMA.

2.5 Temporal

Temporal information is also important in financial domain. The day most investors focusing on is the one with high volatility. For example, the day releasing earning information or the day announcing economics data. Thus, to capture the temporal information could help us capture the important date and time that many investors focus on. Here, the numerals in Temporal category are separated into two subcategories, “date” and “time”. (T11) is an example for “date”, and (T12) is for “time”.

(T11) @DrCooper: \$GDX \$NUGT \$DUST Buying on Weakness (06/30/2015)
(T12) \$AMRN So what was that @ 11 a.m.?

2.6 Quantity

Quantity information can help us know the position of an investor, and we can give the large weighting to the opinions held by persons who have large positions. Furthermore, the amount of sales is also the important information in accounting, which has also been classified into this category. (T13) is an example for Quantity category.

(T13) \$RSOL bought 3500 shares today!

2.7 Product/Version

The opinion toward iPhone 4 and iPhone 8 may have different impacts toward Apple’s security. Thus, to capture the Product/Version Number is one of the important tasks in understanding the topic discussed. (T14) is an example of this category.

(T14) iPhone 6 may not be as secure as Apple thought.. \$AAPL

Table 1. Statistics of FinNum 2.0

Category	Subcategory	Train	Dev.	Test	Total	Ratio
Monetary		2467	261	459	3187	35.94%
	money	589	52	95	736	8.30%
	quote	792	89	152	1033	11.65%
	change	143	8	25	176	1.98%
	buy price	319	36	60	415	4.68%
	sell price	103	10	22	135	1.52%
	forecast	270	33	52	355	4.00%
	stop loss	25	4	6	35	0.39%
	support or resistance	226	29	47	302	3.41%
Percentage		838	105	170	1113	12.55%
	relative	585	70	112	767	8.65%
	absolute	253	35	58	346	3.90%
Option		169	11	22	202	2.28%
	exercise price	113	5	14	132	1.49%
	maturity date	56	6	8	70	0.79%
Indicator		167	22	27	216	2.44%
Temporal		2364	253	401	3018	34.03%
	date	2079	223	351	2653	29.92%
	time	285	30	50	365	4.12%
Quantity		741	87	154	982	11.07%
Product/Version		114	14	22	150	1.69%
		6860	753	1255	8868	100.00%

3 Task Setting and Data

3.1 Task Formulation

In the FinNum task, the position of a numeral in a tweet is given in advance. Participants are asked to disambiguate its category. This task is further separated into two subtasks defined as follows.

- **Subtask 1:** Classify a numeral into 7 categories, i.e., Monetary, Percentage, Option, Indicator, Temporal, Quantity and Product/Version Number.
- **Subtask 2:** Extend the classification task to the subcategory level, and classify numerals into 17 classes, including Indicator, Quantity, Product/Version Number, and all subcategories.

Micro-averaged F-score and macro-averaged F-scores are adopted for evaluating the classification performance of participants’ runs.

3.2 Corpus Creation

We collected the data from StockTwits⁵, one of the social trading platform for investors to share their ideas and strategies. Two experts were involved in

⁵ <https://stocktwits.com/>

6 CC. Chen et al.

the annotating process. The dataset, FinNum 2.0⁶, used in this shared task, contains only the numerals in full agreement. There are 4,072, 457, and 753 tweets in training set, development set, and test set, respectively. Note that, there are at least one cashtag, e.g., \$AAPL is a cashtag stands for the stock of Apple Inc., and at least one numeral in each tweet. There are total 8,868 annotated numerals in FinNum 2.0. The statistics of FinNum 2.0 is shown in Table 1. The annotations are licensed under the Creative Commons Attribution-Non-Commercial-ShareAlike 4.0 International (CC BY-NC-SA 4.0) license, and provided for academic usage.

4 Methods in the Official Runs

4.1 Preprocessing and Features

Some common preprocessing techniques, including replacing URL, hashtag, cashtag by special tokens; and converting characters into lowercase, are used by participants.

