

# **UNDERSTANDING NETWORK STRUCTURE, DIFFUSION PROCESS, OPINION LEADERS, AND SENTIMENTS IN SOCIAL NETWORKS: A CASE STUDY IN HEALTHCARE**

*a Thesis submitted by*

ABHINAV CHOUDHURY

(D15049)

*for the award of the degree of Doctor of Philosophy*



SCHOOL OF COMPUTING & ELECTRICAL ENGINEERING

INDIAN INSTITUTE OF TECHNOLOGY MANDI

January 22<sup>nd</sup>, 2021



## THESIS CERTIFICATE

This is to certify that the work contained in the thesis entitled “Understanding network structure, diffusion process, opinion leaders, and sentiments in social networks: A case study in healthcare” being submitted by Mr. Abhinav Choudhury (Enroll. No: D15049) has been carried out under my supervision. In my opinion, the thesis has reached the standard fulfilling the requirement of regulation of the Ph.D. degree. The results embodied in this thesis have not been submitted elsewhere for the award of any degree or diploma.

A handwritten signature in blue ink, appearing to read "Varun Dutt", is placed over a blue horizontal line.

January 22<sup>nd</sup>, 2021

Dr. Varun Dutt (Supervisor)  
Associate Professor  
School of Computing and Electrical Engineering  
School of Humanities and Social Sciences  
Indian Institute of Technology Mandi  
Kamand, Himachal Pradesh, India  
Email: varun@iitmandi.ac.in

### **Declaration by the Research Scholar**

I hereby declare that the entire work embodied in this Thesis is the result of investigations carried out by me in the *School of Computing and Electrical Engineering*, Indian Institute of Technology Mandi, under the supervision of **Dr. Varun Dutt**, and that it has not been submitted elsewhere for any degree or diploma. In keeping with the general practice, due acknowledgments have been made wherever the work described is based on findings of other investigators.



Place: IIT Mandi, Kamand

Date: January 22<sup>nd</sup>, 2021

Name: Abhinav Choudhury

**Declaration by the Research Advisor**

I hereby certify that the entire work in this Thesis has been carried out by *Abhinav Choudhury*, under my supervision in the *School of Computing and Electrical Engineering*, Indian Institute of Technology Mandi, and that no part of it has been submitted elsewhere for any Degree or Diploma.

Signature:

A handwritten signature in blue ink, appearing to read "Varun Dutt".

Name of the Guide: Dr. Varun Dutt

Date: January 22<sup>nd</sup>, 2021

## ABSTRACT

India has a significant shortage of trained physicians nationwide with a physician to population ratio of only 1:1674. Thus, identifying critical physicians is imperative as the proper diffusion of medical information to these physicians is of utmost importance. To identify critical physicians, one first needs to understand the network structure of the physician's social network as well as the underlying dynamics of adoption. The main contribution of this thesis is to understand how the network structure and the underlying network dynamics affect the diffusion of innovation that takes place inside physician social networks, and to develop novel strategies for identifying critical physicians to accelerate the diffusion process. In the first experiment, a binary approach and a weighted approach was proposed for creation of physician social networks. These approaches relied upon the similarity between physician attributes to assign weights to relationships. Results indicated that the weighted approach had a higher accuracy compared to the binary approach in predicting a physician's medicine adoption.

While network structures play a pivotal role in modeling information diffusion, one also needs to understand the underlying network dynamics. In the second experiment, the effect of network dynamics on the diffusion of innovation was investigated. First, innovation diffusion was analyzed from a time-series perspective. Diffusion (Roger's and Bass model), statistical (seasonal autoregressive integrated moving average), and machine learning (linear regression and gradient boosting regression) models were developed for predicting growth in the number of adopters of pain medications. The best performance was obtained by the time-series models, followed by machine learning and diffusion models. Second, innovation diffusion process inside a physician's social network was analyzed by predicting information cascades using multi-layer perceptron (MLP) and long short-term memory (LSTM) neural

networks. A systematic evaluation of different graph embedding techniques and the effect of embedding dimensions on the prediction of information cascades was also performed. Results indicated that the embedding techniques that preserved both the first-order and second-order proximity performed better in cascade prediction compared to those that only preserved one of the proximity measures. Furthermore, MLPs performed better compared to LSTMs in predicting information cascades.

Next, the third experiment investigates the effect of communication channels like electronic word of mouth (eWOMs) communications (e.g., tweets) on the innovation diffusion process. A sentiment analysis of tweets of followers and non-followers of two rare disease medication manufacturers and their respective medications was performed to understand their effect on the overall perception. Results indicated a counter intuitive finding: there was no significant difference in the average sentiment of followers and non-followers regarding rare disease medications, indicating that followers may not possess positive sentiments.

