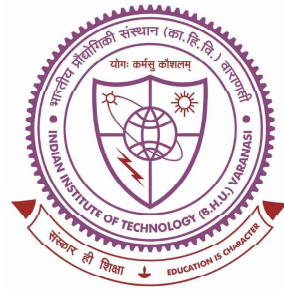


The Impact of Social Media Text: A study on e-Governance, Election and Disaster Management



**Thesis submitted in partial fulfilment
for the Award of
DOCTOR OF PHILOSOPHY (PHD)
in
COMPUTER SCIENCE AND ENGINEERING**

by
Anita Saroj

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
INDIAN INSTITUTE OF TECHNOLOGY
BANARAS HINDU UNIVERSITY
VARANASI - 221 005

ROLL NUMBER
16071006

YEAR OF SUBMISSION
2022

I would like to dedicate this thesis to my loving parents ...

Certificate

It is certified that the work contained in this thesis entitled **“The Impact of Social Media Text: A study on e-Governance, Election and Disaster Management“** by **“Anita Saroj“**, Roll Number **16071006** has been carried out under my supervision and that this work has not been submitted elsewhere for a degree.

It is further certified that the student has fulfilled all the requirements of Comprehensive Examination, Candidacy and SOTA for the award of **Ph.D. Degree in Computer Science and Engineering**.



Supervisor

Dr. Sukomal Pal

Department of Computer Science and Engineering

Indian Institute of Technology (Banaras Hindu University)

Varanasi-221005

Declaration

I, **Anita Saroj**, certify that the work embodied in this thesis is my own bona-fide work and carried out by me under the supervision of (**Dr. Sukomal Pal**) from July 2016 to July 2022 at the **Department of Computer Science and Engineering**, Indian Institute of Technology (BHU), Varanasi. The matter embodied in this thesis has not been submitted for the award of any other degree/diploma. I declare that I have faithfully acknowledged and given credits to the research workers whenever and wherever their works have been cited in my work in this thesis. I further declare that I have not wilfully copied any others' work, paragraphs, text, data, results, etc., reported in journals, books, magazines, reports dissertations, theses, etc., or available at websites and have not included them in this thesis and have not cited as my own work.

Date 14 July 2022


Signature

Place Varanasi

Anita Saroj

Certificate by the Supervisor

It is certified that the above statement made by the student is correct to the best of my knowledge.


Supervisor

Dr. Sukomal Pal

Department of Computer Science and Engineering

IIT (BHU)

Varanasi


Signature of the Head of the Department

(Computer Science and Engineering)

Professor & Head

समयक विज्ञान एवं अभियांत्रिकी विभाग
Department of Computer Sc. & Engg

भारतीय प्रौद्योगिकी संस्थान

Indian Institute of Technology

(वाराणसी हिन्दू विश्वविद्यालय)

(Banaras Hindu University)

वाराणसी- 221005 / Varanasi-221005

Copyright Transfer Certificate

Title of the Thesis : **The Impact of Social Media Text: A study on e-Governance, Election and Disaster Management**

Name of the Student : **Anita Saroj**

Copyright Transfer

The undersigned hereby assigns to the Institute of Technology (Banaras Hindu University) Varanasi all rights under copyright that may exist in and for the above thesis submitted for the award of the **Doctor of Philosophy** in **Computer Science and Engineering**.

Date 14 July 2022

Signature anita

Place Varanasi

(Anita Saroj)

Note: However, the author may reproduce or authorize others to reproduce material extracted verbatim from the thesis or derivative of the thesis for author's personal use provided that the source and the Institute's copyright notice are indicated.

Acknowledgements

This thesis is the end of my journey in obtaining my Ph.D. I have not traveled in a vacuum in this journey. This thesis has been kept on track and been seen through to completion with the support and encouragement of numerous people including my well wishers, my friends, colleagues and various institutions. At the end of my thesis, I would like to thank all those people who made this thesis possible and an unforgettable experience for me. At the end of my thesis, it is a pleasant task to express my thanks to all those who contributed in many ways to the success of this study and made it an unforgettable experience for me.

Firstly, I would like to express my sincere gratitude to my advisor Dr.Sukomal Pal for the continuous support of my Ph.D study and related research, for his patience, motivation, and immense knowledge. His guidance helped me in all the time of research and writing of this thesis. I could not have imagined having a better advisor and mentor for my Ph.D study.

I specially thank the members of my Research Progress Evaluation Committee, Dr. A. K. Singh, Dr. Sandip Ghosh , DPGC Convener Dr. Bhaskar Biswas, Dr. Pratik Chattopadhyay for their invaluable suggestions regarding the thesis, their insightful comments and encouragement, but also for the hard question which incented me to widen my research from various perspectives.

