

Sarcasm Detection System for Hinglish Language (SDSHL)

A dissertation is presented towards the
fulfilment of the requirements for M.Sc. in Data Science

Hari Thapliyal
M.Sc in Data Science

Faculty Advisor: Dr. Anil Vuppula

Liverpool John Moore University, UK
Date:

Acknowledgements

Abstract

Hinglish is third¹ most spoken language on the planet. (Wikipage, n.d.) states that 65% of Indian population is under 35 years age. Several disruptions like low cost mobile phone, extremely cheap data, digital India initiatives by government of India has cause huge surge in Hinglish language content. This content is available in audio, video, images, and text format. We can find Hinglish content in comment box of product, new articles, service feedback, WhatsApp, social media like YouTube, Facebook, twitter etc. To engage with consumer, it is extremely important to analyse the sentiments, but to perform sentiment analysis it is not possible to read every comment or feedback using human eyes. With increasing number of education and sophistication people in Indian society it is evident that people do not say negative things directly even when they want to say. Educated and advance mind is more diplomatic than less educated or village people who are not exposed enough to the world. Due to this reason people use more sarcastic language, they say negative things in positive words. Thus, it becomes necessary to identify the true sentiments in this kind of conversation. In this paper we are demonstrating a system which can help in automatic sarcasm detection. In this work we are extracting text from Hindi twitter handles and Hindi blogs. We take all the tweets which are written in Roman or Devanagari scripts, but words can be from any Indian language or English. Because the text written can also have Roman letters therefore, we are converting that text to Devanagari script. We are performing series of activities to clean the text; so that it can be used for ML work. We know there are not enough good size corpus for Hindi language therefore we will use 2-3 the best available Hindi language corpus for embedding purpose. Not much work has been done in Hindi Language sarcasm detection therefore we do not know which algorithm will give good results. To address that we are going to use SVM, Stochastic Gradient Descent, Adaboost, Random Forest, Naïve Bayesian, Logistic Regression and CNN algorithms to develop different models. We will also deploy Transfer Learning technique for our work. To measure the performance of models we will use AUC, F1 score, Accuracy, Recall, Precision.

¹ https://en.wikipedia.org/wiki/List_of_languages_by_total_number_of_speakers (Accessed on 27-Aug-20)

Table of Contents

| | |
|---|------------|
| ACKNOWLEDGEMENTS | II |
| ABSTRACT | III |
| LIST OF FIGURES | VI |
| LIST OF ABBREVIATIONS | VII |
| CHAPTER 1: INTRODUCTION | 1 |
| 1.1. BACKGROUND OF THE STUDY | 2 |
| 1.1.1. WHAT IS HINGLISH? | 2 |
| 1.1.2. ORIGIN OF HINGLISH | 3 |
| 1.1.3. WHAT IS SARCASM? | 4 |
| 1.1.4. WHY SARCASM DETECTION IS CRITICAL? | 7 |
| 1.1.5. WHY SARCASM DETECTION IS CRITICAL IN ELECTRONIC MEDIA? | 8 |
| 1.1.6. SARCASM DETECTION IN HINGLISH | 9 |
| 1.1.7. CHALLENGES IN PROCESSING HINGLISH LANGUAGE | 9 |
| 1.1.8. COMMON CHALLENGES IN SARCASM DETECTION | 12 |
| 1.1.9. CONTEXT UNDERSTANDING A CHALLENGE IN SARCASM DETECTION | 13 |
| 1.1.10. CHALLENGES IN SARCASM DETECTION IN HINGLISH | 14 |
| 1.1.11. DEGREE OF SARCASM | 16 |
| 1.1.12. POSITIVE SIDE OF HINGLISH | 16 |
| 1.2. PROBLEM STATEMENT | 17 |
| 1.3. AIM AND OBJECTIVES | 17 |
| 1.4. RESEARCH QUESTIONS | 18 |
| 1.5. SCOPE OF THE STUDY | 18 |
| 1.6. SIGNIFICANCE OF THE STUDY | 19 |
| 1.6.1. APPLICATION OF SARCASM DETECTION SYSTEM | 19 |
| 1.6.2. MOTIVATION FROM SELECTED DOMAINS | 20 |
| 1.7. STRUCTURE OF THE STUDY | 21 |
| CHAPTER 2: LITERATURE REVIEW | 23 |
| 2.1. SARCASM DETECTION SYSTEMS (SDS) | 23 |
| 2.2. HISTORY OF SARCASM DETECTION SYSTEMS | 25 |
| 2.3. GENERIC TEXT-BASED SARCASM DETECTION SYSTEM (GTSDS) | 25 |
| 2.4. FEATURE ENGINEERING IN SARCASM DETECTION SYSTEMS | 25 |
| 2.5. APPROACHES TO DEVELOP SDS | 26 |
| 2.5.1. RULE BASED APPROACHES | 27 |
| 2.5.2. CLASSICAL MACHINE LEARNING WITH EMBEDDING & OTHER FEATURES | 27 |
| 2.5.3. CLASSICAL MACHINE LEARNING WITH EMBEDDING | 28 |
| 2.5.4. CLASSICAL MACHINE LEARNING WITHOUT EMBEDDING | 28 |
| 2.5.5. DEEP LEARNING WITH EMBEDDING & OTHER FEATURES | 28 |
| 2.5.6. DEEP LEARNING WITH WORD EMBEDDING | 29 |
| 2.5.7. TRANSFORMER BASED | 29 |
| 2.5.8. TRANSFER LEARNING APPROACHES | 30 |
| 2.5.9. DEEP LEARNING WITHOUT EMBEDDING FEATURES | 31 |
| 2.6. APPROACHES TO HANDLE KEY CHALLENGES IN SARCASM DETECTION | 31 |
| 2.6.1. HANDLING FIGURATIVE LANGUAGES | 31 |
| 2.6.2. HANDLING LIMITED DATA IN SARCASM | 32 |
| 2.6.3. HANDLING OUT OF VOCABULARY ISSUES (OOV) | 32 |

| | | |
|--|---|-----------|
| 2.7. | EMBEDDING | 33 |
| 2.7.1. | ABSOLUTE EMBEDDING | 34 |
| 2.7.2. | CONTEXTUAL EMBEDDING USING FULL WORD | 34 |
| 2.7.3. | CONTEXTUAL EMBEDDING USING SUBWORDS | 35 |
| 2.8. | TYPES OF SARCASM DETECTION SYSTEM | 35 |
| 2.9. | SUMMARY | 36 |
| CHAPTER 3: RESEARCH METHODOLOGY | | 38 |
| 3.1. | DATASET | 38 |
| 3.1.1. | ABOUT DATASET | 38 |
| 3.1.2. | DATA SOURCING | 38 |
| 3.1.3. | DATASET CLEANING | 38 |
| 3.1.4. | SENTENCE LABELLING | 39 |
| 3.1.5. | LANGUAGE TREATMENT OF WORDS | 39 |
| 3.1.6. | DATASET STRUCTURE | 40 |
| 3.1.7. | FEATURE ENGINEERING | 40 |
| 3.1.8. | FINAL DATASET WITH ALL FEATURES | 41 |
| 3.2. | OVERVIEW OF OUR APPROACH | 41 |
| 3.3. | MODEL BUILDING | 42 |
| 3.3.1. | TEST-TRAIN SPLIT | 42 |
| 3.3.2. | HANDLING SMALL DATASET SIZE | 42 |
| 3.3.3. | ALGORITHMS, ARCHITECTURE FOR MODELING | 42 |
| 3.4. | EVALUATION METRICS & REPORTING | 43 |
| 3.4.1. | EVALUATION METRICS & REPORTING | 43 |
| 3.4.2. | REPORTING EXPERIMENTS FORMAT | 44 |
| 3.4.3. | RESULT COMPARISON FORMAT | 45 |
| 3.5. | DEVELOPMENT TOOLS | 45 |
| 3.6. | SUMMARY | 46 |
| REFERENCES | | 48 |
| APPENDIX A: LIST OF OTHER DOCUMENTS | | 51 |
| APPENDIX B: RESEARCH PLAN | | 52 |
| APPENDIX C: SARCASM DETECTION SYSTEMS RESULTS OF PAST WORK. | | 53 |
| APPENDIX D: RESEARCH PROPOSAL | | 1 |

LIST OF FIGURES

| | |
|---|----|
| Figure 1: Evolution of Hinglish..... | 4 |
| Figure 2: Sarcasm & Satire Relationship | 6 |
| Figure 3 : Steps to Creating Dataset | 40 |
| Figure 4: Steps for Labelling Sentences | 40 |
| Figure 5: Overall Approach | 42 |

LIST OF ABBREVIATIONS

1. ACC - Accuracy
2. AFINN - Affective dictionary by Finn ° Arup Nielsen (word with polarity between -5,5)
3. ANEW - Affective Norms for English Word (Emotional Rating)
4. BN - Bayesian Network
5. Chi - χ Square Test
6. CNN - Convolutional Neural Network
7. CORR - Correlation
8. DAL - Dictionary of Affective Language with degree of these three dimensions Activation, Imagery and Pleasantness .
9. DT - Decision Tree
10. EmoLex - A map between words 8 emotions sadness, joy, disgust, anger, fear, surprise, trust, anticipation
11. EmoSN - EmoSenticNet (IR assigns WordNet Affect emotion labels to SenticNet concepts)
12. EWN - The EffectWordNet lexicon (sense-level lexicon created on the basis of WordNet)
13. GBC - Gradient Boost Classifier
14. GI - The Harvard General Inquirer (Words are labelled with a total of 182 dictionary categories and subcategories + Positive Negative)
15. GR - Gain Ratio
16. GRU - Gated Recurrent Unit
17. HL - The Hu-Liu's lexicon (List of Negative and Positive words)
18. IG - Information Gain
19. KNN - K Nearest Neighbour
20. LFS- Linguistic Features of the Sentence
21. LIWC - Linguistic Inquiry and Word Counts dictionary (psycho-linguistic features in texts.)
22. LMT - Logistic Model Tree
23. LR - Logistic Regression
24. LSTM - Long Short Term Memory
25. MLP - Multilevel Perceptron

27. PMI- Pointwise mutual information
28. POS - Part of Speech
29. RBF - Radial Basis Function
30. RF - Random Forest
31. SCUBA - Sarcasm Classification Using a Behavioral modeling Approach
32. SDS - Sarcasm Detection System
33. SDSHL - Sarcasm Detection System for Hinglish Language
34. SGD - Stochastic Gradient Descent
35. SN - SenticNet
36. SS- SentiSense (It attaches emotional meanings to concepts from the WordNet lexical)
37. SVM - Support Vector Machine
38. SWN - SentiWordNet (word along with POS and sentiment score between 0,1)
39. TIM - Topic-Irony Model
40. TL - Transfer Learning
41. TSTAT - T Statistical Test

CHAPTER 1: INTRODUCTION

Mobile phones came to India in 1995² and Internet was launched in India by VSNL in 1995³. Initially the cost of the technology was remarkably high⁴, so it was available only to business class, research labs, high level bureaucrats and politicians. With the increase of literacy and decreasing cost of internet services and mobile phone device internet, it is so common that people started thinking that Internet is our fundamental right. As per the World Economic Forum (WEF), in 2019, about 60% of Indian internet users viewed content in vernacular. WEF also says 75% of this 60% is below 35 years of age (Wikipedia, n.d.). According to the same Wikipedia page, by 2030, 1.1 billion Indian will have access to Internet and 80% will access the content on mobile devices. The WEF also estimated that 80% of the users will be consuming content in vernacular languages.

When Government of India is going for full blown Digital India program and bringing every citizen of India on the internet platform for purchase, payment and government fund transfer then how the citizens are going to provide feedback about the services which they use? As of today, it is easier to perform sentiment analysis of the feedback given in English, but feedback given in Hindi is not easy to analyse. It means voice of Hindi speaking people is not being considered for service improvement. Till the time somebody is not too angry and do some crime or come on the road to do Dharana or protest we do not know what is happening and why.

Many Hindi news portals, book, blogs, chat bot/WhatsApp conversations, YouTube channels, Twitter & Facebook pages are full of content in Hinglish language. People openly express themselves online using Hinglish language which is mix of Hindi, English, Urdu and other Indian languages. Volume of the online content is increasing at unprecedented rate and it is responsibility of the government, business community, professionals, NGO and others to

²https://en.wikipedia.org/wiki/Telecommunications_in_India#:~:text=In%20August%201995%2C%20then%20Chief,launched%20in%20Kolkata%20in%202012. (Accessed 24-Jun-20)

³https://en.wikipedia.org/wiki/Internet_in_India#:~:text=The%20first%20publicly%20available%20internet,not%20permitted%20in%20the%20sector. (Accessed 24-Jun-20)

⁴<https://www.news18.com/news/tech/20-years-of-internet-in-india-on-august-15-1995-public-internet-access-was-launched-in-india-1039859.html#:~:text=The%20Gateway%20Internet%20Access%20Service,organisations%20at%209.6%20kbp%20speed.> (Accessed 27-Aug-20)

understand the feeling of public and respond accordingly. But the biggest challenge is how to analyse the content which is written in mix of Indian languages. It is impossible to analyse the Hinglish language text manually or using traditional systems.

This section is organized as 1.1. Background of The Study, 1.1.1. What is Hinglish?, 1.1.2. Origin of Hinglish, 1.1.3. What is Sarcasm?, 1.1.4. Why Sarcasm Detection is Critical?, 1.1.5. Why Sarcasm Detection is Critical in Electronic Media?, 1.1.6. Sarcasm Detection in Hinglish, 1.1.7. Challenges in Processing Hinglish Language, 1.1.8. Common Challenges in Sarcasm Detection, 1.1.9. Context Understanding A Challenge in Sarcasm Detection, 1.1.10. Challenges in Sarcasm Detection in Hinglish, 1.1.11. Degree of Sarcasm, 1.1.12. Positive Side of Hinglish, 1.2. Problem Statement, 1.3. Aim and Objectives, 1.4. Research Questions, 1.5. Scope of The Study, 1.6. Significance of the Study, 1.6.1. Application of Sarcasm Detection System, 1.6.2. Motivation from Selected Domain

1.1. Background of the Study

1.1.1. What is Hinglish?

There was time when Hindi was a language which is used by majority of Hindi speaking people when they are communicating (writing, speaking) with each other. But in 21st century, most of the Hindi speaking population who express themselves on social media use Hinglish language. Hinglish is a new lingo of Hindi speaking population. Hinglish sentences follow Hindi grammar and most of the word are taken from Hindi but there is no hesitation of taking words from other languages like English, Urdu, Punjabi, Marathi etc. Hinglish language spoken by different people have different amount of words from different languages. For example, those people who know Urdu good enough for them Hinglish is mix of Hindi, Urdu, English. Those who know Avadhi for them Hinglish is mix of Hindi, Avadhi, English. Those who know Marathi very well for them Hinglish is mix of Hindi, Marathi, English. Thus, in Hinglish Language we have words from Hindi, English and various other Indian languages and written in Devanagari & Roman together.⁵ (Sinha and Thakur, 2005) Hindi and English language mixed is called Hinglish. Hinglish is not limited to Hindi & English mix but it includes Punjabi, Gujarati, Marathi, Urdu. Phrase construct happens in Roman and Devanagari script.⁶

⁵ Latin is Region and Rome is part of that reason. Over the period of time Roman empire become famous and script was called Roman but Latin is also used simultaneously. <https://www.quora.com/Why-is-the-language-of-the-ancient-Romans-called-Latin-and-not-Roman> (Accessed 28-Jun-20)

⁶ <https://en.wikipedia.org/wiki/Hinglish> (Accessed 24-Jun-20)

1.1.2. Origin of Hinglish

Before Internet Era in India people use to communicate with each other in much cleaner format of the language and there was not much mix of other language or English and for writing Hindi they were using Devanagari script. But, with the penetration of internet in the society a new language started taking shape. Initially when Devanagari keyboards were not available people were using Roman letters to write Hindi email, SMS.