Several Features are extracted from tweets:

- **Part-of-speech (POS) Tags:** Ait Azzi and Bouamor [1] and Liang and Su [7] extracted POS features with CMU ARK Twitter POS Tagger [10] and CoreNLP [8], respectively.
- **Keywords:** Ait Azzi and Bouamor [1] adopted the keywords from Chen et al. [4]. Liang and Su [7] proposed patterns for some (sub)categories.
- **Topic:** Spark [13] used Latent Dirichlet Allocation (LDA) [2] to extract the features for tweets’ topics.
- **Position:** The position of the target numeral in the tweet is considered [13].
- **Named Entity:** Named entity extracted by CoreNLP [8] is used [7].
- **Term Frequency:** Term Frequency-Inverse Document Frequency (TF-IDF) is adopted for the context information [15].
- **Format:** Integer (float) format information is also adopted as a feature [13, 16]. Several co-occurrence format information is extracted by the given patterns [16].
- **Numeral Information:** Spark [13] not only used the raw value of the numeral, but also the log of raw value and the normalized raw value.
- **Bag-of-Characters:** The near n characters of the target numeral are considered [13].
- **Prefixes/Suffixes:** Prefixes and Suffixes are used in Wu et al. [16]
- **Brown Cluster:** The j -character prefix of the Brown cluster [3] is considered as features [16].
- **Recognizers-Text Type:** The text type extracted by Microsoft Recognizers is apoted [16].

⁶ <http://nlg.csie.ntu.edu.tw/nlpresource/FinNum/>

4.2 Representations

Most submissions used word embeddings and character embeddings to represent tweets. Skip-gram [9], Glove [11], ELMo [12], and BERT [6] are used for representing token information. Ait Azzi and Bouamor [1] used one-layer convolutional neural network (CNN) to create character embeddings. Tian and Peng [14] initialized the character embedding with *glove.840B.300d-char.txt* and used long short-term memory (LSTM)-based model to fine-tune the character embedding with char2id scheme. Some submissions concatenated the token embeddings with the feature embeddings to represent the information of the tweet and used these embeddings as the input of their models.

4.3 Models

Some participants considered the tasks as sequential labeling tasks [1, 7], while some participants formulated the tasks as classification tasks [13, 15, 14, 16]. The abstracts of the participants' models are shown as follows.

- Ait Azzi and Bouamor [1] proposed a CNN-based model with enriched word representation, called E-CNN. They used a fusion model to integrate the fine-tuned model for subtask 1 into E-CNN for subtask 2.
- Spark [13] used two-layer rectified linear units (ReLU) as a classifier with both tweet features and number features.
- Liang and Su [7] developed a recurrent neural networks (RNN) model with CNN filter, and made comparison with both CNN and RNN models.
- Wang et al. [15] used SVM model in the formal run, and adopted the BERT model after the evaluation results released.
- Tian and Peng [14] constructed an attention-based LSTM model for the shared task.
- Wu et al. [16] used multi-layer perceptron (MLP) for target numeral and used LSTM for the preceding context and the posterior context of the target numeral.

5 Results and Discussions

In FinNum shared task, each participant can submit at most two results for evaluation. Table 2 shows the experimental results of participants' submissions in the formal run for subtask 1 and subtask 2, including Fortia1 [1], DeepMRT [16], ASNLU [7], aiai [14], WUST [15], and BRNIR [13]. Fortia1 [1] achieved the first place with their E-CNN model in both tasks. WUST [15] showed that BERT model pretrained with Microsoft Research Paraphrase Corpus (MRPC) can obtain the best performance, (94.50, 88.62) and (87.25, 83.07), in both tasks.

