

# **Reliable Sentiment Analysis in Social Media**

**by Qian Li**

Thesis submitted in fulfilment of the requirements for  
the degree of

**Doctor of Philosophy**

under the supervision of Qiang Wu, Jian Zhang

University of Technology Sydney  
Faculty of Engineering and Information Technology

July, 2020



## CERTIFICATE OF ORIGINAL AUTHORSHIP

I, *Qian Li* declare that this thesis, submitted in partial fulfilment of the requirements for the award of Doctor of Philosophy, in the declare that this thesis, is submitted in fulfilment of the requirements for the award of *Doctor of Philosophy*, in the *Faculty of Engineering and Information Technology* at the University of Technology Sydney.

This thesis is wholly my own work unless otherwise referenced or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

This document has not been submitted for qualifications at any other academic institution. This research is supported by the Australian Government Research Training Program.

Production Note:

SIGNATURE: Signature removed prior to publication.

Qian Li

DATE: 29 July, 2020



## ABSTRACT

Sentiment analysis in social media is critical yet challenging because the source materials (i.e., reviews posted in social media) are with high complexity, low quality, and uncertain credibility. For example, words and sentences in a textual review may couple with each other, and they may have heterogeneous meanings under different contexts or in different language locales. These couplings and heterogeneities essentially determine the sentiment polarity of the review but are too complex to be captured and modeled. Also, social reviews contain a large number of informal words and typos (a.k.a., noise) but a rare number of vocabularies (a.k.a., sparsity). As a result, most of the existing natural language processing (NLP) methods may fail to represent social reviews effectively. Furthermore, a large proportion of social reviews are posted by fraudsters. These fraud reviews manipulate social opinion, and thus, they disturb sentiment analysis.

This research focuses on reliable sentiment analysis in social media. It systematically investigates the sentiment analysis techniques to tackle three major challenges in social media: high data complexity, low data quality, and uncertain credibility. Specifically, this research focuses on two research problems: general sentiment analysis in social media and fraudulent sentiment analysis in social media. The general sentiment analysis targets on tackling high data complexity and low-quality of social articles that are credible. The fraudulent sentiment analysis handles the uncertain credibility issue, which is common and profoundly affects the precise sentiment analysis in social media. Based on these investigations, this research proposes a series of methods to achieve reliable sentiment analysis: It studies the polarity-shift characteristics and non-IID characteristics in general paragraphs to capture the sentiment more accurately. It further models multi-granularity noise and sparsity in short text, which is the most common data in social media, for robust short text sentiment analysis. Finally, it tackles the uncertain credibility problem in social media by studying fraudulent sentiment analysis in both supervised and unsupervised scenarios.

This research evaluates the performance and properties of the proposed reliable sentiment analysis methods by extensive experiments on large real-world data sets. It demonstrates that the proposed methods are superior and reliable in social media sentiment analysis.



## DEDICATION

*To my son Jewell Khu...*





## ACKNOWLEDGMENTS

First of all, I want to thank all the contributions my supervisor Prof. Qiang Wu make during my Ph.D experience. I appreciate it very much for his time, ideas, discussions, and paper polishing to make my Ph.D. studies impressive and stimulating. I am deeply influenced by his professionalism, profound ideas and hard work. Also, his optimism attitude and enthusiasm were contagious and inspiring for me, even during the difficult times in the pursuit of my doctoral degree.

I am also greatly indebted to Prof. Jian Zhang, my co-supervisor, whose constant suggestions, incisive comments and encouragements helped me a lot during my Ph.D. studies.

My sincere thanks also go to all of my colleagues in the Global Big Data Technologies Centre, for their enthusiastic help: Dr. Peng Zhang, Dr. Zongjian Zhang, Xunxiang Yao, Dr. Xiaoshui Huang, Yan Huang, Lu Zhang, Yongshun Gong, Xiaofei Liu, Anan Du, Lingxiang Yao. I appreciate our intensive discussions, in-depth collaboration, and all the fun we have had in the last four years.

I thank the colleagues in the student union. They broadened my horizon and shared unforgettable times with me in my study abroad. In particular, I am grateful to Mr. Wenqi Niu, Mrs. Jing Zhao, Mr. Yongfeng Hou, Dr. Haodong Chang, Dr. Yu Suo, Dr. Yu Chen, Dr. Yan Wang, Dr. Kunpeng Zhao and Zhixun Zhao for their precious advice and encouragement in my Ph.D. research.