An example of late 20th century text in Hinglish language. “Main is doorbhash ka prayog karna nani janta”. This is Hindi in Roman script. We need to keep in mind that people do not follow any IAST or other map for writing Hinglish letters in Roman. Mobile phone and Internet were available to elite, educated journalist, professionals. They started realising they are typing in Roman but some words in English so translating them and then typing in Roman is painful. So, text became like this “Main is phone ko use karna nahi janta”. Roman script with Hindi and English words.

Over the period of time when Devanagari keyboards were easily available people started using Devanagari keyboards for writing Hindi, but by that time so much English has come in day to day conversation that they felt it is uncomfortable to use Hindi words. So, they write like this. “मैं इस फोन को यूज करना नहीं जानता”. Devanagari script with Hindi and English words. Over the period of time people started realizing it is becoming difficult to know which word is Hindi and which one is English therefore a word which come from English root should be written in Roman and word which are from Hindi root should be written in Devanagari. So, they started writing like this. “मैं इस phone को use करना नहीं जानता”. Devanagari & Roman mixed for Hindi and English words.

Evolution of Hinglish from Hindi

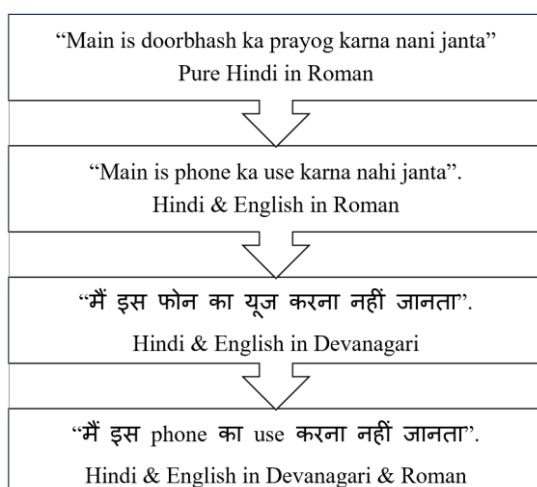


Figure 1: Evolution of Hinglish

Today if you read any Hindi speaker's WhatsApp, twitter or Facebook message you will find they use words from different Indian languages like Urdu, Marathi, Bangla, Punjabi and write either in Devanagari or in Roman. “अमी मौजूलिका. अमी राजा को जरूर मारबो 😊, but why you want to kill him?”. Here Hindi, Bangla, Urdu and English 4 languages used along with emoticon and written in two scripts Devanagari and Roman. This is Hinglish.

Today Hindi social media, Hindi comment boxes of product, Hindi news articles are full of this kind of language, Hinglish. Therefore, this work using Hinglish language is high value from the angle of practical usage.

1.1.3. What is Sarcasm?

Your friend come to you and speak something to you, from the tone of his language, his body language, choice of his words, time and situation he is speaking you realised that the real meaning of what he is saying is completely opposite. It may be easier for you to detect this opposite sense if you are aware about the complete context but if you are not aware about the context then even as intelligent human you may miss the real meaning of what is being said.

For example, you open the door for your friend, and he says wow! your looking handsome in this T-shirt. You know that this is an old T-shirt and many times your friend has seen this. But still not aware of full context, you hesitantly say thank and you invite him inside. After 15

minutes you check yourself in the mirror and realised that you are wearing T-shirt flip side. Now you are embarrassed for your “Thank you” response.

What your friend did was sarcastic remark on your dressing and you being unaware of the full context could not respond properly. In the absence of full context, understanding sarcasm is difficult task and most of the time we take literal meaning of the words or some other time get confused that why someone has made that remarks which was completely out of the context. In English language this type of grammatical construct which has completely opposite meaning than what is said, it called sarcasm.

As per merriam-webster dictionary, sarcasm is⁷

1: a sharp and often satirical or ironic utterance designed to cut or give pain

2a: a mode of satirical wit depending for its effect on bitter, caustic, and often ironic language that is usually directed against an individual

2b: the use or language of sarcasm

In Hindi it has several name and synonyms like कटाक्ष (Kataksha), तंज (Tanja), व्यंग/ व्यङ्ग (Vyanga), टोंट (Tonta)

Ten forms of humour are *irony, satire, sarcasm, overstatement, self-deprecation, teasing, replies to rhetorical question, clever replies to serious statements, and transformations of frozen expressions*. All these are functions of humour and found in the sitcom (situational comedy). What one finds hilarious or pun may be completely opposite to another person in another country or in other situation. Interpretation is filtered by cultural context. (Anggraini, 2014)

⁷ <https://www.merriam-webster.com/dictionary/sarcasm>

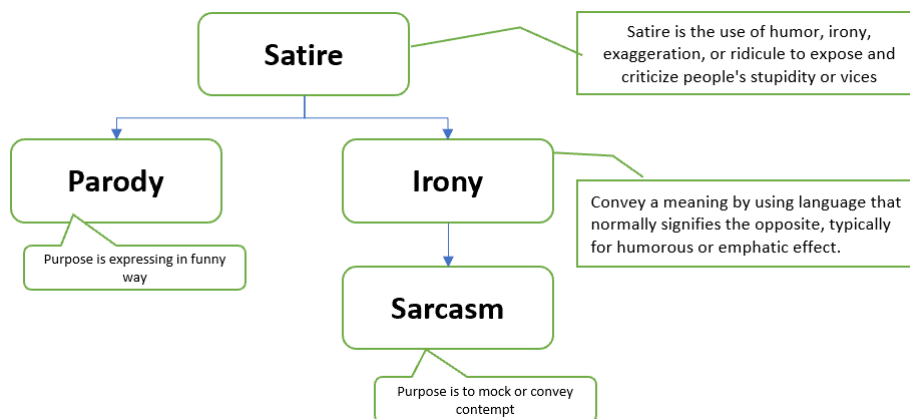


Figure 2: Sarcasm & Satire Relationship

In their work “A Pragmatic Analysis of Humor in Modern Family” (Anggraini, 2014) mentions 11 type of humours. Sarcasm is one type of humour. Let’s understand them with example. We are writing examples in English so that English readers can also understand the important of the this work.

1- Satire:

Rahul: It looks big accident on the road, let’s call police.

Jay: Oh, are you sure? I think police of our state is too busy in catching buffalo of MLAs

2- Irony:

Rahul: Why people steal when there are enough opportunities to work hard and earn.

Jay: Oh, you mean those who steal are doing any less hard work?

3- Sarcasm:

Boss: Why do you work so hard, take leave, enjoy life, have some fun after all life is more than work.

Junior: Oh really! Do you know since last one year we are working in Syria? Come with me tomorrow we will go to have fun in this local Jihadi market.

4- Clever replies to serious statements

Rahul: Jay, why you didn't invite me for your birthday party last night?

Jay: I was not sure you will bring any gift for me.

5- Replies to rhetorical questions

Husband: Today is Sunday, why don't you switch off that alarm?

Wife: So that you get up and help me.

6- Teasing

Boyfriend: Where were you when God was distributing brain?

Girlfriend: I was waiting outside for you.

7- Self-deprecation

"They all left the room when I started singing"

8- Overstatement and Understatement Overstatement

Driver: Please pay me 40 dollar for the service.

Passenger: Because of you I missed my flight, your car had problem. First you pay me \$200 for the missed flight.

9- Double Entendres

Patient: I am having pain in my right hand.

Doctor: But can you raise your right hand?

Patient: You are nice person, why should I raise my hand before you?

10- Transformations of frozen expression Transformations

"Despite of being hare you are not hearing"

11- Pun

Most people don't use God's most valuable gift to them, their mind. The reason for that is they want to make their God happy by returning His gift as is.

In their work "The Differential Role of Ridicule in Sarcasm and Irony" (Lee et al., 2009) says sarcasm and irony are similar because they are both form of reminder yet they are different because sarcasm is about ridiculing a specific person however this is not required in case of irony. Sarcasm plays more important role than irony in ridiculing a specific victim. A speaker is more sarcastic when he reminds the listener somebody else's prediction and less sarcastic when he reminds his own mistake.

In our work we will not pay much attention to these specific aspects of humour. Our intention is to detect a sentence which is not carrying the normal meaning. However, most of the records in our dataset which are labelled as sarcastic are sarcastic but they can have other variation of humour as well.

1.1.4. Why Sarcasm Detection is Critical?

If we do not understand the real intent of the speaker then we cannot respond him properly. Response can be physical action or verbal reply to the speaker or even no action. Sarcasm is like a double edge sword of communication. At one end you can enjoy and another end you can hurt deepest to the opponent. If you do not handle this properly effect can be completely

opposite. Similarly, when other people are sarcastic at us and we are not able to understand the real meaning then other have fun and we ridicule ourselves unknowingly.

Few examples where not understanding the real intent of the person can be catastrophic.

- In face to face communication with your customer when you miss his intent. Result is customer disengagement.
- In live program when you are listening a response or question from the audience in hall or live TV or Radio program or speaking over phone or video conferencing tool and you miss the intent. Result is dent on your reputation.
- In offline communication when you publish some content on blog, news, product selling page and receive some comment from the public. Someone expresses his opinion over your post or tweet, and you are not able to understand that properly or not able to read. All other people read that comment and think that either you are dumb or do not care or accept what is being said. Result you know very well.

When you are dealing with your known people, friends, relatives and not responding properly in that situation, it will have lessor impact because they know your real nature and potential. But in public places, where you do not know the person to whom you need to respond, can cause huge dent on your image and brand.

1.1.5. Why Sarcasm Detection is Critical in Electronic Media?

India a great vibrant democracy so freedom of speech is natural to us. Most of the people in India communicate in Hinglish Language. In democratic societies people have opinion on everything irrespective of their educational qualification and experience. We are a country where public tells how Amibabh Bachchan should act, Virat Kohli should play cricket and how Narendra Modi should run the government. We have view and opinion on everything from politics to religion to product to government functioning to service delivery and what not. Many people choose to remain positive but express their negative feeling in sarcastic way. With the advancement of online sales of products, social media and online blogs, new portals there is huge surge of online feedback. Post COVID19 pandemic there are clear trends of shifting in this direction. People prefer buying, reading, expressing, engaging online. This justifies the need of sophisticated real time sarcasm detection system.

1.1.6. Sarcasm Detection in Hinglish

English⁸ is 1st most spoken language in the world and many researchers across the world are working for sarcasm detection in English. But, Hindi is 3rd most spoken language in the world and not much significant work is happening in sarcasm detection in Hindi. Unfortunately, nobody speaks in pure Hindi and it is considered pride unlike English, where people are shamed for not speaking or writing proper English. On social media and public forums very few Hindi speaker use Devanagari to express what they think. They use Hinglish Language. Due to this reason many of the feedback given on twitter, Facebook, product page, online news goes unnoticed and unanalysed.

Sarcasm is one kind of feedback and if we do not use this to improve our response then we prove ourselves foolish and customer shift to different product, service, or platform. Similar things happen when people change their party or group. Therefore, we feel it is extremely important to detect the sarcastic feedback given by those people who write in Hindi.

1.1.7. Challenges in Processing Hinglish Language

A. Complexity due to English words in Hindi

Observe the variation of a sentence “I have purchased tickets” in devnagari.

मैंने (टिकटें/ टिकटें/ टिकटे/ टिकिट) खरीद (ली/ लीं) (है/हैं). This simple sentence can be spoken in 16 different ways if written in Devanagari. If we mix Roman script in between then number of permutations goes beyond our normal imagination. Here we need to make note that Ticket is English word, and people are making plural of that as they do with any Hindi word.

Let us see another sentence “She has boiled the rice”

उसने राइस बोइल कर दिया है

From the above Hinglish sentence you cannot figure out whether the doer is female or male. Secondly, राइस and बोइल are not words in Hindi dictionary. Sometime people will write letter in Roman like

उसने Rice बोइल कर दिया है / उसने Rice Boil कर दिया है / उसने राइस Boil कर दिया है /

उसने Rice बोयल कर दिया है

⁸ https://en.wikipedia.org/wiki/List_of_languages_by_total_number_of_speakers

Like Guru, Karma are Hindi words and they are part of English dictionary. We do not have Hinglish dictionary which has word like यूज, गुड, नाइस, क्वीन in that dictionary. Without transliterating words like Tickets, Boil into Devanagari and telling system that टिकिटें = टिकटें = टिकटे = टिकिट, बोइल = बोयल embedding will not give good results.

B. Mix Other Indian Language with Hindi

Observe the sentence below, Bangla written in Devanagari and clearly understandable by any Hindi speaking person. Most of the words in the sentence below are from Bangla language but written in Devanagari.

अमी मौजूलिका.अमी राजा को मारबो दीदी ने केजरीवाल को भी पीछे छोड़ दिया. जि तो कमालई कर दओ ददू

India's business film Industry in Mumbai make film in Hindi. Rarely any film use as good Hindi as Hollywood uses English. Adoption of words from other language is not a problem. The problem is quantity of the words taken from other languages, availability of the updated vocabulary of the language. Many famous dialogues or songs from Hindi films which are taken different language or dialects. This increases complexity of sarcasm detection in Hinglish. We do not have comprehensive dictionary which we can call Hinglish dictionary which has all the word being used by the Hinglish speakers.

Without telling system that अमी (Bangla word) = मैं, मारोबो (Bangla word) = मारुंगी = मारुंगा = मारना no embedding is going to help

C. Complexity of Synonyms in Hindi

For this let us understand what Synonyms is. A word or phrase that means exactly or nearly the same as another word or phrase in the same language⁹, for example “shut” is a synonym of “close”. Few examples of synonyms

- The East = The Soviet Union (<https://www.lexico.com/en/definition/synonym>)
- Country of rising sun = Japan, Dragon Country = China,
- Fridge = Refrigerator
- Happy = Joyful, Cheerful, Contented, Jolly, Gleeful, Carefree

⁹ <https://www.lexico.com/en/definition/synonym>.