We further analyze the errors of all participants' submissions. Figure 1 shows the confusion matrix of subtask 1. We find that models can achieve about 90% accuracy in Monetary, Percentage, and Temporal categories. The most challenging one is the Product/Version category, since many product names are low

Table 2. Experimental results

Subtask 1			Subtask 2		
Submission ID	Micro F1 (%)	Macro F1 (%)	Submission ID	Micro F1 (%)	Macro F1 (%)
Fortia1 - 1	93.94	90.05	Fortia1 - 2	87.17	82.40
Fortia1 - 2	93.70	88.98	Fortia1 - 1	86.53	80.49
DeepMRT - 1	91.87	87.94	DeepMRT - 1	83.03	77.90
DeepMRT - 2	91.16	84.72	DeepMRT - 2	81.27	75.59
ASNLU - 2	89.72	80.93	aiai - 1	80.24	74.11
ASNLU - 1	89.40	79.96	aiai - 2	80.64	73.43
ZHAW - 2	86.45	79.27	ASNLU - 1	79.12	72.51
Fortia2 - 1	89.88	79.26	ASNLU - 2	77.37	70.09
Fortia2 - 2	87.73	78.59	Fortia2 - 2	77.05	68.86
aiai - 1	86.45	78.09	Fortia2 - 1	79.28	68.33
aiai - 2	87.41	78.04	ZHAW - 2	75.54	66.44
ZHAW - 1	84.78	75.40	ZHAW - 1	72.67	64.84
WUST	74.02	63.71	Stark - 1	69.08	56.83
BRNIR - 1	74.27	63.53	WUST	60.88	52.93
Stark - 1	78.01	61.75	BRNIR - 1	63.67	51.90
BRNIR - 2	72.91	58.54	BRNIR - 2	61.99	47.14
word-based CNN [4]	55.90	51.67	char-based CNN [4]	43.75	31.12

frequent words in the training set. Some of them are even the out of vocabulary words. Quantity category may be confused with Monetary and Temporal categories. The tailor-made categories, i.e., Option and Indicator, for financial social media also get lower accuracy. Option category has two subcategories: exercise price and maturity date. Exercise price can be considered as one kind of monetary information, and maturity date can be seen as one kind of temporal information. That is the reason why models may not distinguish Option category from Monetary and Temporal categories. The value of the indicator may sometimes be related to the price of the target financial instrument, and the parameter of the indicator is always related to temporal. That results in some wrong predictions to Monetary and Temporal categories.

Figure 2 shows the confusion matrix of subtask 2. The models perform better on the subcategories of Percentage and Temporal categories than on other subcategories. This is because the subcategories have some patterns such as “+/-” sign + target numeral + “%” sign and “YYYY/MM/DD”. Same as the results in subtask 1, most errors of “exercise price” subcategory happen in “quote” and “date”, and that of “maturity date” subcategory occurs in the “date” subcategory. The in-depth analysis to Monetary category is need, because some subcategories in Monetary category contain fruitful fine-grained crowd opinions. “quote” is the major subcategory in Monetary category. It becomes the most confusing subcategory, because models tend to label the numeral as “quote” when the uncertainty is high. As we mentioned in Section 2.1, “buy price”, “sell price”, “forecast”, “stop loss” and “support or resistance” are the most important subcategories, because these subcategories can be used to capture the

Overview of the NTCIR-14 FinNum Task 9

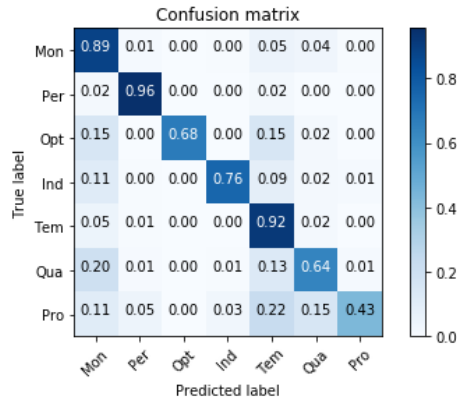


Fig. 1. Confusion matrix of subtask 1.