I gratefully acknowledge the funding from China Scholarship Council. Without its funding, my Ph.D work would not be possible.

Lastly, I would like to thank my family for all their selfless love and support. Without their encouragement, it would not be possible for me to conduct this research. In particular, my grateful thanks goes to my husband Dr. Chengzhang Zhu whose spiritual and academic support are so appreciated. Thank you.

Qian Li  
*University of Technology Sydney*  
July 2020



## LIST OF PUBLICATIONS

### RELATED TO THE THESIS :

#### Published Papers :

1. **Li, Q.**, Wu, Q., Zhu, C., and Zhang, J., 2019, September. Bi-level masked multi-scale CNN-RNN networks for short text representation *In International Conference on Document Analysis and Recognition*. IEEE.
2. **Li, Q.**, Wu, Q., Zhu, C., Zhang, J. and Zhao, W., 2019, July. An inferable representation learning for fraud review detection with cold-start problem. *In 2019 International Joint Conference on Neural Networks*. IEEE.
3. **Li, Q.**, Wu, Q., Zhu, C., Zhang, J. and Zhao, W., 2019, April. Unsupervised user behavior representation for fraud review detection with cold-start problem. *In Pacific-Asia Conference on Knowledge Discovery and Data Mining* (pp. 222-236). Springer, Cham.
4. **Li, Q.**, Wu, Q., Liu, X., 2019. Multi-scale and hierarchical embedding for polarity shift sensitive sentiment classification. *Lecture Notes in Computer Science(including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 11633LNCS, pp.227-238.
5. Da, Q., Cheng, J., **Li, Q.**, Zhao, W., 2019. Socially-attentive representation learning for cold-start fraud review detection. *Communications in Computer and Information Science*, 1069, pp.76-91.

#### Submitted Manuscripts:

6. **Li, Q.**, Wu, Q., Zhu, C., Zhang, J., Zhao, W., 2020. Inferable representation learning for cold-start fraud review detection. *Neurocomputing* (Accepted).
7. **Li, Q.**, Wu, Q., Zhu, C., Cao, L., Zhang, J., 2020. Deep Non-IID Paragraph Representation for Sentiment Analysis. *IEEE Access* (Under Review).

- 
8. Zhu, C., Cao, L., and **Li, Q.**,2020. Unsupervised hierarchical and heterogeneous coupling learning on dynamic categorical data. *Journal of Machine Learning Research* (Under Review).

**OTHERS :**

9. Zhao, W., **Li, Q.**, Zhu, C., Song, J., Liu, X. and Yin, J., 2018. Model-aware categorical data embedding: a data-driven approach. *Soft Computing*, 22, pp.3603-3619.
10. Zhu, C., Zhao, W.,**Li, Q.**, Li, P., Da, Q., 2019. Network embedding-based anomalous density searching for multi-group collaborative fraudsters detection in social media. *Computers, Materials and Continua*, 60(1), pp.317-333.
11. Xiang, L., Zhao, G., **Li, Q.**, Hao, W., Li, F., 2018. TUMK-ELM: A fast unsupervised heterogeneous data learning approach. *IEEE Access*, 6, pp.35305-35315.
12. Zheng, Y., Zhao, W., Sun, C., **Li, Q.**, 2019. Drug side-effect prediction using heterogeneous features and bipartite local models. *Computers, Materials and Continua*, 60(2), pp.481-496.
13. Cui, J., Long, J., Min, E., Liu, Q., **Li, Q.**, 2018. Comparative study of CNN and RNN for deep learning based intrusion detection system. *ICCCS 2018: Cloud Computing and Security (LNCS)*, 11067, pp. 159-170.

