

# **Classification of Comments Supporting the Farmer Protest**

*Project report submitted in partial fulfillment*

*of the requirements for the degree of*

***Master of Technology***

*in*

***Computer Science and Engineering***

***(Specialization: Information Security)***

*by*

***Ajay Biswas***

(Roll Number: 220CS2184)

*based on research carried out*

*under the supervision of*

***Prof. Tapas Kumar Mishra***



October, 2021

Department of Computer Science and Engineering  
**National Institute of Technology Rourkela**

# Contents

<b>List of Figures</b>	<b>iii</b>
<b>List of Tables</b>	<b>iv</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Introduction . . . . .	1
1.2 Applications of Active Learning . . . . .	1
1.3 Motivation . . . . .	2
1.4 Objectives . . . . .	3
<b>2 Literature Survey</b>	<b>4</b>
2.1 Overview . . . . .	4
2.2 Related Works . . . . .	4
<b>3 Classification of Comments Supporting the Farmer Protest</b>	<b>6</b>
3.1 Overview . . . . .	6
3.2 Dataset . . . . .	6
3.2.1 Video and comment statistics . . . . .	6
3.2.2 User level analysis . . . . .	7
<b>4 Conclusion</b>	<b>8</b>
<b>References</b>	<b>9</b>
<b>Index</b>	<b>10</b>

# List of Figures

1.1	Pool-based Active Learning . . . . .	2
-----	--------------------------------------	---

# List of Tables

3.1	Most Frequently Occurring Word in Cloud . . . . .	6
3.2	Manually Selected Seed Word . . . . .	7
3.3	Sentiment of Media Houses Towards Farmers . . . . .	7

# Chapter 1

## Introduction

### 1.1 Introduction

The Farm Bills, or the Indian agriculture acts of 2020, are three acts initiated by Parliament of India during September 2020. The three farm acts are as follows: "Farmers' Produce Trade and Commerce (Promotion and Facilitation) Act, 2020; Farmers (Empowerment and Protection) Agreement on Price Assurance and Farm Services Act, 2020; and Essential Commodities (Amendment) Act, 2020". Although the bills were supposed to be beneficial for the farmers, they led to a mass protest, which gained momentum in September 2020 [1]. The Protest is not limited to grounds, but also spread across various social media websites like YouTube, Reddit, Facebook, Instagram and Twitter. Our main goal is not to argue who is right or wrong but to identify the comments that are in favor of the Indian Farmers' protest.

Our proposed approach is inspired from the work done in [2]. This paper proposed an Active Learning (AL) based classifier that can classify comments supporting the Rohingyas. Further discussions related to this paper will be dealt in the literature survey section. Active learning is a technique for reducing manual annotation effort during training phase of machine learning. The annotation is done by a human (called oracle) which helps AL systems to achieve high accuracy with few labelled instances. For problems having large collection of unlabeled data, pool-based sampling is used [3]. Figure 1.1 shows working of pool-based AL. AL is very useful in classifying comments which involves a person's opinion, belief or political interest, as it's too complicated to be dealt with plain unsupervised machine learning. Also, there are numerous challenges to be dealt with before any classification could take place. Some of the challenges are (i) Dealing with multiple languages having different levels of grammatical accuracy, (ii) Un-structured data, (iii) ambiguous sentences, (iv) Unrelated comments, etc.

### 1.2 Applications of Active Learning

Active learning is becoming a burning topic in the field of machine learning due to its high performance and wide range of uses. Some of the areas where active learning can be useful

are as follows:

- *Speech Recognition.* One by one labeling speech utterances can be very challenging [4] as well as time consuming. As speech may have several languages or multiple dialects, trained linguists are needed for this purpose.
- *Information Extraction.* To extract high quality information, hours of manual labor is required. For highly specialized job, professional are required like Genomic information retrieval requires Phd-level candidates [5].