Finding critical physicians may be summarized as an influence maximization (IM) problem, which aims to select a subset of physicians from an influence social graph, such that the diffusion of information is maximized. In the fourth experiment, a reinforcement learning (RL)-based framework for solving the IM problem was proposed. The RL framework consisted of an edge-based graph neural network (GNN) that generated the node embeddings, which were then fed to a double deep Q-network (DDQN) to learn a Q-function to predict the solution set. The edge based GNN was an ensemble architecture comprising of structure2Vec, multi-headed self-attention, and edge enhanced graph neural network (EGNN). The framework was trained on social graph with 20% of its edges randomly removed and tested on the whole graph. Results revealed that the framework was able to

generalize to an unknown graph and gave a spread difference of 8% compared to a heuristic IM algorithm implemented on the whole graph.

In the fifth experiment, the problem of online influence maximization (OIM) was investigated, where OIM is a variant of the IM problem. The objective of OIM is to identify critical physicians in the absence of influence probabilities in a social graph. A new explore-exploit ensemble approach based on the Exponential-weight, Exploration, and Exploitation using Expert advice (EXP4) algorithm was proposed for solving the OIM problem. Results indicated the ensemble approach performed better compared to other current OIM algorithms.

Lastly, in the sixth experiment, the problem of volume maximization (VM) and different algorithmic solutions to the VM problem were investigated. The objective of the VM problem was to select a set of physicians who are both influential as well as have frequent interactions with patients. For solving the VM problem, these two frameworks were proposed: a reinforcement learning framework that developed Q-learning and SARSA models; and an instance-based learning (IBL) framework that developed IBL models. A greedy based algorithm called weighted-CELF was also proposed for the VM problem. The proposed models were compared with multiple IM algorithms like PMIA. Results revealed that the weighted-greedy and weighted-CELF algorithm gave the best weighted influenced spreads while PMIA gave the best influence spread. In contrast, the RL-frameworks: Q-learning and SARSA gave good weighted influence spread, and influence spread different initial seed set sizes ( $k$ ). This thesis highlights the utility of the proposed approaches in identifying critical physicians and shows the effect of network structure and the underlying dynamics on the diffusion of innovation in a physician's social networks.

**Keywords:** *Social network analysis, Network structure, Influence maximization, Online influence maximization, Healthcare, Reinforcement learning, EXP4, Q-learning, SARSA, Graph neural network, Sentiment analysis, Volume maximization, Instance based learning*

## **Acknowledgments**

It is with immense gratitude and profound thanks that I acknowledge the support of my Ph.D. supervisor, Dr. Varun Dutt, in completing this thesis. I am incredibly grateful for his immense patience while dealing with me. His invaluable guidance, constructive feedback, and constant encouragement are the bases for the success of this research. I want to express my sincere gratitude towards him for welcoming me into his research team, providing me an excellent research atmosphere. His mentorship has significantly influenced, the way I understand and approach research work.

I have great pleasure in acknowledging my gratitude towards the co-chairs of my dissertation committee: Dr. Aditya Nigam, Dr. Sriram Kailasam, and Dr. Devika Sethi for generously offering their time, guidance, and goodwill throughout my Ph.D. duration. I would also like to acknowledge Dr. Samar Agnihotri and Dr. Bharat Singh Rajpurohit for their feedback on my work. Without their constant push, I would not have shaped my thesis in its current form.

I warmly thank my friend and colleague Shruti Kaushik for creating a cordial working environment and supporting me in difficult times. I am also thankful to my fellow lab members (Palvi Aggarwal, Zahid Maqbool, and Akash K Rao) for their invaluable support in the lab.

I am extremely grateful to my parents Dr. Suraja Choudhury and Mr. Ajit Choudhury, for their unconditional support and my brother (Arpan) for being the best brother in the world. Without their love, support, and constant encouragement throughout my life, I would not have completed this thesis or anything else for that matter.

Lastly, I would like to express my gratitude towards Mr. Larry A. Pickett, Jr., and Mr. Sayee Natarajan from RxDataScience Inc. for providing me with the funding for my research work.