D. Influence of Sanskrit

All the synonyms have different spelling, different pronunciation but almost same meaning and part of the same language. l'eau (French word for water) is not synonyms of water because they are two different languages.

Unlike other world languages, all Indian languages (except Tamil, this is debatable) heavily borrow words from Sanskrit.

Let's take English word "Water" and see how many words are available in Sanskrit for "water" जल = पानी = तनि = नीरु = आपः = वाः = वारि = सलिलं = पयः = तोयं = मेघपुष्पं = घनरसः = पाणी. So all these words are synonyms of water in Sanskrit.

Because all Indian languages have root in Sanskrit therefore most of the time, they take word from Sanskrit for communication. For example, Kannada uses नीरु, Bangla use पानी, Hindi uses पानी, सलिलं, मेघपुष्पं. If not regular, they are used in poetical or sometimes in sarcastic language. Because in sarcasm or poetry we often use loaded words.

In Hindi language, can we say नीरु is synonym of पानी? No, because नीरु word is normally is used in Kannada and Sanskrit and not in Hindi. As per the definition of synonym another equal word should be from the same language and we know Hindi is not Kannada nor it is Sanskrit. The answer is yes also; because Sanskrit being mother of Hindi language, it borrows words freely from Sanskrit. Thus, we see synonym in Hinglish is not the way it is understood in the context of English.

Therefore, to be build a complete Hinglish dictionary we have take words from all other Indian languages and frequently used English words as well. Thus it should be like this.

जल = पानी = तनि = नीरु = आपः = वाः = वारि = सलिलं = पयः = तोयं = मेघपुष्पं = घनरसः = वाटर

E. Variation in Spelling of Same Word

In Hindi same word spoken and written with different spelling. Observe the spelling of the same word how they are varying. This kind of problem we do not have in English. As discussed earlier, synonym of Happy is Jolly. They both are not same, neither in spelling, nor in

pronunciation, nor in full sense, but “happy” is close to “jolly”. That is why they are synonyms. But below all “=” signs are referring to the same thing.

विष्णु = बिशणु = विशणु = बिष्णु = विष्नु = बिष्नु,

दरसन= दर्शन= दर्सन = दरशन

करता = कर्ता,

यज्ञ = जग्य,

योग = जोग,

हरि=हरी,

We need to keep in mind Hindi is not Devanagari, nor Hindi is Avadhi or Marathi. Hindi is written in Devanagari script but it is heavily inflicted by other languages like Awadhi, Bhojpuri, Rajasthani, Urdu etc.

Unless we have a dictionary which tells विष्णु = बिशणु = विशणु = बिष्णु = विष्नु = बिष्नु, embedding will not help.

1.1.8. Common Challenges in Sarcasm Detection

Detecting Sarcasm is difficult if sentences are having following characteristics.

- A. **Idioms and Phrases:** Sarcasm detection become more difficult when people speak in idiomatic language. For example: “What a wise man! what he did is nothing other than an axe to grind.” “कितना समझदार आदमी है जो उसने किया वो अपने पैर पर कुल्हाड़ी मारने के सिवा कुछ और नहीं है”
- B. **Speaking with Hint:** When people do not talk directly and use examples which are completely different than context. For example: “You are behaving like Mir Jafar.” “तुम्हारा व्यवहार मीर जाफर जैसा है”
- C. **Culture:** Different languages have different degree of challenges in sarcasm detection. For example, English is spoken all over the world but the way American express their feeling is different than the way British express. The reason for that is the work and social culture of England and United States is hugely different. In English language what is call sarcasm in England may be considered a normal statement or abusive in US and vice versa.

D. Datasource: Sarcasm can be present in any kind of communication platform like whatsapp, twitter, facebook, reddit, linkedin, product review, movie review, news review, blog review etc. But because of the type of audience, type of input interface, awareness of topic, command over language, character limit, text formatting possibility etc content available on the various platform has different characteristics. For example, twitter content is short and full of acronyms, words without vowel, scripting language mixed. On the other hand whatsapp group communications are full of links and forwards with little text written by sender.

1.1.9. Context Understanding a Challenge in Sarcasm Detection

Since the time human child take birth, baby has environment to learn from. Various types of formal or informal environment, social or business or cultural background forces human to think and learn. Either at physical or emotional or intellectual level if human fail to learn then his survival is challenged by the nature around. In this kind of environment, it is easy for any human to understand the context. If we are alert and interested in the topic then we need not to struggle hard to understand the context. But context understanding is extremely difficult in the case of Machine learning. Let us analyze one sarcastic tweet. “#JIO का सच नीता अंबानी ने मन्नत मांगी थी कि अनंत अम्बानी अपना वजन कम कर लेगा तो गरीबों में 3 महीने Net or call का भंडारा करवाऊँगी”

People living in India can understand that this is sarcasm. Because we know the full context. That

- Mukesh Ambani is owner of #Jio
- Neeta Ambani is Wife of Mukesh Ambani
- Anant Ambani is son of Neeta Ambani
- Anant Ambani has 200+ Kg body weight
- Normal body weight of human is around 70 kg
- Anand Ambani is overweight as per the normal standard
- Neeta Ambani desired that her son should have normal weight
- #Jio has launched 3 Month free internet package
- There is no direct connection between Anand Ambani weight reduction and 3-month free internet package

(Joshi et al., 2018) in their work “Investigation on Computational Sarcasm” says there are three type of context, Author Specific context, Conversational Context, Topical Context

We need to understand that keeping all the facts in mind we can say a statement is sarcasm and not normal statement. Even a human, who does not have all this information will fail to classify a statement as sarcasm. It is not easy to give all this information to a system to make a classification decision

1.1.10. Challenges in Sarcasm Detection in Hinglish

A. Script used for writing

70% of the world population uses 26 letters of Roman script to write their language. The Roman alphabet is also used as the basis for the International Phonetic Alphabet, which is used to express the phonetics of all languages.¹⁰ Due to this reason when people are writing different language like English, French, Indonesian, Tagalog, German, Turkish they need not to change much around the letters, so most of the cases script remain Roman. This advantage is not available to Devanagari script and Hindi language.

“Badhai ho kongressi Pappu ki vajah se #मोदी चुनाव फिर जीत गये” This entire sentence is in Hindi but notice script used is Devnagari and Roman. Not only that note the spelling of congress. Because that how native speaker think when he think about the sound of “क” or “K”.

While typing feedback people write @account_name. Most of the time @account_name are proper name and written in Roman like @harithapliyal, @eating_point, @banarasi. Similarly, hashtag, which helps us understanding the context of the feedback, is also written in Roman script #Election2019 #COVID19 #Philosopy #Motivation #NarendraModi.

B. Language mixed

An average westerner knows and speaks one language so written and verbal expression most of the time is that one language. An average Indian speaks minimum 2 languages, one is language of his state, plus national language, or English. In southern part of India, it is not uncommon when you find a taxi or truck driver who can speaker 3 or 4 languages but they cannot speak in English. This, one language- one script, advantage is not available for any Indian and they communicate in multiple language without realising that they have shifted language and borrowing words from different language.

¹⁰ <https://www.worldatlas.com/articles/the-world-s-most-popular-writing-scripts.html> Accessed on 23-Jun-20

“रहने दो उसको, उसके food preparation speed itna fast hai ki जितनी देर में राजधानी रेस्तरां वाले खाना घर पर डिलिवरी कर जायेंगे” This is sarcastic sentence about the laziness of the other person.

But analyze the words and language this

“रहने दो उसको, उसके” script Devanagari, language Hindi

“food preparation speed” script Roman, language English

“itna ---- hai ki” script Roman, language Hindi

“रेस्तरां, डिलिवरी” script Devanagari, language English

No matter how big corpora we use for tokenization and what kind of technique we use for tokenization until we have this kind of mix corpora for training sarcasm prediction in these kind of sentences is always going to be challenging.

C. Missing Context

“I love working hard” It looks normal sentence. But, if you add a context “my brother trying to still sleeping at 9am and saying” then meaning of the original statement is not what the speaker is saying. Thus the missing context or context not fully defined lead to issues of sarcasm detection in the sentences.

D. Limitation of Written Languages

Let's take one sentence “I didn't say he beats his wife”. It is simple statement by the speaker, where he is making a point about what he knows. But, how it is understood also depends in what tone it is said. If he emphasis on “his” then it looks like “I didn't he beats HIS wife” it can imply that he beats but not his wife. Written language has its own limitation message may not be expressed properly and tone of speech, body language, eye contact, facial expression etc which are part of audio-visual domain of communication has lot hidden in it. So, the message still may be sarcastic, but it is not part of the written words.

E. Usage of Idioms & Phrases

आ गया ऊंट पहाड़ के नीचे?

There is nothing special in the words of this sentence. But this is idiomatic phrase and you use it in some context and with interrogation marks then it is sarcasm on someone. It is not easy to know whether sentence contains idiomatic phrase or normal phrase.

F. Sentences containing Emoticon, Interjections etc.

अरे वा! इनको इस महान कार्य के लिए तो कम से पद्मश्री award मिलना ही चाहिए. 😊 😂

This looks normal sentence but emoticon and interjection is sarcastic

ओ साहेब, क्या समझ रखा है इतनी मेहनत के बात पद्मश्री award नहीं labour मजदूरी मांग रहे है 😞

Looks this second sentence it also has emoticons and interjections but it is not sarcastic. It is challenging task comprehend the meaning that too when text is mixed with emoticon and interjections.

G. Different Numerals

Many times, people use non-English numerals like १, २, ३, ४, ५. Depending upon the regional language people use different numerals for writing the same numbers.

1.1.11. Degree of sarcasm

Although how a person perceive & responds to a sarcasm it also depends upon him, yet we need to know all sarcastic statements are not equally intense or powerful to generate pain to the listener or reader. Here are few examples of different degree of sarcasm.

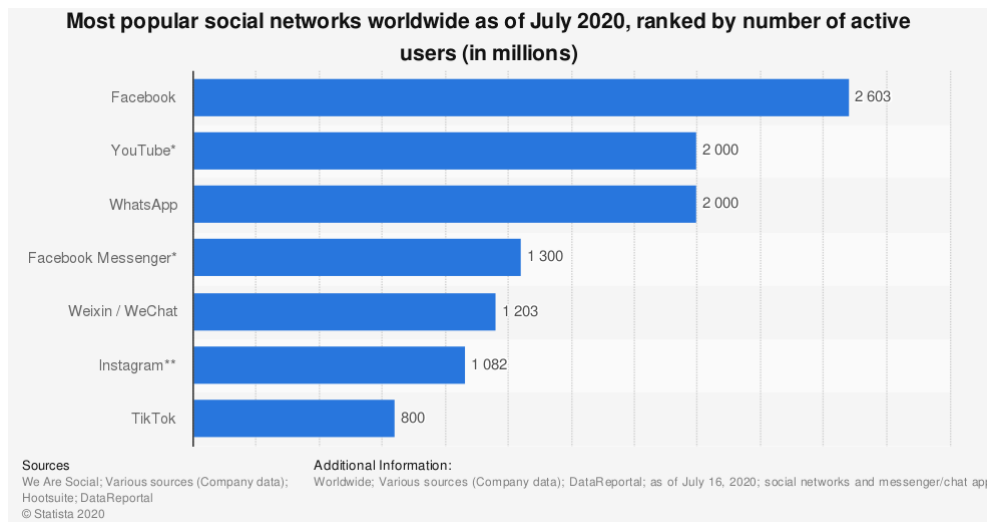
- A. ओ भाई कचोरी समोसे की दुकानें खुल तो गयी है लेकिन ध्यान रखे कचोरी समोसे के चक्कर में आप की ही पूड़ी सब्जी न बट जाये #Covid_Unlock (Least Intense)
- B. NDTV की हैडलाइन एक बेजुबान अल्पसंख्यक भैंस को डूबा कर मारने की कोशिश करती बहुसंख्यक चिड़िया (Lessor intensity)
- C. कोरोना का दवा न होना यह एक साइंस है, और दवा न होते हुए भी बिल लाखों मे आना ये एक आर्ट है !! (Moderate Intensity)
- D. ये शुक्र है जंगल में आरक्षण नहीं, बहोत नहीं तो जंगल का राजा शेर नहीं गधा होता. आरक्षण खत्म करो 70 साल हो गये यार #आरक्षण_भीख_है (Sharp Intensity)

1.1.12. Positive Side of Hinglish

Although India is big country with 1.35 billion people with different culture, religion, tradition but there is some common aspect in India culture and this does not change no matter where a Indian is living on the earth. That common culture helps us understanding the context and intent easily. Although there are many languages in India but because of one overarching culture it is easier to understand the meaning, a simple translation is good enough. Unlike English where Australian struggle to understand what American gentlemen want to say in English.

1.2. Problem Statement

More than 4.5 billion people now use the internet, while social media is used by approximately 3.8 billion users. Nearly 60 percent of the world's population is already online, and the recent trends highlighting that more than fifty percent of the world's total population will use social media by the middle of 2020.¹¹ IT companies like google, Facebook, twitter, amazon, Alibaba, Linkedin, Instagram, Quora dominate the content on Internet.



Keeping this volume, demand and need in mind, we want to develop a sarcasm detection system for Hinglish language which can work for all social media content, reviews, comments, and feedbacks.

1.3. Aim and Objectives

The aim of this research is to propose a model, which can predict sarcasm in a given Hinglish language sentence with highest possible accuracy.

Based on the above primary goal, objectives of this research are as following.

- A. To create Hinglish language dataset with minimum 2000 sentences, which can be used for training and testing a sarcasm detection model of Hinglish Language
- B. To develop a sarcasm detection models
- C. To check the effective of Transfer learning for this work.

¹¹<https://wearesocial.com/blog/2020/01/digital-2020-3-8-billion-people-use-social-media#:~:text=More%20than%204.5%20billion%20people,the%20middle%20of%20this%20year.>

1.4. Research Questions

- A. To study how sarcasm detection is done by other researchers for English and any other Indian languages?
- B. To determine which word embedding & linguistic features works best for sarcasm detection in our Hinglish dataset?
- C. How to do transliteration from Roman to Devanagari? Many options are available for reverse translation. For example “एकीकरण” => “Ekikaran” is easy and many options are there but “Ekikaran” => “एकीकरण” is not easy. Because Hindi speaking population is not aware about IAST¹² and nor they use it for transliteration. So confusion is how to convert word of Roman script to Devanagari, for example (a) “ra”=> र or र्, (b) n=> न or न् or ण or ण् or ञ or ज् or ड or ड्, (c) ki=> कि or की or क्ि or क्ी or क्इ or क्ई
- D. Is transfer learning useful for our work?