## TABLE OF CONTENTS

<b>List of Publications</b>	<b>ix</b>
<b>List of Figures</b>	<b>xvii</b>
<b>List of Tables</b>	<b>xix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background . . . . .	1
1.2 Current Work and Gap Analysis . . . . .	2
1.3 Research Problems and Objectives . . . . .	4
1.4 Thesis Contributions . . . . .	7
1.5 Thesis Organization . . . . .	8
<b>2 Literature Review</b>	<b>11</b>
2.1 Introduction . . . . .	11
2.2 General Sentiment Analysis Methods . . . . .	12
2.2.1 Hand-Craft Features for Sentiment Analysis . . . . .	12
2.2.2 Homogeneous Sentiment Analysis . . . . .	12
2.2.3 Heterogeneous Sentiment Analysis . . . . .	16
2.2.4 Domain-level Heterogeneous Sentiment Analysis . . . . .	20
2.2.5 Short Text Sentiment Analysis . . . . .	22
2.3 Fraudulent Sentiment Analysis Methods . . . . .	23
2.3.1 Text-Content-Based Fraudulent Sentiment Analysis . . . . .	23
2.3.2 Rating-Distribution-Based Fraudulent Sentiment Analysis . . . . .	25
2.3.3 Social-Relation-Based Fraudulent Sentiment Analysis . . . . .	26
2.3.4 Cold-Start Fraudulent Sentiment Analysis . . . . .	27
<b>3 Polarity-Shift Sentiment Analysis</b>	<b>29</b>
3.1 Introduction . . . . .	29

## TABLE OF CONTENTS

---

3.2	Sentiment-Polarity-shift Characteristics . . . . .	32
3.2.1	The Multiple Forms of Polarity-Shift . . . . .	32
3.2.2	The Hierarchy of Polarity-Shift with Long-Term Dependence . . .	33
3.3	Multi-Scale and Hierarchical Network for Polarity-Shift Sentiment Analysis	33
3.3.1	Sentence-Level Embedding Module . . . . .	34
3.3.2	Paragraph-Level Embedding Module . . . . .	35
3.4	Experiments and Evaluation of Multi-Scale and Hierarchical Network . .	36
3.4.1	Data sets . . . . .	37
3.4.2	Experimental Settings . . . . .	38
3.4.3	Evaluation on Sentiment Classification Performance Enabled by Multi-Scale and Hierarchical Network . . . . .	39
3.4.4	Evaluation on Polarity-Shift Capturing of Multi-Scale and Hierar- chical Network . . . . .	40
3.5	Summary . . . . .	42
<b>4</b>	<b>Sentiment Analysis on Non-Independent and Identical Distributed Para- graph</b> . . . . .	<b>43</b>
4.1	Introduction . . . . .	43
4.2	Non-Independent and Identically Distributed Paragraph Representation Design . . . . .	46
4.2.1	Characteristics of Non-Independent and Identically Distributed Paragraph . . . . .	46
4.2.2	Objectives of Non-Independent and Identically Distributed Para- graph Representation . . . . .	47
4.2.3	Architecture for Representing Non-Independent and Identically Distributed Paragraph . . . . .	49
4.3	Implement of Non-Independent and Identically Distributed Paragraph Representation Design . . . . .	51
4.3.1	Attentive Multi-Scale Convolutional-Recurrent Neural Network .	52
4.3.2	Explicit Features used by the Multi-Scale and Hierarchical Deep Neural Network . . . . .	54
4.4	Experiments and Evaluation of the Multi-Scale and Hierarchical Para- graph Representation Performance . . . . .	56
4.4.1	Data Sets for the Performance Evaluation of Non-Independent and Identically Distributed Representation Representation . . . . .	57