ABHINAV CHOUDHURY

## **TABLE OF CONTENTS**

|  |              |
|--|--------------|
| <b>THESIS CERTIFICATE .....</b>                  | <b>ii</b>    |
| <b>DECLARATION BY THE RESEARCH SCHOLAR .....</b> | <b>iii</b>   |
| <b>DECLARATION BY THE RESEARCH ADVISOR.....</b>  | <b>iv</b>    |
| <b>ABSTRACT.....</b>                             | <b>v</b>     |
| <b>ACKNOWLEDGMENTS .....</b>                     | <b>ix</b>    |
| <b>LIST OF TABLES .....</b>                      | <b>xviii</b> |
| <b>LIST OF FIGURES .....</b>                     | <b>xix</b>   |
| <b>Chapter 1 .....</b>                           | <b>1</b>     |
| 1.1 Background.....                              | 7            |
| 1.2 Objective .....                              | 15           |
| 1.3 Contribution of Thesis .....                 | 16           |
| 1.4 Thesis Layout.....                           | 18           |
| <b>Chapter 2 .....</b>                           | <b>20</b>    |
| <b>Computational Background.....</b>             | <b>20</b>    |
| 2.1 Social Network.....                          | 20           |
| 2.1.1 Social Network.....                        | 20           |
| 2.1.2 Social Influence Graph .....               | 20           |
| 2.1.3 Degree Centrality .....                    | 21           |
| 2.1.4 Eigenvector Centrality .....               | 21           |
| 2.1.5 Diffusion Degree.....                      | 22           |
| 2.1.6 Maximum Influence Degree .....             | 22           |
| 2.1.7 Small World Phenomenon .....               | 24           |
| 2.1.8 Scale Free Networks .....                  | 25           |
| 2.2 Diffusion of Innovation.....                 | 26           |
| 2.2.1 Diffusion .....                            | 26           |
| 2.2.2 Theory of Diffusion of Innovation.....     | 26           |
| 2.2.3 Bass Diffusion Model .....                 | 27           |
| 2.3 Influence Maximization.....                  | 27           |
| 2.3.1 Diffusion Model.....                       | 27           |
| 2.3.2 Influence Spread .....                     | 28           |
| 2.3.3 Monotone Function .....                    | 28           |
| 2.3.4 Submodular Function.....                   | 28           |

|   |           |
|---|-----------|
| 2.3.5 Independent Cascade (IC) Diffusion Model .....                | 28        |
| 2.3.6 General Threshold Diffusion Model .....                       | 29        |
| 2.3.7 Influence Maximization .....                                  | 29        |
| 2.3.8 Greedy Algorithm .....  | 30        |
| 2.3.9 Cost-Effective Lazy Forward (CELF) .....                      | 30        |
| 2.3.10 Prefix Excluding Maximum Influence Arborescence (PMIA) ..... | 31        |
| 2.3.11 Two-phase Influence Maximization (TIM) .....                 | 32        |
| 2.3.12 Online Influence Maximization (OIM).....                     | 33        |
| 2.4 Graph Embedding .....   | 33        |
| 2.4.1 Graph Embedding .....   | 33        |
| 2.4.2 DeepWalk .....  | 34        |
| 2.4.3 Node2Vec .....  | 34        |
| 2.4.4 Structured Deep Network Embedding (SDNE) .....                | 36        |
| 2.5 Graph Neural Network.....                                       | 38        |
| 2.5.1 Graph Neural Network (GNN) .....                              | 38        |
| 2.5.2 Structure2Vec .....   | 39        |
| 2.5.3 Edge Enhanced Graph Neural Network (EGNN) .....               | 40        |
| 2.6 Artificial Neural Network.....                                  | 43        |
| 2.6.1 Multi-layer Perceptron (MLP) .....                            | 43        |
| 2.6.2 Long Short-Term Memory (LSTM) .....                           | 44        |
| 2.7 Reinforcement Learning .....                                    | 45        |
| 2.7.1 Q-Learning .....  | 45        |
| 2.7.2 SARSA.....  | 46        |
| 2.7.3 Deep Q-Learning (DQN) .....                                   | 47        |
| 2.7.4 Double Deep Q-Learning (DDQN) .....                           | 48        |
| 2.7.5 Prioritized Experience Replay (PER) .....                     | 49        |
| 2.7.6 Upside-Down Reinforcement Learning (UDRL) .....               | 50        |
| 2.8 Multi-headed Self-attention .....                               | 51        |
| 2.9 Instance-based Learning Theory (IBLT) .....                     | 52        |
| 2.9.1 Instance-based Learning Framework .....                       | 53        |
| <b>Chapter 3 .....</b>  | <b>55</b> |
| <b>Physician Network Creation .....</b>                             | <b>55</b> |
| 3.1 Introduction.....   | 55        |