Most importantly, my deepest gratitude is for my family for their constant support, inspiration, guidance, and sacrifices. My parents were constant source of motivation and inspiration. Their affection and guidance was instrumental in me choosing Engineering and eventually continuing on to my Ph.D.

Very special gratitude goes out to my colleagues and friends with special mention to Mr. Rajesh Kumar Munodtiya, Mr. Siba Shankar Shau, and Mr. Supriya Chanda for being great friends and the best advisors I could be ever have. Their advice, encouragement, and critics were always sources of inspiration. This thesis would not have been possible without their invaluable remarks and persistent help. We really enjoyed our discussions, which spanned on every topic under the sun other than academics.

I extend special thanks to the non-teaching staff in the department, particularly, Mr. Ravi Kumar Bharti, Mr. Ritesh Singh, Mr. Shubham Pandey, Mr. Prakhar Kumar, and Mr. Manoj Kumar Rai, Mr. Viplav Biswas, and Mr. Akhilesh Kumar Pal.

I take this opportunity to sincerely acknowledge the Ministry of Human Resource Development, Government of India for providing Ph.D. Fellowship for financial assistance. It is said that it is no coincidence who you meet on the journey. I have had the fortunate of meeting so many wonderful people on the journey; capable and ready to guide, help, and advice me. I sincerely thank all of them who contributed in for helping me to see the light at the end of every scary tunnel during my Ph.D.

-Anita Saroj

Preface

Social media (SM) is part of our online existence today. We express ourselves on social media every day on different issues like politics, government policies, disasters, commercial products, celebrity statements, controversies and so on. Numerous kinds of social media platforms are available that host blogs, wikis, or multimedia content like YouTube, Instagram and Flickr, or enable social networking like Facebook and LinkedIn, micro blogging services like Twitter, Weibo etc. With Web 2.0, web-users are increasingly creating content than just consuming it. How this trend is entering into and thereby affecting our life and society is an interesting area of research. We chose, in particular, the social and political life and governance. An analysis and review of the adoption of social media by as well as its effects on the public and government sector is of due significance - examining the role and implications of these new technologies with a focus on different emotion detection and analysis in the chosen field from an Indian perspective is attempted in the thesis.

Many in the West believe that ICT-enabled public services positively impact economic growth, inclusion, and quality of life. With the ease-of-use, social media platforms are convenient vehicles for all forms of information exchange and continuously opening up newer avenues of governance, education, healthcare, entertainment and business as well as ushering in social changes. Even though information content in the media can be multi-modal involving image, audio and video, text is the predominant data-type. Social media text carries a lot of information. User-generated content provides diverse and unique information in the form of comments, posts, and tags. Features like friends, followers, connections provide a lot of contextual, structural and social network-related information about the users. The valuable information hidden in the text resources of social media provides opportunities for researchers of different disciplines to mine interesting patterns and metadata that might not be obvious. We conduct three studies to see the nature and

effect of social media communication through text in Indian context, keeping in mind two important stakeholders of the society: Government and general public. SM has been used in the government sector of the West for quite some time. In India, however, use of SM in governance is slowly but steadily increasing. We study the pattern of use, topics of discussion, interaction among different ministries and public participation in the governance in the first task. We look at another event of public-private interaction through social media. Elections are events when governments, political parties and public engage intensely with each other, especially in India. People's interaction with the candidates in social media posts reflects many social trends in a charged atmosphere. People's likes and dislikes of political parties and their leaders, their manifesto, promises and public comments often trigger hate and offensive posts. Sometimes irony is used while expressing opinions. We look at these varied sentiments in social media posts during a parliamentary election of India. Thirdly, we study the public reaction and sentiment from a different perspective like when a crisis or disaster happens. Crawling the SM channels during COVID-19 pandemic to collect data for the study and conducting different sentiment analysis related experiments is another task reported in the thesis.

First, we study the adoption and penetration of social media (SM) in government sector. The study explores the adoption of social media (SM) by different ministries of the Government of India (GoI) (actually 46 ministries). We analyze how government ministries use SM, particularly Twitter and Facebook, to disseminate information, engage with different stakeholders and take feedback on government initiatives. The study based on three years of SM data captures activities, discussion topics, inter-connectedness, public engagement, and popularity of GoI ministries through SM. Out of 8 different SM channels that GoI ministries use, almost all are active in Twitter and Facebook. Primary topics of discussion are meetings and projects like Prime Minister's (PM) projects, Railways projects, GST, and different development schemes. The Ministry of Railways and the

Ministry of Agriculture post maximum in SM, while PM's account is followed the most. The Ministry of Information and Broadcasting serves as a good coordinator. The study involves crawling the data from SM channels, pre-processing raw data, content analysis, time series analysis, citation analysis among different users.