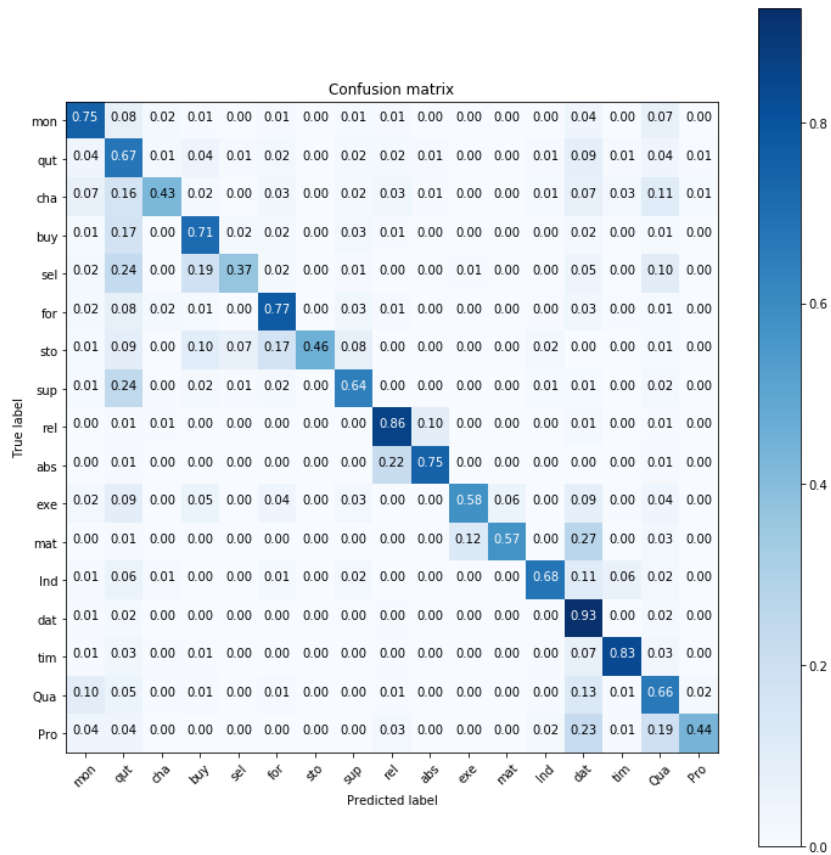


Fig. 2. Confusion matrix of subtask 2.

fine-grained opinions of the crowd. “forecast” achieves the best performance of these five subcategories, and we had showed the usefulness of this subcategory in our previous work [4]. Due to the small number of occurrences of “selling price” and “stop loss”, these two subcategories perform worse than others.

Note that, even with the E-CNN model [1], the accuracy of the “selling price”, “stop loss”, and “support or resistance” subcategories are only 50%, 50%, and 70%, respectively. The accuracy of “exercise price” and “maturity date” are 71% and 62%. The accuracy of Indicator category is 81%. These results indicate that there are still room for improving the performance of capturing fine-grained opinions from the crowd in financial social media data.

In sum, the overall experimental results show that although we now perform well in traditional (sub)categories such as Monetary, Percentage and Temporal categories, there are still some challenges in domain specific (sub)categories like Option, Indicator, “selling price”, “stop loss”, and “support and resistance”. The issues for Quantity and Product/Version categories also need to be resolved.

6 Conclusion and Future Work

According to World Economic Forum 2015, social trading is one of the crucial trends in FinTech tide. In the FinNum task, We presented novel and important issues in analyzing the numerals in financial social media data in a fine-grained way, and provided a large and high quality dataset to lead the new track of numeral understanding. The proposed tasks are the vanguard of in-depth opinion mining for financial social media data.

Participants proposed several features for understanding the meaning of the given numerals, and leveraged by neural network models. The results of E-CNN and BERT show that embeddings play an important role for FinNum shared task. Although the micro-averaged F1 of 12 runs in subtask 1 and 6 runs in subtask 2 is higher than 0.80, there are still spaces to improve for some critical (sub)categories.

The proposed taxonomy can be extended to other kinds of documents, including analyst reports and financial reports. For a future edition of this task, we will present the annotations related to different fine-grained aspects of financial social media data.