|  |           |
|--|-----------|
| <b>3.2 Methodology.....</b>  | <b>59</b> |
| 3.2.1 Data .....   | 59        |
| 3.2.2 Data Preprocessing.....  | 61        |
| 3.2.3 Binary physician social network .....  | 62        |
| 3.2.4 Social Network Creation.....   | 62        |
| 3.2.5 Weighted Physician Social Network with Equal Weights.....                                      | 63        |
| 3.2.6 Social Network Creation.....   | 63        |
| 3.2.7 Weighted Physician Social Network with Unequal Weights .....                                   | 64        |
| 3.2.8 Data Preprocessing.....  | 64        |
| 3.2.9 Social Network Creation.....   | 65        |
| 3.2.10 Influence Probabilities .....   | 71        |
| 3.2.11 Bernoulli distribution .....  | 71        |
| 3.2.12 Jaccard Index .....   | 72        |
| 3.2.13 Partial Credit .....  | 72        |
| 3.2.14 The Diffusion Process using the General Threshold Model.....                                  | 74        |
| 3.2.15 Model Evaluation.....   | 75        |
| <b>3.3 Expectation .....</b>   | <b>77</b> |
| <b>3.4 Results.....</b>  | <b>77</b> |
| 3.4.1 Analysis of the Network .....  | 78        |
| 3.4.2 Mapping the Social Network with Directed Diffusion Graph .....                                 | 80        |
| 3.4.3 Evaluating the Diffusion inside the Social Network Using the Different Probability Models..... | 89        |
| <b>3.5 Discussion and Conclusions .....</b>  | <b>91</b> |
| <b>Chapter 4 .....</b>   | <b>97</b> |
| <b>Understanding the Dynamics of Diffusion in Healthcare .....</b>                                   | <b>97</b> |
| <b>4.1 Understanding Diffusion Using Diffusion of Innovation Models .....</b>                        | <b>98</b> |
| 4.1.1 Introduction.....  | 98        |
| 4.1.2 Methodology .....  | 99        |
| 4.1.2.1 Rogers' Diffusion of Innovation Model .....  | 99        |
| 4.1.2.2 Bass Model.....  | 100       |
| 4.1.2.3 Time series forecasting using ARIMA models.....  | 100       |
| 4.1.2.4 Linear Regression .....  | 103       |
| 4.1.2.5 Gradient Boosting Regression .....   | 104       |

|  |            |
|--|------------|
| 4.1.3 Data .....   | 104        |
| 4.1.4 Time-series Analysis.....  | 106        |
| 4.1.4.1 Stationarity .....   | 106        |
| 4.1.5 Model Calibration .....  | 108        |
| 4.1.5.1 Rogers' Diffusion of Innovation Model .....                              | 108        |
| 4.1.5.2 Bass Model.....  | 109        |
| 4.1.5.3 Time-series Analysis.....  | 109        |
| 4.1.5.4 Linear Regression .....  | 110        |
| 4.1.5.5 Gradient Boosting Regression .....                                       | 113        |
| 4.1.6 Evaluation Metrics .....   | 114        |
| 4.1.7 Expectation .....  | 115        |
| 4.1.8 Results.....   | 115        |
| Rogers' Diffusion of Innovation Model .....                                      | 115        |
| Bass Model.....  | 117        |
| Time series forecasting using ARIMA .....  | 118        |
| Linear Regression .....  | 120        |
| Gradient Boosting Regression .....   | 123        |
| 4.1.9 Model Comparison.....  | 125        |
| 4.1.10 Discussion and Conclusions .....  | 127        |
| <b>4.2 Evaluating Cascade Prediction via Different Embedding Techniques.....</b> | <b>130</b> |
| 4.2.1 Introduction.....  | 130        |
| 4.2.2 Methodology .....  | 133        |
| 4.2.3 Dataset.....   | 134        |
| 4.2.4 Model Calibration for Graph Embedding Techniques.....                      | 135        |
| 4.2.5 Cascade prediction .....   | 135        |
| 4.2.6 Model Calibration for Neural Network Models.....                           | 136        |
| 4.2.7 Evaluation Metrics .....   | 137        |
| 4.2.8 Expectation .....  | 138        |
| 4.2.9 Results.....   | 138        |
| Cascade Prediction using MLP.....  | 138        |
| Cascade Prediction using LSTM .....  | 139        |
| 4.2.10 Model Comparison.....   | 141        |
| 4.2.11 Discussion and Conclusions .....  | 143        |