1.5. Scope of the Study

- A. This research is not related to any specific domain like philosophy, politics, history, current affair new etc. Rather it is trying to detect sarcasm in day to day informal conversation.
- B. Sarcasm in our communication can be expressed and experienced at Visual (facial express, body language), Vocal (tone, pace of speech, emphasis on certain word) and text (book, newspaper, articles, social media tweets, comments and feedback box on internet. Visual sarcasm is more universal than vocal. Because voice uses language and there are 7000+ languages on the earth so there is no universal vocal language of expressing sarcasm. But pause, pitch, pace, modulation between words, while speaking, are more universal like Visual. In this paper we are deal only with text-based sarcasm.
- C. Only Roman and Devanagari scripts are considered.
- D. Only Hindi and English language words are considered. If heavily used words from other languages which are part of day to social communication, then we will consider an opportunity to expand Hinglish vocabulary.
- E. No analysis of degree of sarcasm.
- F. We know to understand the context datetime plays a critical role. Our base dataset does not have datetime. And lots of the text in the dataset is coming from non-tweet sources which does not have datetime chronology of communication. Therefore, we ignored

¹² https://en.wikipedia.org/wiki/International_Alphabet_of_Sanskrit_Transliteration

context which is coming from datetime. We want our system to be indifferent of datetime metatag.

1.6. Significance of the Study

We didn't find one place which has done research and can say with conviction that approximately these are the number of Hindi speaker in the world. Different sources reveal different numbers. As per a lingoda.com¹³ and babbel.com¹⁴ after English and Mandarin Hindi is 3rd most spoken language on earth. It is spoken by 615mn people. As per Wikipedia 176 million people speak Urdu.¹⁵

Culture of Hindi speaking population and Urdu speaking population resembles a lot. While speaking or writing Hinglish many words of Urdu are spoken or written unknowingly. Therefore, any sarcasm analysis system in Hinglish will benefit Urdu speaking community as well.

With current trend of increasing online content in Hindi, it is practically not possible to read each and every review, even if you try it is very expensive and not worth work. We know, even one negative feedback or abuse which goes unnoticed can cause huge problem for the brand of the company, product, or person. Therefore, performing sentiment analysis on every feedback makes a perfect sense and it can be done automatically almost in real time.

Sarcasm is one type of sentiment and we are trying to discuss overall benefits of sentiment analysis keeping Sarcasm at the centre of discussion.

1.6.1. Application of Sarcasm Detection System

A. Sentiment analysis has a broad range of applications like understanding whether a feedback is Sarcasm, Warning, Love Emotion, Hate Emotion, Advertisement of some other product, Contradicting statement, Pun, Abuse, Inspiring Quote, Sensational Revelation, Pleasant Surprise, Allegation, Poetry/Dohe/Chands etc.

¹³ <https://blog.lingoda.com/en/most-spoken-languages-in-the-world-in-2020> Accessed on 22-Jun-20

¹⁴ <https://www.babbel.com/en/magazine/the-10-most-spoken-languages-in-the-world> Accessed on 22-Jun-20

¹⁵ <https://en.wikipedia.org/wiki/Urdu> Accessed on 22-Jun-20

- B. Government, NGO, religious leaders, product sellers are able to perform the sarcasm analysis against some product, political party, ideology, religion, company etc then they will be able to control the situation in much better way with minimum damage.
- C. Sarcasm analysis can be used to analyse the feedback on airlines service, travel service, bus or taxi service, telecom, health, government service, new articles, personal blog, food delivery, insurance service, personality page, book page are good places where sentiment analysis plays a critical role.
- D. In multinational companies it becomes exceedingly difficult to use humour to communicate the idea, crack joke or sarcasm, even if all the team member can speak English. The reason for that is different cultural background and different level of comprehension of English by non-native speakers. But when Hindi speaking people connect over video, telephonic or chat conversation it is easy for them to use idioms, joke, sarcasm and ensure that idea is understood. There is different kind of joy of working in lesser formal and light-hearted environment. When Indian people are speaking to each other using Hinglish we can perform sarcasm analysis to know the feeling of the group.

1.6.2. Motivation from Selected Domains

Below are examples of motivation written in English language. We have taken examples of sarcasm enabled chatbots. Answers given below by a chatbot is possible only if chatbot understand that input given is sarcasm and nor normal text.

A. Motivation in Travel Domain

Passenger: #ac_not_working. I love to get roasted in heat.

Chatbot: Sorry for the inconvenience. Our service engineer will call you.

B. Motivation in Hospital Business

Attendant: #expensive_treatment. We come to your hospital for this expensive treatment so that we can talk to your cute nurses.

Chatbot: We understand your concern about treatment cost. Our billing manager will call you.

C. Motivation in Restaurant Business

Customer: Last time, your food was so good that since last 2 days I am taking rest.

ग्राहक: पिछली बार आपका खाना इतना अच्छा था कि मैं दो दिन से आराम कर रहा हूँ

Chabot: I am sorry to hear that.

चैटबाट: यह सुनकर बहुत दुःख हुआ

D. Motivation in Learning Portal

Learner: What a great content. I am still trying to understand the head and tell of that since last 30 min video.

Chatbot: Sorry, can you please share with us what difficulty you faced ?

E. Motivation in News Portal

Reader: What a great story! Did you read it after writing?

Chatbot: We are sorry that you didn't like this story.

F. Motivation in Airlines Business

Traveler: First time in my life I got such a wonderful service from any airlines. I reached to the destination one day before my check-in baggage.

Chatbot: We are sorry to hear that. We hope your baggage reached safe to you.

G. Motivation in Dialogue Analysis Work

A dialogue from a Hindi Film "Sholey"¹⁶

मौसी मेरा दोस्त इतना अच्छा है कि वह शराब को कभी न नहीं बोल पाता। पीने के बाद जुआ खेलना उसकी खूबी है इसमें उसका कोई दोष थोड़ी है मौसी। बस हारने के बाद थोड़ा मारपीट करता है और घर में आ के मेरे को गाली देता है। पर मेरा दोस्त दिल का बहुत अच्छा है मौसी आप अपनी बेटी की शादी मेरे दोस्त से पक्की कर दो

This is a pure sarcasm paragraph. These kind of dialogues makes movie interesting.

1.7. Structure of the Study

Structure of the study is as following.

2.1. Sarcasm Detection Systems (SDS)

2.2. History of Sarcasm Detection Systems

¹⁶ <https://en.wikipedia.org/wiki/Sholey>

- 2.3. General Purpose Text-Based Sarcasm Detection System
- 2.4. Feature Engineering In Sarcasm Detection Systems
- 2.5. Approaches to Develop SDS
- 2.6. Approaches to Handle Key Challenges in Sarcasm Detection
- 2.7. Embedding
- 2.8. Types of Sarcasm Detection System

CHAPTER 2: LITERATURE REVIEW

Lot of work has been done in English language sarcasm detection and authors mentioned different challenges in sarcasm detection, although results are not that great as for any other classification problems. Challenges exists because of context understanding, missing context, domain, culture, different words, or expression used by people to flip the meaning etc. There is not much work done in Hinglish Language Sarcasm detection. Hinglish language has a separate set of challenge like mixing script, mixing language, highly morphological words, using same morphology on English language words, meagre size of corpus etc.

Let us take one English verb “do”, in Hindi, it can be used like कर्ता (noun) , करता (verb with male), करती (verb with female), करूंगा (future tense with male), करूंगी (future tense with female), करेंगे(future tense with plural), किया (did, done), करो (request, must do) करें (please do) etc. these all are with different gender, mood and tenses. However, in English we have inflection like do, does, did, done. These inflections are such that even without using pronoun sentence is meaningful. For example, करता है = वह करता है. Even without pronoun वह sentence is correct, complete and meaningful. While this is not true in the case of English language.

Now, let’s take another example but this time we take noun “Ram”. राम का, राम ने, राम को, राम द्वारा, राम में, राम पर, राम के लिए, राम पर and many times you will see letters are written together. We never see any word like “ByRam” in English but in Hindi रामने and राम ने both have same meaning.

2.1. Sarcasm Detection Systems (SDS)

Sarcasm is perception of the human receiver about some inputs. “Input” can be of four types. First type of input is text format written in social media, book, newspaper etc. Second type of input can be vocal tone, expressed in some voice communication over phone, face to face meeting, stage show, etc. Third kind of input can be image appearing on some public hoarding, newspaper article, blog post, social media etc. Forth kind of input can be body language of human during face to face interaction or in video.

To understand a message correctly following conditions should be met successfully.

- Speaker speaks in the language which listener can understand
- Listener understand the background
- Listener has technical knowledge about the subject

Beauty is in the eye of beholder. If receiver missed the intent of input due to any reason, then will you call that statement sarcastic? This is philosophical debate and in our work, we will be focusing on text which is marked as sarcastic by different annotators. From receiver's perspective input received can be any of the following four types.

All Weather Sarcastic: Every civilized person will treat those statements as sarcastic. For example, "I like when you treat me like a slave". No matter what the context is, what language is used to communicate this text everybody will say this is sarcastic statement. No other information is required, sentence has complete information and almost all human agree to this.

Conditionally Sarcastic I: More information is required to classify a sentence is sarcastic or not. In the presence of that important information we can confidently say this is normal or sarcastic sentences. This more information may be related to profession, culture, rules, law of the land etc. For example, "I love to beat drum at 5 am in the morning". Some cultures, profession forces their follower, community members to do eating, praying, singing, playing activities at a time so in that profession or community's context it may be normal. Otherwise it is sarcastic.

Conditionally Sarcastic II: Sometimes we need *individual event and person specific information* to detect sarcasm in sentence and this information cannot be generalized even for the same person at other time. For example "First thing in morning I like to do is cleaning potty of Ruby". If receiver know the context that the speaker is mother and Ruby is new born baby then speaker may say it is sarcastic or non-sarcastic depends upon receivers individual like or dislike. Here even after understanding the full context it is depending upon the receiver who does the classification. But if receiver knows that that speaker is busy CXO and Rudy is his pet name. Then receiver will say it is definitely a sarcastic statement.

Non-Sarcastic: Normal sentences with straight forward meaning without hiding any intention and no scope of different interpretation.

Sarcasm detection system is one which can flag a "input" provided to the system as sarcastic or non-sarcastic. In the context of our project "input" mean text and no other type of input like body language, image, video, speech etc. Even with text as "input" we are particularly dealing with one or two liner text appearing on social media or day to day communication. We are not

dealing with long chain of text like a paragraphs, a page, a chapter or a book. We are interested in developing state of art sarcasm detection system for Hinglish language. Systems takes input as one or two sentences with full context and returns True if Hinglish sentence is sarcastic else returns false. If the context is missing, then system may fail to predict correctly.

2.2. History of Sarcasm Detection Systems

We have prepared a separate report on [History of Sarcasm Detection](#). If you are interested in the chronology of the development you can check it from [github](#) link.

2.3. Generic Text-based Sarcasm Detection System (GTSDS)

There are many dimensions of complexity in any sarcasm detection. General purpose text-based sarcasm detection system means a system which can detect sarcasm in any text. Before building a GTSDC we need to answer following question.

- Can we develop a SDS which can detect sarcasm in any human language like English, Hindi, Japanese, Chinese, Spanish?
- Can we develop a SDS where text written in any script like Devanagari, Roman, Hebrew, Chinese?
- Can we develop a SDS where text written in simple language vs figurative language which is full of proverbs and coded words?
- Can we develop a SDS where text using words from any business domain like politics, philosophy, medical practitioners, lawyers etc?
- Can we develop a SDS text written by the people of different culture like British, North American, Indian, Japanese etc?

Building a general-purpose text-based sarcasm detection is extremely complex task. In this work we are trying to develop a general purpose SDS where the text of two scripts and words from multiple language like Hindi, Sanskrit, Urdu, Punabi, Marathi, Bhojpuri, Avadhi are used. We are aware this is not a complete *generic purpose text-based sarcasm detection system* and but a step towards that.

2.4. Feature Engineering in Sarcasm Detection Systems

Researchers have extracted various features from the given text to detect whether sentence is sarcastic or not. These features can be grouped under following categories.

- **Lexical:** unigram, bigram. These can be created using words or characters.

- **Pragmatic:** These features are created using emoticons, punctuations and capital letters used
- **Incongruity:** Incongruity in the sentences is detected using novel approaches.
- **Polarity:** Polarity of the noun, adverb, adjectives are counted
- **Syntactical:** These features are based on POS (part of speech)
- **Idiosyncratic:** Sentences are analysed for repetition of any specific word by the speaker. Many people have habit of say words like “I know”, “you know”, “yah yah”, “absolutely”, “like” etc.
- **Prosodic:** Analysing pattern of words in the sentence, how a specific words is written to emphasis something. For example “It is soooooooooooooooooo beauuuuuuuuuuuuuutiful”
- **Features based on the Author’s or reader’s profile data:** Gender, nationality, religion, education, ideology, familiarity of language etc.
- **Features based on the environment:** Datetime, current news, messages in past, present state of mind etc.
- **Hashtag & @users:** Different hashtags used and different users tagged in the message
- **Slang:** Number of slang used, ratio of slag to normal words, nature of slang word etc.
- **Profanity:** Any dirty, abusive, naughty, offensive words

There are many creative ways to create hundreds of features under above categories. We will refer all these features as Linguistic Features of the Sentence (LFS)

In their work, (Joshi et al., 2018) have used 3 types of features POS, Named Entities, Unigram to predict the disagreement. (Sharma et al., 2014) in their work “A Sentiment Analyzer for Hindi Using Hindi Senti Lexicon” suggests using bootstrap approach to extract senti words from Hindi Wordnet. It has given encouraging results of 87% accuracy in sentiment analysis. We are going to test usefulness of this approach in sarcasm detection.

We have prepared a [“Summary of Papers on Sarcasm Detection”](#). This presents a summary of these features used by different researchers and the performance reported by them. If you interested to read more, you can refer to the github repository.

2.5. Approaches to Develop SDS

Over the period of last 20 years different approaches are adopted by different researchers. Broadly these can be categorized into following categories. In the following subsections we are analyzing features explored, algorithms used, and results gained by the different researchers. If you want to more about these then you can refer to our work [“Summary of Papers on Sarcasm](#)

[Detection”](#) and [History of Sarcasm Detection](#). Table below presents the summary of approaches used to develop SDS. Numbers written in the cells of the table are section number following the table.

| | | Feature Types | | |
|---------------------|-------------------------|---------------|-----------|-------|
| | | LFS | Embedding | Both |
| Classification Type | Rule Based | 2.5.1 | x | x |
| | Classical ML Algorithms | 2.5.4 | 2.5.3 | 2.5.2 |
| | CNN | 2.5.9 | 2.5.6 | 2.5.5 |
| | Transformers | x | 2.5.7 | x |
| | Transfer Learning | x | 2.5.8 | x |

2.5.1. Rule based Approaches

In this approach researcher depends upon the content and context-based of the text. They extracted various Linguistic Features of the Sentence (LFS). Some experimenters have demonstrated a good performance on sarcasm detection work without using any machine learning algorithm. “Lexicon-Based Sentiment Analysis in the Social Web” by (Asghar et al., 2014) didn’t use any classical or neural network based algorithm for this work. They could achieve 95% accuracy by using a) Lexical features- unigram using chi-square test, (b) Pragmatic- emoticons, punctuation marks, capital words, (c) Explicit congruity- related to polarity changed, and (d) Implicit incongruity features.