Elections are events that causes boom in the use of SM. In India, governmental agencies, political parties and public engage intensely with each other during elections that sees waves of emotive posts in SM. People often use accusations, counter-accusations, hate speeches, ironies, sarcasms. Verbal irony, an utterance that conveys a spirit completely opposed to the surface meaning expressed, is usually understood by body language and the context of the conversation. However, it is challenging to automatically detect irony in a limited amount of text like in SM posts. To study the issue, we crawl the data from social media posts during the 2019 general election in India and make a standard collection with annotations (Indian General Election 2019 or IGE 2019 dataset). We then focus on automatic irony detection using various machine learning and deep learning (Bidirectional Encoder Representations from Transformers (BERT) and Embeddings from Language Models (ELMo)) models to classify them into irony and non-irony. We propose an ensemble model of machine learning and deep learning approaches. The classifiers are trained using a combined word embedding representation obtained from both BERT and ELMo. A series of experiments on irony detection then are performed including a domain adaptation with SemEval-2018 Task-3 (Sub-task A) dataset (SE-2018 T3 data). Our experiments on IGE 2019 and SemEval-2018 Task-3 data show results comparable to the state-of-the-art performance for irony detection. We also created a dataset on hate speech and offensive content identification during parliamentary election 2019 of India (PEI data-2019). It is a hierarchical classification task where we explore multi-task learning (MTL), a machine learning technique that accomplishes many tasks in a single model by exploiting commonalities and differences across different tasks. On the created dataset

we conduct hate speech and offensive content (HOF) identification and classification task, using MTL with the convolution network (MTL-CN) method. The MTL-CN extracts the shared and private latent features from the text. Our experiments on three different text classification tasks, demonstrate benefits of our approach. We also show that the knowledge learned by the proposed model here can be shared with new tasks. Finally, results on different standard datasets like LREC-2020 (PEI-2019), FIRE-2019, FIRE-2020, and SemEval-2019 datasets are shown to establish that multi-task learning yields better results for different classification tasks, irrespective of the languages.

When a disaster (man-made ones like shootings, bomb-blasts, war or natural calamities like floods, cyclones, tsunami, earthquakes) strikes, people suffer and/or get panicked and react emotionally. SM see a sudden increase with flight of emotions. During the COVID-19 pandemic, the whole world saw SM full with emotive posts. Social distancing caused lack of social interactions. The physical void led to increased online interaction among users on social media platforms. Sentiment analysis of such interactions can help us analyze the general public psychology during the pandemic. However, the lack of data in non-English and low-resource languages like ‘Hindi’ makes it difficult to study it among native and non-English speaking masses. We create a collection of ‘Hindi’ tweets on COVID-19 during the pandemic containing 10,011 tweets for sentiment analysis (named sentiment analysis for Hindi or SAFH dataset). We describe the process of collecting, creating, annotating the corpus, and performing sentiment classification task.

Table of contents

List of figures	xxiii
List of tables	xxv
Nomenclature	xxviii
1 Introduction	1
1.1 What are Social Media?	1
1.2 Why do we use Social Media?	3
1.2.1 Use of Social Media in E-governance	7
1.2.2 Social Media during Elections	8
1.2.3 Social Media during Disaster Management	9
1.3 Motivation and Challenges	10
1.4 Dissertation Overview	12
1.5 Research Goals	13
1.6 Contribution and Impact	14
1.6.1 A study on use of SM by Indian ministries	15
1.6.2 A study on people's opinions during IGE 2019	16
1.6.3 A study on people's reaction during IGE 2019	17
1.6.4 A study on COVID-19	18
1.7 Structure of the Thesis	19

2	Literature Review	21
2.1	Use of Social Media in e-governance	22
2.2	Social Media during Elections	26
2.3	Use of Social Media in Disaster Management	31
2.3.1	Extraction	33
2.3.2	Event Detection	39
2.3.3	Summarization	41
2.3.4	Classification	43
3	Background	47
3.1	Information Retrieval	47
3.1.1	Hyperlink-Induced Topic Search (HITS) Algorithm	50
3.1.2	Latent Dirichlet Allocation (LDA)	51
3.2	Traditional Machine Learning	53
3.2.1	Support vector machine	53
3.2.2	Multinomial Naive-Bayes (MNB)	54
3.2.3	Stochastic Gradient Descent (SGD)	55
3.2.4	Linear Regression (LR)	56
3.2.5	Random Forest	56
3.2.6	Decision Tree	57
3.2.7	XG Boost	57
3.2.8	k-Nearest Neighbour (kNN)	58
3.3	Deep Learning	58
3.3.1	Convolutional Neural Network	59
3.3.2	BiLSTM	61
3.4	Transfer Learning	62
3.5	Distributional Vector Representation	64