4.4.2	Experimental Settings . . . . .	58
4.4.3	Evaluation on Sentiment Classification Performance Enabled by the Non-Independent and Identically Distributed Representation Method . . . . .	59
4.4.4	Evaluation on Representation Quality of the Non-Independent and Identically Distributed Representation Method . . . . .	61
4.4.5	Evaluation Significance of Hierarchically Integrating Explicit Fea- tures with Implicit Features . . . . .	63
4.5	Summary . . . . .	65
<b>5</b>	<b>Sentiment Analysis on Short Text</b>	<b>67</b>
5.1	Introduction . . . . .	67
5.2	Bi-level Masked Multi-Scale Convolutional-Recurrent Neural Network for Short Text Sentiment Analysis . . . . .	70
5.2.1	The Architectural of Bi-level Masked Multi-Scale Convolutional- Recurrent Neural Network . . . . .	70
5.2.2	The Transformation of Short Textual Data . . . . .	70
5.2.3	Masked Convolutional Neural Network . . . . .	71
5.2.4	Multi-Scale Convolutional-Recurrent Neural Network Structural .	73
5.3	Experiments and Evaluation of Sentiment Analysis on Short Text . . . . .	74
5.3.1	Data Sets . . . . .	75
5.3.2	Experimental Settings . . . . .	75
5.3.3	Evaluation on Bi-level Masked Multi-Scale Convolutional-Recurrent Neural Network Performance on Short Text Classification . . . . .	76
5.3.4	Evaluation on Bi-level Masked Multi-Scale Convolutional-Recurrent Neural Network Performance on Short Text Retrieval . . . . .	77
5.3.5	Evaluation on Short Text Representation Quality of Bi-level Masked Multi-Scale Convolutional-Recurrent Neural Network . . . . .	78
5.3.6	Summary . . . . .	79
<b>6</b>	<b>Fraudulent Sentiment Analysis</b>	<b>81</b>
6.1	Introduction . . . . .	81
6.2	The Framework of Inferable Representation Learning for Fraudulent Sentiment Analysis . . . . .	84
6.2.1	Preliminaries . . . . .	84
6.2.2	Framework . . . . .	85

## TABLE OF CONTENTS

---

6.3	Inferable Representation Space Building . . . . .	86
6.3.1	Representation Learning Networks . . . . .	86
6.3.2	Co-occurrence-based User Reviewing Behavior Learning . . . . .	86
6.3.3	User Representation Inferring in Cold-Start Problem . . . . .	89
6.4	Supervised Fraud-Sensitive Information Embedding . . . . .	90
6.4.1	Social Relation Embedding . . . . .	91
6.4.2	Statistical Information Embedding . . . . .	95
6.5	Supervised Fraud Detection Method . . . . .	95
6.6	Unsupervised Fraud-Sensitive Information Embedding . . . . .	96
6.6.1	Social Relation Mining . . . . .	97
6.6.2	Fraud-Sensitive Information Embedding . . . . .	99
6.6.3	Dynamic Re-Weighting Strategy . . . . .	100
6.7	Unsupervised Fraud Detection Method . . . . .	100
6.8	Experiments and Evaluation of Fraudulent Sentiment Analysis . . . . .	101
6.8.1	Data Sets . . . . .	101
6.8.2	Evaluation Metrics . . . . .	102
6.8.3	Parameters Settings . . . . .	103
6.8.4	Effectiveness on Supervised Cold-start Fraudulent Detection . . . . .	104
6.8.5	Effectiveness on Supervised General Fraudulent Sentiment Analysis . . . . .	105
6.8.6	Effectiveness on Unsupervised Cold-start Fraudulent Sentiment Analysis . . . . .	106
6.8.7	Effectiveness on Unsupervised General Fraudulent Sentiment Analysis . . . . .	107
6.8.8	Evaluating the Effectiveness of User Reviewing Behavior and User/Item Social Relations for Fraudulent Sentiment Analysis . . . . .	109
6.9	Summary . . . . .	109
<b>7</b>	<b>Conclusions and Future Directions</b>	<b>111</b>
7.1	Conclusions . . . . .	111
7.1.1	Paragraph Sentiment Analysis . . . . .	112
7.1.2	Short Text Sentiment Analysis . . . . .	112
7.1.3	Supervised Fraudulent Sentiment Analysis . . . . .	112
7.1.4	Unsupervised Fraudulent Sentiment Analysis . . . . .	113
7.2	Future Directions . . . . .	113



7.2.1	Exploiting Other Support Techniques for Reliable Sentiment Analysis . . . . .	113
7.2.2	Studying Multi-Granular Reliable Sentiment Analysis . . . . .	114
7.2.3	Exploring the Interpretability of Reliable Sentiment Analysis . . .	115
<b>A</b>	<b>Appendix</b>	<b>117</b>
A.1	List of Notations . . . . .	117
	<b>Bibliography</b>	<b>121</b>