Just using rule based approaches (Bharti et al., 2018) achieved 87% accuracy on Hindi language tweets and (Sharma et al., 2014) could achieve 85-89.5% accuracy on Hindi language product reviews.

2.5.2. Classical Machine Learning with Embedding & Other Features

Using this approach, we create LFS along with word embedding for every sentence. (Kumar et al., 2019) used tokens using Classical language toolkit, unigram, bigram. They also used

fastText and TF-IDF embedding. Authors used SVM linear kernel, LR, RF, Shallow CNN + Bi-Directional LSTM for classification purpose. In their work “BHAAV- A Text Corpus for Emotion Analysis from Hindi Stories” they were trying to classify emotions in Hindi language sentences. They claim that they could get an accuracy of 62%.

2.5.3. Classical Machine Learning with Embedding

In this approach we need not to explore any LFS. Using word embedding word vectors are created and they can be used for creating classification model. During literature survey we could not find any papers which solely rely on word embedding for creating the models.

2.5.4. Classical Machine Learning without Embedding

In this approach we can use LFS but classification is done using the classification machine learning algorithms like LR, SVM, RF etc. Many experiments are done using this approach.

(Fafias et al., 2016) demonstrated 73-96% accuracy using classifiers like SVM, DT, NB and feature engineering approaches. (Suhaimin et al., 2017) shown 82.5% accuracy with non-linear SVM. Both of these experiments are done on English tweets.

(Sundararajan and Palanisamy, 2020) used English twitter data and shown 86.61% to 99.79% accuracy using classifiers like Random Forest, Naive Bayes, Support Vector Machine, K-Nearest Neighbor, Gradient Boosting, AdaBoost, Logistic Regression, and Decision Tree. They extracted 20 features from the dataset.

In an another interesting work on Instagram image (English text), (Kumar and Garg, 2019) has developed a sarcasm detection system with 73% to 88% accuracy. They extracted features like Number of negative words, number of positive words, POS tag, hashtag from the dataset.

2.5.5. Deep Learning with Embedding & Other Features

In this approach, researchers created features using both the word embedding and LFS. For classification researchers used deep learning networks like CNN. In “CARER: Contextualized Affect Representations for Emotion Recognition” (Saravia et al., 2020) used BoW, char n-gram, TF-IDF, Word2Vec, fastText(ch), word-cluster, enriched patterns, Twitter-based pre-trained word embeddings and reweight them via a sentiment corpus through distant supervision. Authors used

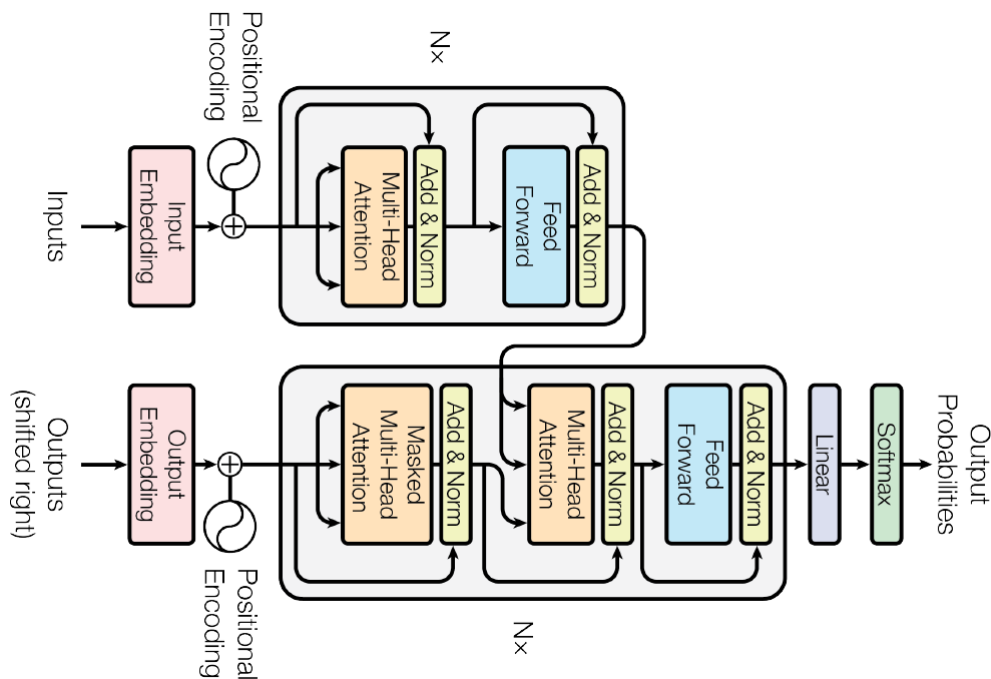
CNN for the classification and claimed an accuracy of 81% using their novel architecture named CARER.

2.5.6. Deep Learning with Word Embedding

Deep learning approaches includes those experiments where experimenters have used CNN, RNN, GRU, LSTM or any variation of neural network. They transformed the text input into vectors using different embedding techniques like TF-IDF, word2vec, fastText etc. (Subramanian et al., 2019) used GRU on English language twitter and facebook dataset and show 89.36% accuracy on twitter dataset and 97.97% accuracy on facebook dataset.

2.5.7. Transformer Based

(Vaswani et al., 2017) in their work “Attention Is All You Need” proposed a novel architecture which named Transformer Model Architecture. A transformer has two units first is encoder, and second is decoder. Subcomponents of transformer architecture are positional encoding, multi-headed attention, feed forward network, masked multi-headed attention, fully connected dense layer and finally Softmax layer.



Source: (Vaswani et al., 2017)

Several companies are taking lead and exploiting this architecture to build state of art model for NLP tasks. Below is the list of some selected transformer-based models by various companies.

1. [GPT](#) from OpenAI by (Peters et al., 2018) in their paper “[Improving Language Understanding by Generative Pre-Training](#)”
2. [BERT](#) from Google by (Devlin et al., 2019) in their paper “[BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#)”
3. [XLNet](#) from Google & CMU by (Yang et al., 2019) in their work “[XLNet: Generalized Autoregressive Pretraining for Language Understanding](#)”
4. [ALBERT](#) from Google Research by (Lan et al., 2019) in their work “[ALBERT: A Lite BERT for Self-supervised Learning of Language Representations](#)”
5. [T5](#) (from Google) by (Raffel et al., 2019) in their work “[Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer](#)”
6. [ELECTRA](#) from Google Research & Stanford University by (Clark et al., 2020) in their work “[ELECTRA: Pre-training text encoders as discriminators rather than generators](#)”
7. [RoBERTa](#) from Facebook by (Liu et al., 2019b) in their work “[Robustly Optimized BERT Pretraining Approach](#)”
8. [DialoGPT](#) from Microsoft Research by (Zhang et al., 2020) in their work “[DialoGPT: Large-Scale Generative Pre-training for Conversational Response Generation](#)”
9. [DistilBERT](#) from HuggingFace by (Sanh et al., 2019) in their work “[DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter](#)”.

(Potamias et al., 2020) developed novel architecture RCNN-RoBERTa in their work “A transformer-based approach to irony and sarcasm detection”. They developed this architecture using an existing transformer RoBERTa. As mentioned above RoBERTa is developed by facebook research for natural language processing tasks. This novel architecture by authors could predict the sarcasm with 85% to 94% accuracy. In this paper authors has compared performance of various kind of transformers like ELMo , USE, NBSVM, FastText, XLnet, BERT base cased, BERT base uncased, RoBERTa base model, UPF, ClaC, DESC to compare the performance of their novel architecture.

2.5.8. Transfer Learning Approaches

Training a new model from scratch is expensive and time-consuming work. Therefore, recent trends of Transfer Learning are picking up. Some companies or universities who have plenty of resources to develop new models using large amount of data develops the model of various size. They release the models, which need lesser resources to run, for the consumption by other researchers, who have lessor resources & time at their disposale. These released models are called pre-trained models. We can use models as is or with some fine tuning, based on our need.

The pre-trained models are developed using corpus of some language, some task and text from some domain. The beauty of these models is we can finetune them using our data for a task which we want to accomplish. This is called transfer learning. Using transformer-based system we can perform three kind of transfers namely *task transfer*, *language transfer* and *domain transfer*. In the NLP world task means like classification, comprehension, text generation, next word or sentence prediction etc. When we say task transfer it means a model which is created for let us say classification task can be finetuned for next word prediction or any other task. To use the model, we need to convert the text into vectors using embedding provided by the transformer. Recently we see a surge of models of various size in the NLP marketplace. Researcher community is happily adopting those for their experiments and getting good results compare to other approaches and techniques mentioned earlier.

2.5.9. Deep Learning without Embedding Features

In this approach deep learning neural networks are explored for classification but instead of using word embedding Linguistic features of the sentence are used. (Liu et al., 2019a) in their work “A2Text-net: A novel deep neural network for sarcasm detection” used CNN and created a novel architecture for sarcasm detection. They tested their model on different dataset and got different results. The results vary between 71%-90% F1 score.

2.6. Approaches to Handle Key Challenges in Sarcasm Detection

2.6.1. Handling Figurative Languages

Figurative language is the language used by intellectuals or those who have a good command over language. If you take a literal meaning of a sentence written in figurative language you will not get anything useful and meaningful. Many times educated people of the society want to communicate some idea or message but they use simile or old proverbs or chose words which are not directly related to the situation but the gist of that incident or proverb has parallel to the situation in hand. Figurative language is work of intellectual caliber and many times it is not easy even for human to understand the message. For example “My daughter is apple of my eyes” मेरी बेटी मेरे आंख का तारा है” If you miss the presence of figurative language in this sentence then you will miss the meaning of this sentence.

(Potamias et al., 2020) in their work “A transformer-based approach to irony and sarcasm detection” claims their novel architecture RCNN-RoBERTa performs well on the figurative

language. (Nozza et al., 2016) in their work “Unsupervised Irony Detection: A Probabilistic Model with Word Embeddings” claims that if we integrate probabilistic models like TIM with word embedding then we get promising results in detecting irony and sarcasm.

2.6.2. Handling Limited Data in Sarcasm

Although we didn’t find research work which tells how much percentage of our day to day communication is sarcastic, but we know from our day to day communication that percentage is very less. You just observe yourself or family members around for one day and count how many times you used sarcastic language. Due to this reason, we do not have enough good size dataset of sarcastic communication. Hindi & Hinglish being one of the least NLP resource languages has too little data to build a good sarcasm detection system.

Researchers takes either of the two approaches to handle small size sarcastic dataset. In first approaches they do not take more non-sarcastic sentences than they have sarcastic ones. So, the dataset is balanced but of the small size. In second approach they do over sampling of sarcastic sentences.

In either of the cases if dataset size is small for the training purpose we use cross validation techniques. In the technique we create multiple fold of the same dataset using random sampling and then use the fold for the training purpose. Let’s say our data set has only 1000 records and it is balanced dataset. If we create a 5 folds cross validation for the training purpose then 5 folds of 200 records will be created from these 1000 records. These 5 folds should have same distribution of the classes. Every time we create new folds there will be different set of records in those folds. After 5 folds are created, we can use 4 folds for training and 1-fold for validation purpose. Thus, we run train our model 5 times and every folds gets opportunity to become validation set.

2.6.3. Handling Out of Vocabulary Issues (OOV)

To create word embedding vector which represents all the possible words and their possible usage in different context we need a huge corpus. Not only this, if we have huge corpus of political news or short moral stories that will not represent the same words which are used in the context of medical, physics, philosophy, finance etc. For example, “Interest of various stakeholders is increasing in the recent peace talk process”. This is a statement from normal news. But “Banks are continuously increasing interest and it is making capital more costly” is

a statement from financial news. Same word “interest” in financial news has different context than when it is used in normal life. To make sure that final word embedding represent all the possible context we need to include corpus of all the possible domain’s data. But this is difficult task as of today. Because of limited good quality corpus from all the domain of business, science, technology, culture etc.

Due to this reason, at training or prediction time, when we are looking for a word vector for a new context and if word embedding is not available then that word becomes OOV word. When our dataset has many OOV words then training task will not be able to generate a model which can perform NLP, NLU task with good results. Similarly, if word is available at the time of training but it is not available at the time of validation or in real environment then due to OOV NLP, NLU task performance will be poor, and nobody will use that model.

OOV problems becomes serious when we are using a dataset for training which has words from multiple languages and multiple scripts are used to write those words. This is the typical case of Hinglish language especially in social media or whatsapp communication between Indians. Although there is no silver bullet solution for this OOV problem but if do following we can address this problem to a large extend.

- A- Use large corpus
- B- Use corpus of different domains
- C- Use corpus which has text written in multiple scripts
- D- Use corpus which has words from multiple languages
- E- Instead of creating context-based vector for words, create subwords from the word and create context vector of those. This is the approach used by fastText of Facebook.

2.7. Embedding

Computers cannot understand text so we need to convert them into numbers. But how to convert a word, phrase, sentence, dialogue, paragraph, chapter, news article, book or encyclopaedia in number? Broadly there are two approaches one is frequency based and another is prediction based.

TF-IDF: Term frequency inverse document frequency is frequency based embedding approach. This is a numerical statistic technique that is intended to reflect how important a word is in a collection or document. TF-IDF numbers of a word imply a strong relationship with the

document they appear in, it suggests that if that word were to appear in a query, the document could be of interest to the user, (Ramos, 2003) .

CBOW: Continuous bag of words is a prediction-based technique. It predicts the probability of the word if a context is given. Context window is number of words around the word. Context window of size one means one word left and one word right of the main word. (Wang et al., 2017)

Skip-gram: Skip gram is another prediction-based technique. If we want 3 gram one skip, skip-gram from a sentence “I hit the tennis ball” then we get following skip-grams “I hit the”, “hit the tennis”, “the tennis ball”. This gives us good context understanding. However with this approach a problem of sparsity of the word becomes more severe, (Van Brunt, 1987).

2.7.1. Absolute Embedding

Word embedding like TF-IDF are absolute word embedding approaches. In these approaches word meaning is fixed irrespective of the context a word is used. We know from our experience that meaning of same word can change from one domain to another and one context to other. For example, “गया गया गया”. English meaning “Gaya went to Gaya”. First word is subject, second word is verb and the third word is a location. Absolute embedding approaches cannot handle this kind of text and because of wrong vector classification or any NLP task will be incorrect.