Acknowledgments

We greatly appreciate the efforts of all the participants in the FinNum shared task at NTCIR-14.

This shared task was partially supported by Ministry of Science and Technology, Taiwan, under grants MOST 107-2218-E-009-050-, MOST-106-2923-E-002-012-MY3, MOST-107-2634-F-002-011-, MOST-108-2634-F-002-008-, and MOST-107-2218-E-009-050-, and by Academia Sinica, Taiwan, under grant AS-TP-107-M05.

References

1. Ait Azzi, A., Bouamor, H.: Fortial at the ntcir-14 finnum task: Enriched sequence labeling for numeral classification. In: Proceedings of the 14th NTCIR Conference on Evaluation of Information Access Technologies (2019)
2. Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. *Journal of machine Learning research* **3**(Jan), 993–1022 (2003)
3. Brown, P.F., Desouza, P.V., Mercer, R.L., Pietra, V.J.D., Lai, J.C.: Class-based n-gram models of natural language. *Computational linguistics* **18**(4), 467–479 (1992)
4. Chen, C.C., Huang, H.H., Shiue, Y.T., Chen, H.H.: Numeral understanding in financial tweets for fine-grained crowd-based forecasting. In: 2018 IEEE/WIC/ACM International Conference on Web Intelligence (WI). pp. 136–143. IEEE (2018)
5. Chen, C.C., Huang, H.H., Tsai, C.W., Chen, H.H.: Crowdpt: Summarizing crowd opinions as professional analyst. In: Proceedings of the 2019 World Wide Web Conference (WWW) (2019)
6. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. Proceedings of the conference of the North American chapter of the association for computational linguistics (2019)
7. Liang, C.C., Su, K.Y.: Asnlu at the ntcir-14 finnum task. In: Proceedings of the 14th NTCIR Conference on Evaluation of Information Access Technologies (2019)
8. Manning, C., Surdeanu, M., Bauer, J., Finkel, J., Bethard, S., McClosky, D.: The stanford corenlp natural language processing toolkit. In: Proceedings of 52nd annual meeting of the association for computational linguistics: system demonstrations. pp. 55–60 (2014)
9. Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J.: Distributed representations of words and phrases and their compositionality. In: Advances in neural information processing systems. pp. 3111–3119 (2013)
10. Owoputi, O., O’Connor, B., Dyer, C., Gimpel, K., Schneider, N., Smith, N.A.: Improved part-of-speech tagging for online conversational text with word clusters. In: Proceedings of the 2013 conference of the North American chapter of the association for computational linguistics: human language technologies. pp. 380–390 (2013)
11. Pennington, J., Socher, R., Manning, C.: Glove: Global vectors for word representation. In: Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP). pp. 1532–1543 (2014)
12. Peters, M.E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., Zettlemoyer, L.: Deep contextualized word representations. In: Proceedings of the conference of the North American chapter of the association for computational linguistics (2018)
13. Spark, A.: Brnir at the ntcir-14 finnum task: Scalable feature extraction technique for numeral classification. In: Proceedings of the 14th NTCIR Conference on Evaluation of Information Access Technologies (2019)
14. Tian, K., Peng, Z.J.: aiai at the ntcir-14 finnum task: Financial numeral tweets fine-grained classification using deep word and character embedding-based attention model. In: Proceedings of the 14th NTCIR Conference on Evaluation of Information Access Technologies (2019)
15. Wang, W., Liu, M., Zhang, Z.: Wust at the ntcir-14 finnum task. In: Proceedings of the 14th NTCIR Conference on Evaluation of Information Access Technologies (2019)

12 CC. Chen et al.

16. Wu, Q., Wang, G., Zhu, Y., Liu, H., Karlsson, B.: Deepmrt at the ntcir-14 finnum task: A hybrid neural model for numeral type classification in financial tweets. In: Proceedings of the 14th NTCIR Conference on Evaluation of Information Access Technologies (2019)