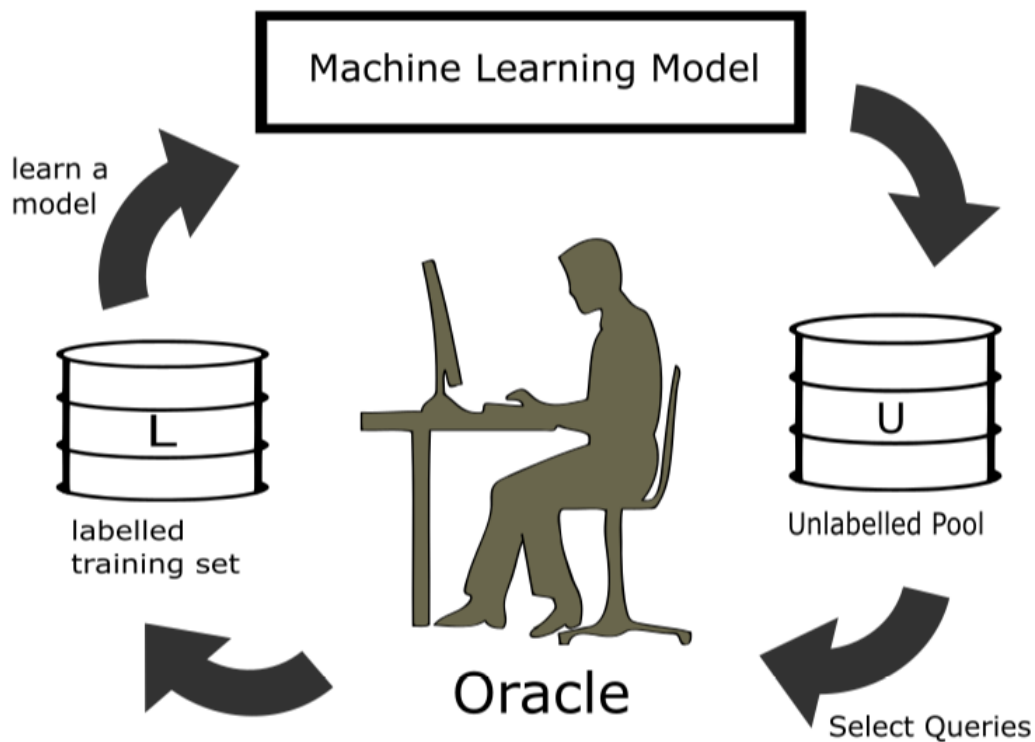


Figure 1.1: Pool-based Active Learning

### 1.3 Motivation

Voting is an important part of democracy as it allows people choose their representative. Just like citizens choose their representative, they also agree or disagree with government policies. By analyzing comments from social media sites, we can understand the sentiment of the country. The motivation behind the selection of Active Learning for comment classification is its high performance when problem is specific.

## 1.4 Objectives

- i. Create a dataset by parsing comments from social media sites like YouTube.
- ii. Filter the dataset by removing unwanted characters/words.
- iii. Construct a seed set containing positive and negative comments.
- iv. Expand the seed set by randomly sampling comments from the unseen corpus.
- v. Obtain real valued embeddings for the comments.
- vi. Find the nearest neighbors (NN) of the seed set and include them in the corpus.
- vi. Expand using minority-class certainty sampling (if any).
- vii. Perform final expansion using uncertainty sampling.

The organization of this report is as follows: Chapter 1 provides brief introduction briefly describes about the ongoing farmer protest and how active learning can be used to classify comments supporting the protest. It also outlines the motivation and objective of this work; Chapter 2 provides a brief summary on the related works are summarized in this chapter; Chapter 3 describes our proposed work. It describes how we built the dataset containing comments fetched from YouTube. Also it briefly summarizes how we are going to apply active learning to classify them; and finally, Chapter 4 presents the conclusion our this work, the limitations and the areas that can be improve in future work.

## Chapter 2

# Literature Survey

### 2.1 Overview

This chapter contains a brief survey of research work done in the field of Sentence Classification and Active Learning.

### 2.2 Related Works

Previously various research were conducted in the field of Active Learning and Sentence Classification. We are focusing on those researches which are related to our work.

- i. Chen et al. [6] proposed a state-of-the-art result on the Stanford Natural Language Inference Dataset using Long Short-Term Memory (LSTM). They employed Bi-directional LSTM (BiLSTM) as one of the building blocks. Later it is used to perform inference composition to construct the final prediction.
- ii. Kim et al. [7] propose a densely-connected co-attentive recurrent neural network to find semantic relation between sentences. To overcome the problem of ever-increasing size of feature vector due to densely connected networks, they also have propose an autoencoder after dense concatenation.
- iii. Nguyen and Smeulders [8] incorporated clustering into active learning. The algorithm first constructs a classifier on the set of the cluster representatives, and then with the help of a local noise model, it passes the classification decision to the other samples. The model allows selecting the most representative samples as well as avoids labelling samples in the same cluster. The paper focuses on discriminative models including logistic regression and Support Vector Machines (SVM) which are less sensitive to training data and hence, good for active learning.
- iv. Ru et al. [9] propose adversarial uncertainty sampling in discrete space (AUSDS) which retrieves informative unlabeled samples more efficiently and is 10x faster when compared to typical uncertainty sampling method for active learning.



- v. Schumann and Rehbein [10] showed that it is possible to use Membership Query Synthesis [5] for generating AL queries for natural language processing, using Variational Autoencoders for query generation, and provides competitive performance to pool-based AL strategies while substantially reducing annotation time.
- vi. Ren et al. [11] propose a novel fake news detection framework "Adversarial Active Learning based Heterogeneous Graph Neural Network for Fake News Detection (AA-HGNN)", which employs a novel hierarchical attention mechanism to perform node representation learning in the HIN. In this paper, the authors model the news content and related entities as a News-HIN. The AA-HGNN utilizes both structural information as well as News-HIN to identify fake news.
- vii. Xiang et al. [12] proposes a novel approach which exploits linguistic regularities in profane language via statistical topic modeling on a huge Twitter corpus, and detects offensive tweets using these automatically generated features. This approach works with various Machine Learning models such as J48 decision tree learning, SVM, logistic regression (LR) and random forest (RF).
- viii. Watanabe et al. [13] proposes a pragmatic approach to collect hateful speech. The proposed approach uses unigram and patterns that are automatically collected from training dataset. Accuracy of 87.4% was achieved on detecting whether a tweet is offensive or not (binary classification) and 78.4% accuracy when detecting a tweet is hateful, offensive, or clean (ternary classification).
- ix. Palakodety et al. [2] proposes a classifier which can classify comments supporting the Rohingyas. This is done by building a corpus from YouTube comments and applying multiple AL strategies based on nearest-neighbors in the comment-embedding space.