2.7.2. Contextual Embedding using full word

Three popular and most used absolute embedding vectors are **glove**, **word2vec** and **freebase**. **Glove840B** is pretrained word vector with 940 billion tokens. This is developed by Stanford university. **Word2vec** [[GoogleNews-vectors-negative300.bin.gz](https://drive.google.com/file/d/0B93t2m34E0rUaGZhd2E0ZWZlZDQ0/view)] is pretrained word vector with 100 million tokens. **Freebase** [[freebase-vectors-skipgram1000.bin.gz](https://drive.google.com/file/d/0B93t2m34E0rUaGZhd2E0ZWZlZDQ0/view)] is pretrained word vector with 1.4 million tokens. Word2vec and Freebase are developed by google using google news dataset. In the contextual embedding different meaning of one word in different context can be represented by different vector of the same word. Contextual embedding is done using skip-gram and CBOW. Full word is used to develop this kind of embedding. Issue with this kind of embedding is OOV. If you create word vector using this embedding post lemmatization of word then context is not fully represented but OOV problem will be less. If you develop word vector using this embedding without lemmatization, then OOV problem will be more and matrix will be too sparse.

2.7.3. Contextual Embedding using subwords

As discussed above contextual Embedding using full word cause OOV problem during the training. To address that problem this technique create subwords from a word and then create word vector of those subwords. Final word vector is sum of all these vectors. **fastText** uses this technique to create word vector. Fasttext treats each word as composed of character ngrams. So the vector for a word is made of the sum of this character n grams. Let's say there is a word "apple" in the sentence so to get the word vector of "apple" we need to sum all vectors of the n-grams of apple "<ap", "app", "appl", "apple", "apple>", "ppl", "pple", "pple>", "ple", "ple>", "le>". Assuming ngram-min is 3 and ngram-max is 6. This embedding technique also uses n-gram and CBOW for creating word vector.

In their paper, "Adaptive GloVe and FastText Model for Hindi Word Embeddings", (Gaikwad and Haribhakta, 2020) states that AGM gives better results than GloVe and FastTextWeb. They also mentioned that FastText embeddings which are trained on FastTextHin (Hindi Monolingual corpus) produce better results than FastTextWeb. Google research has introduced a multilingual BERT which is capable of working with more than 100 languages (Romano, n.d.).

2.8. Types of Sarcasm Detection System

Sarcasm detection systems can be classified in following ways

- **Architecture Used:** Based on the architecture used to develop the system.
 - Rule Based
 - Classical Machine Learning Based
 - Neural Network Based
 - Transformer Based
- **Domain Specific:** Based on the domain it serves.
 - Health
 - Education
 - Travel
 - Social
 - Generic (It is extremely difficult to build a generic SDS)
- **Mode of Communication Based:** This classification is based on the mode of inputs it can accept to perform the classification.

- Text Based Systems: They can process only online or offline text is used as an input.
- Voice Based Systems: They can process only voice signals.
- Video Based Systems: They can process only videos.
- Image Based Systems: They can accept only images.
- Multimodal System: These systems can take any form of input to perform the classification. It is challenging task to build a SDS which can take all type of inputs, as mentioned above.
- **Time of Detection Based**
 - Realtime Systems: In real time message can be classified as sarcasm or not. For example as soon as message is delivered on whatsapp, twitter, facebook receiver get a different kind of tick message that it is sarcastic message.
 - Batch Systems: At the end of day or any other frequency, based on the need, all the messages or text can be processed in batch to know how many of them were sarcastic.
- **Language Script Based:** This classification is based upon spoken Language used to write message and written script used to write message.
 - Language Specific: Only for specific language like German or Japanese or Hindi etc.
 - Multiple Language: Can support any spoken language of the world. It is very challenging work to develop such a model.
 - Script Specific: Only for a specific script like Roman or Devanagari or Chinese etc.
 - Multiple Scripts: Can support any script of the world. It is very challenging work to develop such a model which supports all the scripts of the world

2.9. Summary

Thus, we see many researchers have tried to perform the task of sarcasm detection and achieve different accuracy or F1 score depending upon their experiment setups. They have tried different feature extraction techniques and applied those features on different classification algorithms. Most of the work has happened in English language and results are not consistent because they are depending upon quality of text in dataset, domain, classification techniques used, features used, data source used etc. Some work has been done for Hindi language and other Indian languages. We did not find any work in Hinglish language which is beyond Twitter dataset. If we observe the table in Appendix B we cannot say with certainty that there is any

significant improvement in sarcasm detection results when we use transformers or CNN. We want to experiment with different features and classification algorithms and understand what best results we can achieve when want to detect sarcasm in Hinglish language text.

CHAPTER 3: RESEARCH METHODOLOGY

In this section we are going to discuss a high-level approach to accomplish the research goal. The flow of discussion in this section is as following 3.1. Dataset, 3.2. Feature Engineering, 3.3. Overview of Our Approach, 3.4. Model Building, 3.5. Evaluation Metrics & Reporting, 3.6. Development Tools.

3.1. Dataset

3.1.1. About Dataset

We started building dataset using Hindi tweet dataset.¹⁷ This excel file had total 442 records. But this sheet does not have labelled data. We cleaned this file, removed ambiguous sentences and put data in our required format. For our project we needed data in the two column format 1- Sentences 2- Label. After cleaning this data, we had 300 labelled sentences but this is not sufficient for building a reliable sarcasm detection model for any language. So, we decided to expend this dataset to 2000 sentences with balance data, i.e. 1000 sarcastic sentences and 1000 non-sarcastic sentences. This new dataset has data from tweet as well from normal text or story blogs. All the sources we used to scrap the Hinglish data are available in [github file](#).

3.1.2. Data Sourcing

Sarcasm data in Hindi and Hinglish language on internet is very less. Whatever data is available it is too scattered and painful to extract the data to build a reasonable good size dataset for model building. After lot of surfing on internet we finalized 36 twitter accounts, 22 blogs, and 2 hashtags to scrap the data. To extract the tweets from 36 twitter accounts we wrote some python code using tweepy api. To extract the tweets from 2 hashtags we did little change in the earlier code and could scrap the tweets. To extract sarcasm from blog post we followed two steps. 1- Copied text from the blog post. 2- Read the blog text and break the sentence wherever it looks logical. All the tweets and sentences are put together in one csv file. All the text put in one column "Sentence". Sentence id is generated for each sentence

3.1.3. Dataset Cleaning

Most of the text from blog was clean but twitter had uncleaned, unstructured sentences. We know that tweet text is unclean because it has text from different languages, in different scripts, extra space, emoticons, non-text sign like "~" ":", "<" etc, flag sign, line break, over used words like ".....", "???????", "beau.....tiful", "!!!!!!". Blog text may also have this kind of text but

¹⁷ <https://github.com/rkp768/hindi-pos-tagger/tree/master/News%20and%20tweets> (Accessed on 26-Jun-20)

chances of that is extremely less. Now onwards we will not refer this as tweet or blog text but as sentences. Save all the clean sentences text in a new csv file. We wrote a python script to clean all the text. We used following [checklist](#) to clean the text, this file is available in github.

3.1.4. Sentence Labelling

Because of various reasons as mentioned earlier, a sentence cannot be labelled as sarcastic by all the people. People have different opinion and it varies based on individual's personality, education, environment, mood at a specific time and other human personality factors. Before we proceed with our dataset, we wanted a dataset which has unbiased labels. Final dataset has 2368 sentences. To label these sentences as sarcasm we used three annotators who are native speakers and use Hinglish in their day to day communication. All three annotators labelled each sentence independently. Whatever was the max vote for a sentence that label was finally assigned to the sentence.

3.1.5. Language Treatment of words

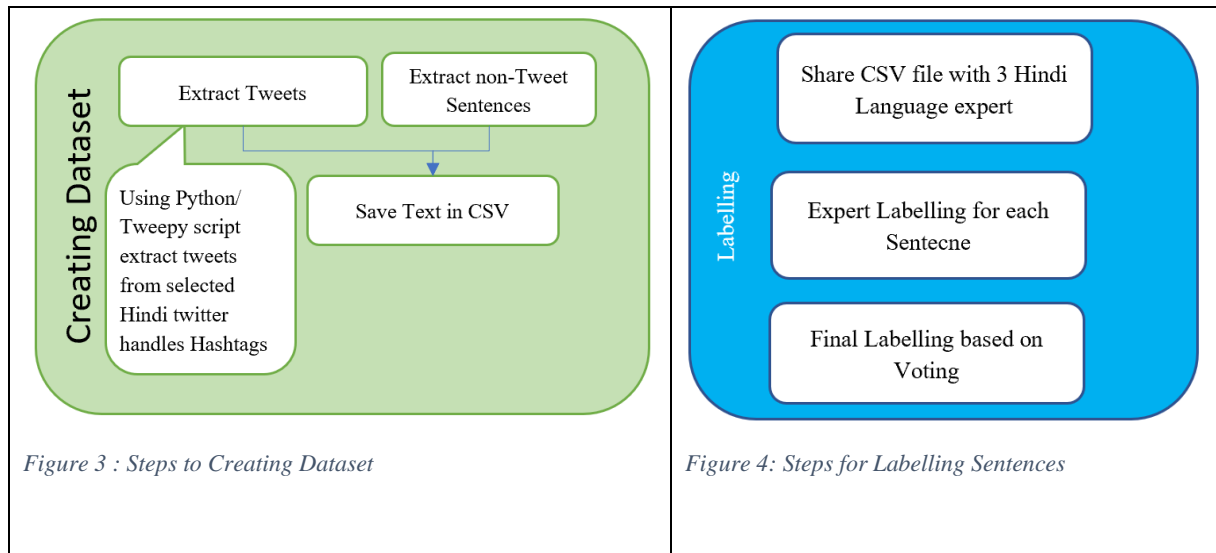
We know Roman typing is much easy compare to typing in Devanagari therefore many time people use Roman letters in between the sentence. This is true especially if it is name of politician, film actor, place name, (#AmitShah, #Modi, #Khan, #India #Bollywood, #Delhi, #Karnataka #Yogi) etc. Because same word will be written in Devanagari and other times in Roman and this is not good for text analysis. To remove these kind of anomalies so that we have the best possible word embedding we took following steps. They are summarised in the table.

- **Hindi word in Roman script words:** Non-English words written in roman scripts like “Aap to Mahan hai”, “tussi great ho ji” is converted to Devanagari.
- **Non-Hindi Indian words in Devanagari:** Words from other language like Urdu, Punjabi, Marathi like खत्म, खल्लास, गजल, सोन्देश are left as.
- **Hindi word in Devanagari:** No change
- **Hindi word in Roman script:** Transliterate to Devanagari
- **English words in Devanagari:** English words written in Devanagari like “राइस” “विन” “ग्रेट” are left as is.
- **English words in Roman:** English words in Roman like “Courtesy to my friend” is left has. But any noun is converted to Devanagari script.

| | | Script | |
|----------|---------------------------|---|---|
| | | Devanagari | Roman |
| Language | English | राइस, ग्रेट [No change] | Great, Win, Fine, Everest, Trump, Tonny [Transliterate to Devanagari, if noun, Translate to Hindi, if verb] |
| | Hindi | ठीक, है, मेरा, नाम [No change] | Aap, Khana, Kha, lo, Hari, Narendra [Transliterate to Devanagari] |
| | Non-Hindi Indian Language | खत्म, खल्लास, गजल, सोन्देश [No change] | Gazal, Neeru, Tanni [Transliterate to Devanagari] |

3.1.6. Dataset Structure

1. Dataset will have 3 columns “Id”, “Sentence”, “Label”
2. Sentence: Sentence is text of the tweet or any normal sentence.
3. Label: This column will have 0 for normal sentence and 1 for sarcastic sentence.



3.1.7. Feature Engineering

Linguistic Features

We will explore following features.

- a) POS based
- b) Hashtag
- c) Emoticon
- d) Polarity Based

Word Embedding

We will explore following word embeddings.

1. TF-IDF
2. fastText
3. BERT

3.1.8. Final Dataset with All Features

The final dataset with all the features will have embedding features as well linguistic features.

While creating model we will create five datasets.

- 1- A dataset only with linguistic features
- 2- TF-IDF Embedding
- 3- fastText Embedding
- 4- BERT Embedding
- 5- A dataset with linguistic + Best Embedding (depends upon results)

All the models we are building will be build using all five datasets. And metrices will be compared to see which model works best on which type of features.

3.2. Overview of Our Approach

We are starting this project with almost zero data in our hands. So the first steps is create a good size dataset which can be used for our project. The details are mentioned in section 3.1.2. The dataset created is not fit for model, so we need to need to clean this dataset. The details are mentioned in section 3.1.3. Following this we need to manually label each record with the help of annotators. The details are mentioned in section 3.1.4. We are dealing with Hinglish language and if we do not treat script and language of sentence properly then word and context cannot be represented correctly by the word vector. Hence, we need to handle the language and script. The details are mentioned in section 3.1.5. After following all above steps we have clean data in place. We can do word embedding and linguistic feature creation from this data. Following that we combine different linguistic feature and embedding based features for the modelling. Then we develop model and evaluate the performance of various models and features. Finally draw a conclusion which modelling technique and feature set is best suited for sarcasm detection for Hinglish language.

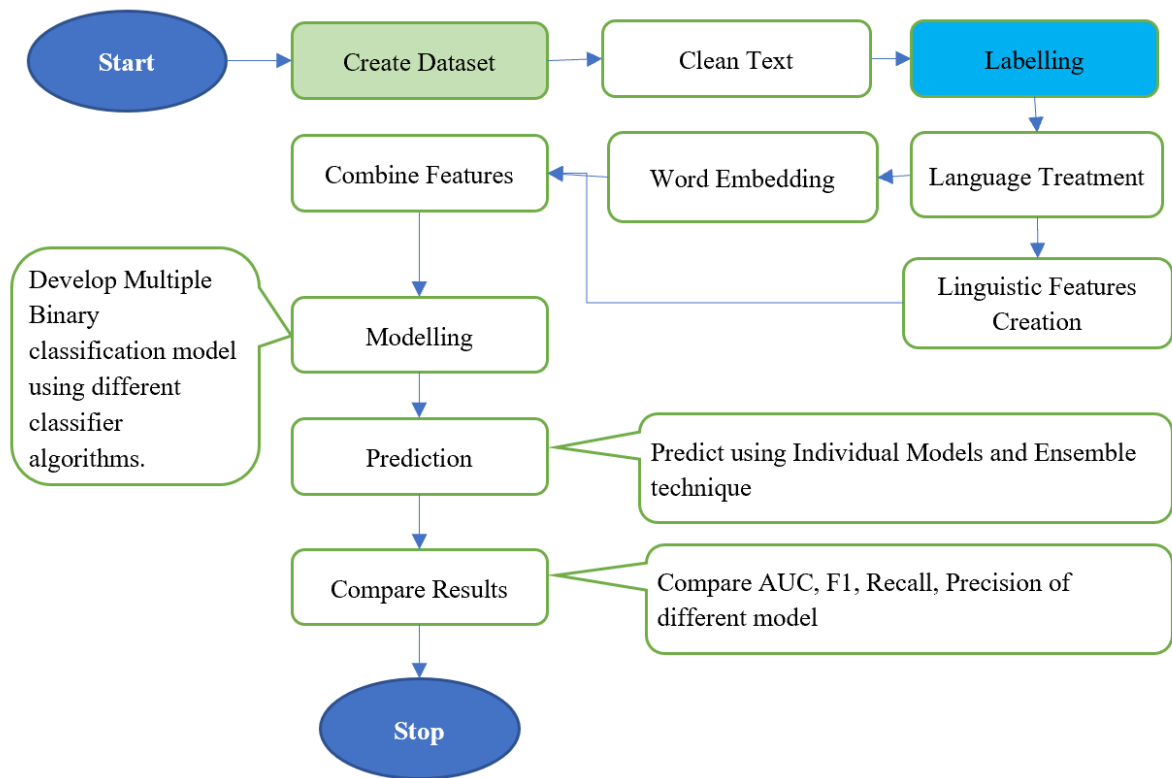


Figure 5: Overall Approach

3.3. Model Building

3.3.1. Test-Train Split

We will use train-test split of 80:20.