|   |            |
|---|------------|
| <b>Chapter 5 .....</b>  | <b>149</b> |
| <b>Influence of Followers on Twitter Sentiments about Rare Disease Medications.....</b> | <b>149</b> |
| 5.1 Introduction.....   | 149        |
| 5.2 Methodology.....  | 150        |
| 5.2.1 Data .....  | 150        |
| 5.2.2 Preprocessing of the tweets.....  | 151        |
| 5.2.3 Sentiment Analysis .....  | 151        |
| 5.3 Expectation .....   | 152        |
| 5.4 Results.....  | 152        |
| 5.4.1 Sentiment Analysis .....  | 152        |
| 5.4.2 Sentiment Analysis using positive and negative words.....                         | 153        |
| 5.5 Discussion and Conclusion.....  | 157        |
| <b>Chapter 6 .....</b>  | <b>160</b> |
| <b>Influence Maximization using Reinforcement Learning.....</b>                         | <b>160</b> |
| 6.1 Introduction.....   | 160        |
| 6.2 Methodology.....  | 162        |
| 6.2.1 Problem Formulation: IM .....   | 162        |
| 6.2.2 Graph Embedding .....   | 163        |
| 6.2.3 Network Architecture.....   | 164        |
| 6.2.4 Reinforcement Learning Formulation.....   | 166        |
| 6.2.5 Sparse Reward .....   | 167        |
| 6.2.6 Learning Algorithm .....  | 169        |
| 6.2.7 Upside-Down Reinforcement Learning Formulation (UDRL).....                        | 172        |
| 6.2.8 Data .....  | 175        |
| 6.2.9 Model Calibration .....   | 176        |
| 6.3 Evaluation Metrics .....  | 178        |
| 6.4 Expectation .....   | 178        |
| 6.5 Results.....  | 179        |
| 6.5.1 Dense reward .....  | 179        |
| Results on the wiki-Vote dataset.....   | 179        |
| Results on the fb-pages dataset.....  | 180        |
| 6.5.2 Sparse reward.....  | 180        |
| Results on the wiki-Vote dataset.....   | 180        |

|  |            |
|--|------------|
| Results on the fb-pages dataset .....  | 181        |
| 6.6 Discussion and Conclusion.....   | <b>182</b> |
| <b>Chapter 7 .....</b>   | <b>184</b> |
| <b>Online Influence Maximization using Experts.....</b>                            | <b>184</b> |
| 7.1 Introduction.....  | <b>184</b> |
| 7.2 Methodology.....   | <b>190</b> |
| 7.2.1 Experts .....  | 191        |
| 7.2.2 Ensemble Approach for Seed Selection.....                                    | 191        |
| 7.2.3 Real-World Feedback .....  | 195        |
| 7.2.4 Updating the Influence Probabilities in the Uncertain Influence Graph.....   | 195        |
| 7.2.5 Local Update .....   | 197        |
| 7.2.6 Global Update .....  | 197        |
| 7.2.7 Kalman Update .....  | 198        |
| 7.3 Data.....  | <b>199</b> |
| 7.4 Expectation .....  | <b>199</b> |
| 7.5 Results.....   | <b>200</b> |
| 7.5.1 Results on NetHEPT .....   | 200        |
| 7.5.2 Updating the Uncertain Influence Graph .....                                 | 203        |
| 7.5.3 Results on PhysicianSN .....   | 208        |
| 7.5.4 Updating the Uncertain Influence Graph .....                                 | 209        |
| 7.6 Discussion and Conclusions .....   | <b>211</b> |
| <b>Chapter 8 .....</b>   | <b>216</b> |
| <b>Beyond Influence Maximization: Volume Maximization in Social Networks .....</b> | <b>216</b> |
| 8.1 Introduction.....  | <b>216</b> |
| 8.2 Method.....  | <b>219</b> |
| 8.2.1 Volume Maximization .....  | 219        |
| 8.2.2 Reinforcement Learning Framework .....                                       | 220        |
| 8.2.3 Instance-based Learning Framework .....                                      | 221        |
| 8.2.4 Weighted-CELF Algorithm .....  | 221        |
| 8.2.5 Baseline Algorithm for comparison.....                                       | 221        |
| 8.2.6 Model Calibration .....  | 222        |
| 8.2.7 Data .....   | 222        |
| 8.3 Expectation .....  | <b>223</b> |

|   |            |
|---|------------|
| 8.4 Result .....                              | 223        |
| 8.5 Discussion and Conclusions .....          | 227        |
| <b>Chapter 9 .....</b>                        | <b>230</b> |
| <b>Conclusion and Future Scope .....</b>      | <b>230</b> |
| 9.1 General Discussion .....                  | 230        |
| 9.2 Research Implications.....                | 238        |
| 9.3 Limitations.....                          | 240        |
| 9.4 Future Scope .....                        | 243        |
| Reference .....                               | 245        |
| <b>List of Publications from Thesis .....</b> | <b>259</b> |
| <b>List of Other Publications .....</b>       | <b>261</b> |