## Chapter 3

# Classification of Comments Supporting the Farmer Protest

### 3.1 Overview

This chapter provides the proposed work done so far in classification of comments supporting the Indian farmers' protest.

### 3.2 Dataset

The dataset was constructed by parsing comments from 863 unique videos from YouTube. Total 45,376 comments were fetched out of which 40,791 comments were in pure English words. Although most of the comments uses English alphabets, it has been observed that, they are mostly Hinglish (Hindi written using English characters), and hence, for now we are taking them in consideration.

#### 3.2.1 Video and comment statistics

On analyzing unigram [kisan] and bigram [kisan nahi], clouds of 788 and 45 unique words were formed. Table 3.1 shows the most frequent words of these two clouds.

Table 3.1: Most Frequently Occurring Word in Cloud

kisan	Frequency	kisan nahi	Frequency
ekta	844	dalal	5
jai	228	aatankwadi	3
majdoor	183	khalistani	3

To start the active learning process, we will require some manual seed words. Table 3.2 shows the manually selected seed words.

Table 3.2: Manually Selected Seed Word

Positive	Negative
support, love, power, zindabad, care, jay jawan jay kisan, proud	evil, shame, dacoit, illegal, riot, khalistani, selfish, greedy, talibani

### 3.2.2 User level analysis

After analyzing various comments and YouTube videos, we have observed that there are two groups of channels that either make videos in favor of farmers or against of farmers. Table 3.3 shows sentiment of media houses towards farmers.

Table 3.3: Sentiment of Media Houses Towards Farmers

For	Against
NDTV, BBC News, Al Jazeera	Zee News, Times Now, Republic TV, Sudharshan News, ABP News, News Nation

## **Chapter 4**

# **Conclusion**

In this report we provided active learning based approach to classify comments supporting the farmers protest. As we are working with our own dataset, it gives much more flexibility and provides latest results. The statistical analysis showed us the sentiments of people towards farmers.

## **Scope for Further Research**

The report is limited to Data pre-processing. As for future work, we will be completing the remaining objectives.

# References

- [1] Shankar, A. Indian agriculture farm acts: 2020. *International Journal of Modern Agriculture* **10**, 2907–2914 (2021).
- [2] Palakodety, S., KhudaBukhsh, A. R. & Carbonell, J. G. Voice for the voiceless: Active sampling to detect comments supporting the rohingyas. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, 454–462 (2020).
- [3] Lewis, D. D. & Gale, W. A. A sequential algorithm for training text classifiers. In *SIGIR'94*, 3–12 (Springer, 1994).
- [4] Zhu, X. J. Semi-supervised learning literature survey (2005).
- [5] Settles, B., Craven, M. & Friedland, L. Active learning with real annotation costs. In *Proceedings of the NIPS workshop on cost-sensitive learning*, vol. 1 (Vancouver, CA:, 2008).
- [6] Chen, Q. *et al.* Enhanced lstm for natural language inference. *arXiv preprint arXiv:1609.06038* (2016).
- [7] Kim, S., Kang, I. & Kwak, N. Semantic sentence matching with densely-connected recurrent and co-attentive information. In *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, 6586–6593 (2019).
- [8] Nguyen, H. T. & Smeulders, A. Active learning using pre-clustering. In *Proceedings of the twenty-first international conference on Machine learning*, 79 (2004).
- [9] Ru, D. *et al.* Active sentence learning by adversarial uncertainty sampling in discrete space. *arXiv preprint arXiv:2004.08046* (2020).
- [10] Schumann, R. & Rehbein, I. Active learning via membership query synthesis for semi-supervised sentence classification. In *Proceedings of the 23rd conference on computational natural language learning (CoNLL)*, 472–481 (2019).
- [11] Ren, Y., Wang, B., Zhang, J. & Chang, Y. Adversarial active learning based heterogeneous graph neural network for fake news detection. In *2020 IEEE International Conference on Data Mining (ICDM)*, 452–461 (IEEE, 2020).
- [12] Xiang, G., Fan, B., Wang, L., Hong, J. & Rose, C. Detecting offensive tweets via topical feature discovery over a large scale twitter corpus. In *Proceedings of the 21st ACM international conference on Information and knowledge management*, 1980–1984 (2012).
- [13] Watanabe, H., Bouazizi, M. & Ohtsuki, T. Hate speech on twitter: A pragmatic approach to collect hateful and offensive expressions and perform hate speech detection. *IEEE access* **6**, 13825–13835 (2018).

# Index

Applications of Active  
Learning, 1

Dataset, 6

Introduction, 1

Motivation, 2

Objectives, 3

Overview, 4, 6

Related Works, 4

User level analysis, 7

Video and comment  
statistics, 6