3.3.2. Handling Small Dataset size

Our dataset has 2300+ records. This is small size dataset. Because dataset is not large enough therefore, we will use cross validation of 5 folds. If time permits, we will also explore text oversampling techniques and build the models using the oversampled data.

3.3.3. Algorithms, Architecture for Modeling

1. Logistic Regression
2. Stochastic Gradient Descent
3. Gradient Boosting
4. Random Forest
5. Adaboost
6. SVC
7. CNN

8. BERT (TL)
9. GPT2 (TL)
10. fastText (TL)

3.4. Evaluation Metrics & Reporting

3.4.1. Evaluation Metrics & Reporting

ROC graphs are useful tool for visualizing and evaluating classifiers. ROC are able to provide a richer measure of performance than accuracy or error rate (Fawcett, 2004). From the Appendix B we can notice that most of the researchers either used Accuracy or F1 score to measure the performance of the sarcasm detection system or sentiment analysis. However, we will also use AUC, Accuracy, F1, Recall & Precision, because they have their relevance depending upon the domain where we use this for sarcasm detection. To understand it better, let's see a sarcasm from hospital, health domain.

A patient says "Hospital administration thinks that I come to hospital because I have lot of money and they have beautiful nurses to chat with" (writing sarcasm in English to make sure more readers understand the impact of choice of evaluation metrics).

Healthcare domain, hospital administrators would like to take a sarcasm seriously and they do not want any sarcasm to be misclassified and they are ready for more False-True (which our system identify sarcastic but in reality they are not). To illustrate the choice of metrics, let's assume there are 1000 sentences in the real time dataset, 150 are sarcasm and 850 are normal sentences. Let us say Model1 predicts 110 are sarcasm and 890 normal and Model2 predicts 140 sarcasm and 860 normal sentences. Let's say accuracy of both the models is 90%. If we select Recall and F1 score, then Model1 is better. If we select precision, then Model2 is better. If we need to detect sarcasm in comment box of YouTube channel of some political party, then we can go for Model1 which is giving recall of 73%. If we are dealing with some more serious product or service like healthcare, airlines service then we can go for Model2 which is giving Precision score of 63%.

| Model1 | | | | | Model2 | | | | |
|------------|-------------|-------|------|------|------------|-------------|-------|------|------|
| Actual | Observation | | | | Actual | Observation | | | |
| | | FALSE | TRUE | | | | FALSE | TRUE | |
| | FALSE | 820 | 30 | 850 | | FALSE | 805 | 45 | 850 |
| | TRUE | 70 | 80 | 150 | | TRUE | 55 | 95 | 150 |
| | | 890 | 110 | 1000 | | | 860 | 140 | 1000 |
| Accuracy | | | | 0.90 | Accuracy | | | | 0.90 |
| Recall | | | | 0.73 | Recall | | | | 0.68 |
| Precesion | | | | 0.53 | Precesion | | | | 0.63 |
| F1 Score | | | | 0.81 | F1 Score | | | | 0.77 |
| Error Rate | | | | 0.10 | Error Rate | | | | 0.10 |

The result of prediction will be compared using AUC, F1, Accuracy, Recall, Precision.

3.4.2. Reporting Experiments Format

In the final thesis we will move this section to chapter 5. At this time, we are keeping a template here so that we know how we need to structure the flow of experiments and result collection. This table below gives an idea what data will be collected. For each metric AUC, F1, Accuracy, Recall, Precision performance of the model will be reported in following format

| | Linguistic Features | TF-IDF | fastText | BERT | Linguistic + Best Embedding |
|-----------------------------------|------------------------|--------|----------|------|--------------------------------|
| Logistic Regression | | | | | |
| Stochastic Gradient Descent | | | | | |
| Gradient Boosting | | | | | |
| Random Forest | | | | | |
| Adaboost | | | | | |
| SVC | | | | | |
| CNN | | | | | |
| BERT (TL) | | | | | |
| GPT2 (TL) | | | | | |
| fastText (TL) | | | | | |

3.4.3. Result Comparison Format

In the final thesis we will move this section to chapter 5. At this time, we are keeping a template here so that we know what data is needed to prepare a summary of the metrics.

This table below gives an idea what data will be collected.

| | AUC | F1 | Accuracy | Recall | Precision |
|---|-----|----|----------|--------|-----------|
| Logistic Regression + Best Feature Set Name | | | | | |
| Stochastic Gradient Descent + Best Feature Set Name | | | | | |
| Gradient Boosting + Best Feature Set Name | | | | | |
| Random Forest + Best Feature Set Name | | | | | |
| Adaboost + Best Feature Set Name | | | | | |
| SVC + Best Feature Set Name | | | | | |
| CNN + Best Feature Set Name | | | | | |
| BERT (TL) + Best Feature Set Name | | | | | |
| GPT2 (TL) + Best Feature Set Name | | | | | |
| fastText (TL) + Best Feature Set Name | | | | | |

Based on this table conclude which model with which set of features gives best results on 5 metrics (AUC, Accuracy, F1, Recall, Precision).

3.5. Development Tools

Language: Python 3.0>

ML Libraries

- Matplotlib
- Seaborn
- Pandas
- Numpy
- Sklearn
- RE

Indian Language Libraries

- NLTK
- iNLTK

Word Embedding

- TF-IDF
- FastText
- Word2Vec

Classical Modelling

- Logistic Regression
- Gradient Boost
- Decision Tree
- Random Forest
- Stochastic Gradient Boost
- Adaboost
- SVM
- Naive Bayesian

Transformer & Architecture

- CNN
- BERT
- GPT2

Framework

- PyTorch
- Tensorflow
- Fast.ai
- Keras

3.6. Summary

We will develop a dataset of 2000+ sentences. Some text will be taken from twitter and some other will be taken from Hindi blogs. Data will be cleaned and labelled with the help of native speakers. For creating features of the dataset, we will extract linguistic features from the sentences. We will also use word embedding like TF-IDF, word2vec, fastText, BERT to create features. For developing models, we will use classical machine learning models like LR, GB, DT, RF, SGB, Adaboost, SVM & NB. We will also explore CNN and transformers like BERT & GPT2. For measuring model performance, we will use 5 metrics F1, AUC, Accuracy, Recall,

and Precision. With all these experiments we will present our finding which type of features, which word embedding, which model works best of different metrics.

REFERENCES

- [1.] Anggraini, S.D., (2014) *A Pragmatic Analysis Of Humor In Modern Family*.
- [2.] Asghar, M.Z., Kundi, F.M., Khan, A. and Ahmad, S., (2014) Lexicon-Based Sentiment Analysis in the Social Web. *J. Basic. Appl. Sci. Res*, 46, pp.238–248.
- [3.] Bharti, S.K., Babu, K.S. and Raman, R., (2018) Context-based Sarcasm Detection in Hindi Tweets. *2017 9th International Conference on Advances in Pattern Recognition, ICAPR 2017*, pp.410–415.
- [4.] Bharti, S.K., Sathya Babu, K. and Jena, S.K., (2017) Harnessing Online News for Sarcasm Detection in Hindi Tweets. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 10597 LNCS, pp.679–686.
- [5.] Van Brunt, J., (1987) A closer look at. *Bio/Technology*, 511.
- [6.] Clark, K., Luong, M.-T., Le, Q. V. and Manning, C.D., (2020) ELECTRA: Pre-training Text Encoders as Discriminators Rather Than Generators. [online] Available at: <https://github.com/google-research/> [Accessed 28 Aug. 2020].
- [7.] Devlin, J., Chang, M.W., Lee, K. and Toutanova, K., (2019) BERT: Pre-training of deep bidirectional transformers for language understanding. *NAACL HLT 2019 - 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies - Proceedings of the Conference*, 1Mlm, pp.4171–4186.
- [8.] Faías, D.I.H., Patti, V. and Rosso, P., (2016) Irony detection in twitter: The role of affective content. *ACM Transactions on Internet Technology*, [online] 163, pp.1–24. Available at: <http://dx.doi.org/10.1145/2930663> [Accessed 17 Aug. 2020].
- [9.] Fawcett, T., (2004) ROC graphs: Notes and practical considerations for researchers. *Machine learning*, 311, pp.1–38.
- [10.] Gaikwad, V. and Haribhakta, Y., (2020) Adaptive glove and fasttext model for Hindi word embeddings. *ACM International Conference Proceeding Series*, pp.175–179.
- [11.] Joshi, A., Bhattacharyya, P. and Carman, M.J., (2018) Investigations in computational sarcasm. *Cognitive Systems Monographs*, 37, pp.137–143.
- [12.] Kumar, A. and Garg, G., (2019) Sarc-M : Sarcasm Detection in Typo-graphic Memes. In: *International Conference on Advanced Engineering, Science, Management and Technology – 2019 (ICAESMT19) Sarc-M*: pp.1–8.
- [13.] Kumar, Y., Mahata, D., Aggarwal, S., Chugh, A., Maheshwari, R. and Shah, R.R., (2019) *BHAAV- A Text Corpus for Emotion Analysis from Hindi Stories*. Available at: <http://arxiv.org/abs/1910.04073> <http://dx.doi.org/10.5281/zenodo.3457467>.
- [14.] Lan, Z., Chen, M., Goodman, S., Gimpel, K., Sharma, P. and Soricut, R., (2019) ALBERT: A Lite BERT for Self-supervised Learning of Language Representations. [online] Available at: <https://github.com/google-research/ALBERT>. [Accessed 28 Aug. 2020].
- [15.] Lee, C.J., Katz, A.N., Lee, C.J. and Katz, A.N., (2009) The Differential Role of Ridicule in Sarcasm and Irony The Differential Role of Ridicule in Sarcasm and Irony. 6488May 2015, pp.37–41.
- [16.] Liebrecht, C., Kunneman, F. and Bosch, A. Van den, (2013) The perfect solution for detecting sarcasm in tweets #not. [online] June, pp.29–37. Available at: <http://www.aclweb.org/anthology/W13-1605>.
- [17.] Liu, L., Priestley, J.L., Zhou, Y., Ray, H.E. and Han, M., (2019a) A2Text-net: A novel deep neural network for sarcasm detection. *Proceedings - 2019 IEEE 1st International Conference on Cognitive Machine Intelligence, CogMI 2019*, December, pp.118–126.
- [18.] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L. and Stoyanov, V., (2019b) RoBERTa: A Robustly Optimized BERT

- Pretraining Approach. [online] Available at: <https://github.com/pytorch/fairseq> [Accessed 28 Aug. 2020].
- [19.] Nozza, D., Fersini, E. and Messina, E., (2016) Unsupervised Irony Detection: A Probabilistic Model with Word Embeddings. In: *Proceedings of the 8th International Joint Conference on Knowledge Discovery, Knowledge Engineering and Knowledge Management*. [online] SCITEPRESS - Science and Technology Publications, pp.68–76. Available at: <http://www.scitepress.org/DigitalLibrary/Link.aspx?doi=10.5220/0006052000680076> [Accessed 16 Aug. 2020].
- [20.] Peters, M.E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K. and Zettlemoyer, L., (2018) Improving Language Understanding by Generative Pre-Training. *OpenAI*, [online] pp.1–10. Available at: https://gluebenchmark.com/leaderboard%0Ahttps://gluebenchmark.com/leaderboard%0Ahttps://s3-us-west-2.amazonaws.com/openai-assets/research-covers/language-unsupervised/language_understanding_paper.pdf.
- [21.] Potamias, R.A., Siolas, G. and Stafylopatis, A.G., (2020) A transformer-based approach to irony and sarcasm detection. *Neural Computing and Applications*.
- [22.] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W. and Liu, P.J., (2019) Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal of Machine Learning Research*, [online] 21, pp.1–67. Available at: <http://jmlr.org/papers/v21/20-074.html>. [Accessed 28 Aug. 2020].
- [23.] Ramos, J., (2003) Using TF-IDF to Determine Word Relevance in Document Queries. In: *Proceedings of the first instructional conference on machine learning*. [online] pp.133–142. Available at: <https://sites.google.com/site/caonmsu/ir/UsingTFIDFtoDetermineWordRelevanceinDocumentQueries.pdf> [Accessed 29 Aug. 2020].
- [24.] Romano, S., (n.d.) *Multilingual Transformers - Towards Data Science*.
- [25.] Sanh, V., Debut, L., Chaumond, J. and Wolf, T., (2019) DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. [online] Available at: <https://github.com/huggingface/transformers> [Accessed 28 Aug. 2020].
- [26.] Saravia, E., Toby Liu, H.C., Huang, Y.H., Wu, J. and Chen, Y.S., (2020) Carer: Contextualized affect representations for emotion recognition. *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, EMNLP 2018*, pp.3687–3697.
- [27.] Sharma, D.S., Sangal, R., Pawar, J.D., Sharma, R. and Bhattacharyya, P., (2014) A Sentiment Analyzer for Hindi Using Hindi Senti Lexicon. In: *NLP Association of India*. NLP AI, pp.150–155.
- [28.] Sinha, R.M.K. and Thakur, A., (2005) Machine Translation of Bi-lingual Hindi-English (Hinglish) Text. *10th Machine Translation summit (MT Summit X)*, pp.149–156.
- [29.] Subramanian, J., Sridharan, V., Shu, K. and Liu, H., (2019) Exploiting emojis for sarcasm detection. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 11549 LNCS April, pp.70–80.
- [30.] Suhaimin, M.S.M., Hijazi, M.H.A., Alfred, R. and Coenen, F., (2017) Natural language processing based features for sarcasm detection: An investigation using bilingual social media texts. *ICIT 2017 - 8th International Conference on Information Technology, Proceedings*, pp.703–709.
- [31.] Sundararajan, K. and Palanisamy, A., (2020) Multi-rule based ensemble feature