3.5.1	Context Independent Embedding	65
3.5.2	Context Dependent Embedding	65
4	A study on use of SM by Indian ministries	69
4.1	Short title	69
4.1.1	Twitter	70
4.1.2	Facebook	70
4.1.3	Social Media in Governance	71
4.2	Contributions of Chapter	72
4.2.1	Data Collection	73
4.3	Content Analysis	75
4.4	Level of Activities of different Government Ministries	78
4.4.1	Time-Series Analysis	80
4.4.2	Moderate Active Group	81
4.4.3	High Active Group	82
4.4.4	Super Active Group	83
4.5	Inter-connection among ministries	84
4.5.1	Hub and Authority scores	85
4.5.2	PageRank (PR)	86
4.6	Public Participation	87
4.6.1	Likes	88
4.6.2	Followers	90
4.7	Influence of Social media and public opinion in Government	93
4.8	Summary	94
5	A study on people's opinions during PEI 2019	97
5.1	Short title	97

5.2	Contribution	98
5.2.1	Multi-task learning with the convolution network (MTL-CN)	99
5.3	Experimental Setup	103
5.3.1	Datasets	104
5.3.2	Pre-processing	106
5.3.3	Implementation details	107
5.4	Results and Analysis	107
5.4.1	Effect of α -value	109
5.5	Summary	111
6	A study on people’s reaction during IGE 2019	113
6.1	Short title	113
6.2	Contribution	115
6.2.1	Characteristics of posts	115
6.2.2	Annotations	117
6.2.3	Preprocessing	117
6.3	Models	118
6.3.1	Machine Learning	118
6.3.2	Deep Learning	120
6.4	Our Approaches	123
6.4.1	Ensemble of Machine Learning techniques (EMLT)	124
6.4.2	Ensemble of BERT and ELMo models (EBEM)	125
6.4.3	Domain adaptation on each ensemble setting (EMLT and EBEM)	125
6.5	Results and Analysis	126
6.5.1	Performance of ML techniques	127
6.5.2	Performance of deep learning techniques	128
6.5.3	Performance of ensembling techniques	129

6.5.4	Results of Domain Adaptation	130
6.6	Summary	131
7	Sentiment Analysis on Hindi Tweets during COVID-19 Pandemic	133
7.1	Short title	133
7.2	Contributions	136
7.2.1	Annotation	137
7.2.2	Data Statistics	137
7.3	Methods	138
7.3.1	BERT	139
7.3.2	FastText	141
7.3.3	Self-attention	141
7.3.4	BiLSTM with BERT + Fasttext	142
7.3.5	BiLSTM+ fastText with self attention	143
7.3.6	BiLSTM with fastText	143
7.4	Experiments	144
7.4.1	Pre-processing	144
7.4.2	Benchmark Systems	145
7.5	Results	147
7.6	Summary	149
8	Discussion And Conclusion	151
8.1	Discussion	151
8.2	Summary and Contributions	152
8.2.1	RQ1: What roles do SM play in e-Governance and the public interaction with the government?	152

8.2.2	RQ2: How politics affect society (A study on people’s opinions during PEI 2019)?	155
8.2.3	RQ3: How do people react to political campaigns?	156
8.2.4	RQ4: How disaster affect people’s life (A study on COVID-19)? .	157
8.3	Conclusion	157
8.3.1	RQ1: What roles do SM play in e-Governance and the public interaction with the government?	158
8.3.2	RQ2: How politics affect society (A study on people’s opinions during IGE 2019)?	158
8.3.3	RQ3: How do people react to political campaigns?	159
8.3.4	RQ4: How disaster affect people’s life (A study on COVID-19)?	160
8.4	Future Scopes	160
	References	163
	Appendix A List of Publications	179