## LIST OF FIGURES

FIGURE	Page
1.1 The research problems and their relations. . . . .	5
3.1 The MUSH Structure: The primary structure of MUSH is a multi-scale and hierarchical neural network that consists of two modules (i.e., sentence-level embedding module and paragraph-level embedding module). The MUSH method learns paragraph representation for polarity-shift-sensitive classification through an end-to-end training process guided by a specific task. . . . .	34
3.2 The Sentence-level Embedding Module Structure: The first layer of this module is a word embedding layer. This word embedding layer then connects to several CNN layers with multiple filter sizes. Each of these CNN layers follows a bidirectional gated RNN layer. The outputs of these bidirectional gated RNN layers are merged by fully connected layers to form a sentence representation. . . . .	36
3.3 The Paragraph-level Embedding Module Structure: The CNN layers with multiple filter sizes followed by bidirectional gated RNN layers are used to capture multi-scale polarity-shift. The fully connected layers merge the outputs of RNN layers to form a paragraph representation. . . . .	37
4.1 An Example of the Non-IID Characteristics in a Paragraph. . . . .	47
4.2 The Non-IID Paragraph Representation Framework for Sentiment Analysis.	50
4.3 The Architecture of the non-IID-Characteristics-Learning Module in MEDEA.	53
4.4 The RMSE@K of Different Methods on Yelp13. . . . .	62
4.5 The Visualization of Different Representations through t-SNE Transformation on Yelp13. . . . .	63
4.6 The Training and Validation Accuracy of MEDEA on Data Set Yelp13. . . . .	65
5.1 The Architecture of Bi-MACRO. . . . .	71

## LIST OF FIGURES

---

5.2	Masked Convolutional Network. The shadow parts refer to masks. . . . .	72
5.3	Short Text Representation of Different Methods. . . . .	79
6.1	The Framework of Inferable Representation Learning for Fraudulent Sentiment Analysis . . . . .	86
6.2	Architecture of the proposed Inferable Representation Space Building Method	87
6.3	The Inferable Representation Space. In this figure, <b>u, t, d, r</b> refer to the representations of user, item, review, and rating, respectively. . . . .	88
6.4	User Representation Inferring Process in Cold-Start Problem . . . . .	90
6.5	The New User Representation Inferring. . . . .	91
6.6	Architecture of the Proposed JESTER Method . . . . .	92
6.7	The Users/Items Social Relation Embedding Workflow. . . . .	92
6.8	Architecture of the Proposed URBER Method . . . . .	97
6.9	Density in a User-Item Bipartite Graph . . . . .	98
6.10	User Representation with Density of Different Methods on Yelp-Hotel and Yelp-Restaurant. The sub-figures (a), (b), (c), (d) contain the user representation information with the ground-truth labels, and the sub-figures (e), (f), (g), (h) show the density in the representation space. S refers to the social relation embedding-based method, and B refers to the behavior embedding-based method. . . . .	110

## LIST OF TABLES

<b>TABLE</b>	<b>Page</b>
3.1 Characteristics of Data Sets in the Experiments: #v refers to the number of vocabulary, #s and #w refers to the number of sentences and words, respectively.	38
3.2 Sentiment Classification Accuracy of Different Methods . . . . .	40
3.3 Polarity-Shift Capturing Effectiveness of Different Methods. The predicted sentiment ratings are reported. A larger rating indicates a stronger positive sentiment, and vice versa. . . . .	41
4.1 Statistical Properties of Data Sets: #s and #w refer to the number of sentences and words, respectively. . . . .	57
4.2 Sentiment Classification Accuracy of Different Methods . . . . .	60
4.3 Sentiment Prediction RMSE of Different Methods . . . . .	61
4.4 Sentiment Analysis Performance based on MEDEA and Its Variants . . . . .	64
5.1 The Data Characteristics of Each Short-Text Data Set. . . . .	75
5.2 The Short Text Classification Performance based on Different Representation Methods. . . . .	77
5.3 The Short Text Retrieval Performance based on Different Representation Methods. . . . .	78
6.1 Statistics of Data Sets for Supervised Fraudulent Sentiment Analysis . . . .	102
6.2 Supervised Cold-start Fraud Detection Performance of Different Methods . .	105
6.3 Supervised General Fraud Detection Performance of Different Methods . . .	106
6.4 Unsupervised Cold-start Fraud Detection of Different Methods . . . . .	107
6.5 Unsupervised General Fraud Detection of Different Methods . . . . .	108