List of figures

1.1	Users of Social Media (in millions)	4
1.2	Usage of SM in general (<i>and in the thesis</i>)	5
1.3	Illustration of thesis structure	13
3.1	Illustration of a CNN architecture for sentence classification [1]	60
3.2	Sentiment Analysis based BiLSTM [2]	61
3.3	Overview of various transfer learning settings [3]	63
3.4	The architecture of embeddings from language models (ELMo)	67
3.5	BERT uses many layers of bidirectional transformers0	68
4.1	Topics discussed by ministries in both the platform (Facebook and Twitter)	77
4.2	Number of posts per user from Twitter and Facebook	79
4.3	Number of ‘like’ of all the ministries on Facebook and Twitter	89
4.4	Number of followers of all the ministries Facebook and Twitter account	90
4.5	News updated by PMO	94
5.1	Overview of the MTL-CN architecture.	100
5.2	Emotion based MTL-CN architecture.	101
5.3	Accuracy score of MTL-CN for English	110
5.4	Accuracy score of MTL-CN for Hindi	110
5.5	F_1 – score of Task A	111

5.6	F_1 – score of Task B	112
5.7	F_1 – score of Task C	112
6.1	An example of irony posts	117
6.2	An example of non-irony posts	117
6.3	Working module of ELMo	120
6.4	Working module of BERT [4]	122
6.5	Machine learning techniques based on Majority Voting Ensemble	124
6.6	Ensemble of BERT and ELMo model	125
6.7	Performance on IGE-2019 and SemEval-2018 in terms of Precision, Recall, F_1 -scores	126
6.8	Accuracy of machine learning models on IGE 2019 and SemEval 2018 datasets	127
6.9	Performance of deep learning models on IGE 2019 with SemEval 2018 dataset	128
6.10	Accuracy of deep learning models on IGE 2019 and SemEval 2018 datasets	128
6.11	Ensemble results of IGE 2019 with SemEval 2018 Task 3 data-set	129
6.12	Accuracy of ensemble of deep learning and machine learning models on IGE 2019 and SE 2018	130
7.1	SAFH data annotation	136
7.2	Architecture of the proposed work	139
7.3	Architecture of the BERT	140

List of tables

2.1	Use of SM in e-Governance	24
2.2	Use of SM in e-Governance cont.	26
2.3	Summary of the existing work on Hate Speech and Offensive content identification	29
2.4	Extraction Techniques	34
2.5	Extraction Techniques continue.	35
2.6	Event Detection Techniques	39
2.7	Summarization Techniques	42
2.8	Classification Techniques	44
4.1	Ministries of India on social media	73
4.2	Table 4.1 (Contd.) Ministries of India on social media	74
4.3	Social Media‘ Post Language Distribution of the Indian government Min- istries in %.	76
4.4	Twitter Content classification of Ministry Based on [5] in %	78
4.5	Facebook Content classification of Ministry Based on [5]	78
4.6	low active ministries from Facebook and Twitter	80
4.7	Moderate active ministries from Facebook and Twitter	82
4.8	High active ministries from Facebook and Twitter	82
4.9	Super active ministries from Facebook and Twitter	83

4.10	Ministries Hub, Authority and PageRank scores	88
4.11	Post Distribution of Ministry’s Social Media Accounts into Personal, Professional and Public in fraction of total post.	92
5.1	Label distribution in each task of PEI and FIRE datasets	104
5.2	Label distribution in each task of SemEval dataset	104
5.3	Comparison of the MTL-CN to baseline on different datasets	108
6.1	Statistics of Datasets	116
6.2	Results of our proposed models with comparison of SemEval (SE) 2018 Task 3 dataset	130
7.1	Distribution of labels combinations in SAFH data.	138
7.2	Precision	147
7.3	Recall	147
7.4	F_1 -score	148
7.5	Accuracy	148
7.6	BiLSTM + FastText with self attention on SAFH COVID-19 dataset	149
7.7	BiLSTM with FastText on SAFH COVID-19 dataset	149
7.8	BiLSTM with BERT + FastText on SAFH dataset	149
7.9	Accuracy of the proposed models on SAFH dataset	150

Nomenclature

Acronyms / Abbreviations

BERT Bidirectional Encoder Representations from Transformers

CNN Convolutional Neural Network

DA Domain Adaptation

DL Deep Learning

ELMo Embeddings from Language Model

HOF Hate Speech and Offensive

IGE-2019 Indian general election 2019

KNN K-Nearest Neighbour

LR Logistic Regression

LSTM Long short-term memory

ML Machine Learning

MNB Multinomial Naive Bayes

MTL-CN Multi-task Learning With the Convolution Network

MTL Multi-task Learning

PEI-2019 Parliamentary Election of India 2019

RF Random forest

SGD Stochastic Gradient Descent

SM Social Media

SVM Support Vector Machine

TL Transfer Learning