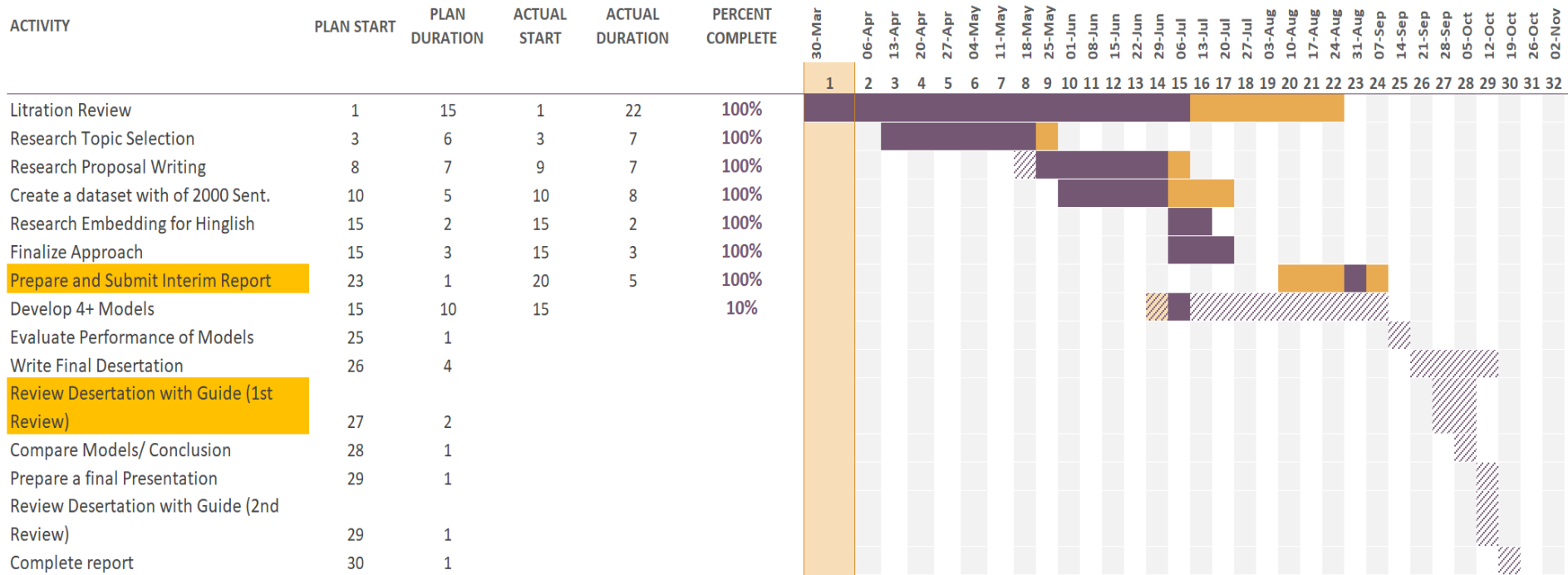
- selection model for sarcasm type detection in Twitter. *Computational Intelligence and Neuroscience*, 2020.
- [32.] Turney, P.D., (2002) Thumbs Up or Thumbs Down? Semantic Orientation Applied to Unsupervised Classification of Reviews. [online] July, pp.417–424. Available at: <http://arxiv.org/abs/cs/0212032>.
 - [33.] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł. and Polosukhin, I., (2017) Attention is all you need. *Advances in Neural Information Processing Systems*, 2017-DecemNips, pp.5999–6009.
 - [34.] Wang, Q., Xu, J., Chen, H. and He, B., (2017) Two improved continuous bag-of-word models. In: *Proceedings of the International Joint Conference on Neural Networks*. [online] Institute of Electrical and Electronics Engineers Inc., pp.2851–2856. Available at: <http://ieeexplore.ieee.org/document/7966208/> [Accessed 29 Aug. 2020].
 - [35.] Wikipage, (n.d.) *Demographics of India - Wikipedia*.
 - [36.] Wikipage, (n.d.) *Internet in India - Wikipedia*.
 - [37.] Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R. and Le, Q. V, (2019) XLNet: Generalized Autoregressive Pretraining for Language Understanding. In: *33rd Conference on Neural Information Processing Systems (NeurIPS)*. [online] Available at: <https://github.com/zihangdai/xlnet> [Accessed 28 Aug. 2020].
 - [38.] Zhang, Y., Sun, S., Galley, M., Chen, Y.-C., Brockett, C., Gao, X., Gao, J., Liu, J. and Dolan, B., (2020) DIALOGPT : Large-Scale Generative Pre-training for Conversational Response Generation. [online] pp.270–278. Available at: <https://github.com/mgalley/> [Accessed 28 Aug. 2020].

Appendix A: List of Other Documents

These are the documents prepared during this project but for the purpose of brevity there are not part of main document. Those who are interested can refer them in github.

1. [History_AutoSarcasmDetection.pdf](#)
2. [Summary-of-Sarcasm-Papers.pdf](#)
3. [Dataset-cleaning-steps.pdf](#)
4. [Datasource-links.pdf](#)

APPENDIX B: Research Plan



APPENDIX C: Sarcasm Detection Systems Results of Past Work.

| Sno | Year | Authors | Model | Features | Metrics |
|-----|------|--------------------------------|--------------|----------|--|
| 1 | 2002 | (Turney, 2002) | Rule based | LFS | Acc: 74.39% |
| 2 | 2009 | (Burfoot and Baldwin, 2009) | Classical ML | LFS | F1: 79.8% |
| 3 | 2010 | (Pak and Paroubek, 2010) | Classical ML | LFS | Not Mentioned |
| 4 | 2010 | (Davidov et al., 2010) | Classical ML | LFS | F1: 78% Amazon |
| 5 | 2010 | (Davidov et al., 2010) | Classical ML | LFS | F1: 83% Twitter |
| 6 | 2011 | (González-Ibáñez et al., 2011) | Classical ML | LFS | Acc: 55.59% to 75.78% depending upon tweet format. |
| 7 | 2013 | (Mittal and Agarwal, 2013) | Rule Based | LFS | Acc: 80.21% |
| 8 | 2013 | (Liebrecht et al., 2013) | Ruled Based | LFS | AUC: 77% |
| 9 | 2013 | (Riloff et al., 2013) | Classical ML | LFS | F1: 51% |
| 10 | 2014 | (Asghar et al., 2014) | Rule based | LFS | Acc: 95.24% |
| 11 | 2014 | (Sharma et al., 2014) | Rule Based | LFS | Acc: 85 to 89.5% |
| 12 | 2015 | (Rajadesingan et al., 2015) | Classical ML | LFS | Acc: 83.46% |
| 13 | 2015 | (Joshi et al., 2015) | Classical ML | LFS | F1: 61% |
| 14 | 2015 | (Bamman and Smith, 2015) | Classical ML | LFS | Acc: 85.1% |
| 15 | 2016 | (Farias et al., 2016) | Classical ML | LFS | Acc: 73-96% depends upon datasets and classifier. |
| 16 | 2016 | (Jha et al., 2016) | Classical ML | LFS | Acc: 92.2% to 100% depending upon unigram or bigram feature and classifier |
| 17 | 2016 | (Desai and Dave, 2016) | Classical ML | LFS | Acc: 84% |
| 18 | 2017 | (Suhaimin et al., 2017) | Classical ML | LFS | Acc: 82.5% |
| 19 | 2017 | (Bharti et al., 2017) | Rule Based | LFS | Acc: 79.4% |
| 20 | 2017 | (Ravi and Ravi, 2017) | Classical ML | LFS | F1: 96.58% (L+T+D features) + GR feature selector + SVM RBF Classifier |
| 21 | 2018 | (Bharti et al., 2018) | Rule Based | LFS | Acc: 87% |
| 22 | 2018 | (Parde and Nielsen, 2018) | Classical ML | LFS | F1: 59% (Twitter) |

| | | | | | |
|----|------|-------------------------------------|-----------------|-----------|--|
| 23 | 2018 | (Parde and Nielsen, 2018) | Classical ML | LFS | F1: 78% (Amazon) |
| 24 | 2018 | (Swami et al., 2018) | Classical ML | LFS | Acc: 78.4% with RF |
| 25 | 2018 | (Van Hee et al., 2018) | Classical ML | LFS | Acc: 67.54% (SVM) |
| 26 | 2018 | (Van Hee et al., 2018) | Classical ML | LFS | Acc: 68.27% (LSTM) |
| 27 | 2019 | (Kumar and Garg, 2019) | Classical ML | LFS | Acc: 73.25% to 87.95% depending upon the classifier used. |
| 28 | 2019 | (Kumar et al., 2019) | Classical ML | Both | Acc: 62% |
| 29 | 2019 | (Subramanian et al., 2019) | GRU | LFS | F1: 89.36% (Twitter) |
| 30 | 2019 | (Subramanian et al., 2019) | GRU | LFS | F1: 97.97% (facebook) |
| 31 | 2019 | (Liu et al., 2019a) | Classical + CNN | LFS | F1: 71% - 90% depending upon dataset with A2Text classifier |
| 32 | 2020 | (Zhang et al., 2020b) | CNN | LFS | Acc: 94.6% on Large dataset of SST2 |
| 33 | 2020 | (Sundararajan and Palanisamy, 2020) | Classical ML | LFS | Acc: 86.61% to 99.79% Depending upon the type of sarcasm. Final classifier is RF |
| 34 | 2020 | (Castro et al., 2020) | Classical ML | LFS | F1: 71.8% |
| 35 | 2020 | (Potamias et al., 2020) | Transformer | Embedding | Acc: 85% to 94% depending upon dataset |
| 36 | 2020 | (Saravia et al., 2020) | CNN | Both | Acc: 81% with CARER |

APPENDIX D: Research Proposal

1. Introduction / Overview

Mobile phones came to India in 1995¹⁸ and Internet was launched in India by VSNL in 1995¹⁹. Initially the cost of the technology was extremely high, so it was available only to business class, research labs, high level bureaucrats and politicians. With the increase of literacy and decreasing cost of internet services and mobile phone device internet, it is so common that people started thinking that Internet is our fundamental right. As per the World Economic Forum (WEF), in 2019, about 60% of Indian internet users viewed content in vernacular. WEF also says 75% of this 60% is below 35 years of age (Wikipage, n.d.). According to the same Wikipedia page, by 2030, 1.1 billion Indian will have access to Internet and 80% will access the content on mobile devices. The WEF also estimated that 80% of the users will be consuming content in vernacular languages.

When Government of India is going for full blown Digital India program and bringing every citizen of India on the internet platform for purchase, payment and government fund transfer then how the citizens are going to provide feedback about the services which they use? As of today, it is easier to perform sentiment analysis of the feedback given in English but feedback given in Hindi is not easy to analyse. It means voice of Hindi speaking people is not being considered in service improvement. Till the time somebody is not too angry and do some crime or come on the road to do Dharana or protest we do not know what is happening and why.

Many Hindi new portals, book, blogs, chat bot/WhatsApp conversations, YouTube channels, Twitter & Facebook pages are full of content in Hindi language. People openly express themselves online using Hinglish language which is mix of Hindi, English, Urdu and other languages. Volume of the online content is increasing at unprecedented rate and it is responsibility of government,

¹⁸https://en.wikipedia.org/wiki/Telecommunications_in_India#:~:text=In%20August%201995%2C%20then%20Chief,launched%20in%20Kolkata%20in%202012. (Accessed 24-Jun-20)

¹⁹https://en.wikipedia.org/wiki/Internet_in_India#:~:text=The%20first%20publicly%20available%20internet,not%20permitted%20in%20the%20sector. (Accessed 24-Jun-20)

business community, professionals, NGO and others to understand the feeling of public and respond accordingly. But the biggest challenge is how to analyse the content which is written in mix of Indian languages. It is impossible to analyse the Hinglish language text manually or using traditional systems.

This section is organized as 1.1 What is Hinglish, 1.2 Origin of Hinglish, 1.3. What is Sarcasm?, 1.4. Why Sarcasm Detection is Critical?, 1.5. Why Sarcasm Detection is Critical in Electronic Media? , 1.6. Sarcasm Detection in Hindi, 1.7. Challenge in Processing Hinglish, 1.8. Common Challenges in Sarcasm Detection, 1.9. Context Understanding a Challenge in Sarcasm Detection, 1.10. Challenges in Sarcasm Detection in Hinglish, 1.11. Degree of Sarcasm, 1.12. Positive Side of Hinglish

1.1 What is Hinglish?

There was time when Hindi was a language which is used by majority of Hindi speaking people when they are communicating (writing, speaking) with each other. But in 21st century, most of the Hindi speaking population who express themselves on social media use Hinglish language. Hinglish is a new lingo of Hindi speaking population. Hinglish sentences follow Hindi grammar and most of the word are taken from Hindi but there is no hesitation of taking words from other languages like English, Urdu etc. Hinglish language spoken by different people have different amount of words from different languages. For example, those people who know Urdu good enough for them Hinglish is mix of Hindi, Urdu, English. Those who know Avadhi for them Hinglish is mix of Hindi, Avadhi, English. Those who know Marathi very well for them Hinglish is mix of Hindi, Marathi, English. Thus, in Hinglish Language we have words from Hindi, English and various other Indian languages and written in Devanagari & Roman together.²⁰ (Sinha and Thakur, 2005) Hindi and English language mixed is called Hinglish. Hinglish is not limited to Hindi & English mix but it includes Punjabi, Gujarati, Marathi, Urdu. Phrase construct happens in Roman and Devanagari script.²¹

²⁰ Latin is Region and Rome is part of that reason. Over the period of time Roman empire become famous and script was called Roman but Latin is also used simultaneously. <https://www.quora.com/Why-is-the-language-of-the-ancient-Romans-called-Latin-and-not-Roman> (Accessed 28-Jun-20)

²¹ <https://en.wikipedia.org/wiki/Hinglish> (Accessed 24-Jun-20)

1.2 Origin of Hinglish

Before Internet Era in India people use to communicate with each other in much cleaner format of the language and there was not much mix of other language or English and for writing Hindi they were using Devanagari script. But, with the penetration of internet in the society a new language started taking shape. Initially when Devanagari keyboards were not available people were using Roman letters to write Hindi email, SMS.

An example of late 20th century text in Hinglish language. “Main is doorbhash ka prayog karna nani janta”. This is Hindi in Roman script. We need to keep in mind that people do not follow any IAST or other map for writing Hinglish letters in Roman. Mobile phone and Internet were available to elite, educated journalist, professionals. They started realising they are typing in Roman but some words in English so translating them and then typing in Roman is painful. So, text became like this “Main is phone ko use karna nahi janta”. Roman script with Hindi and English words.

Over the period of time when Devanagari keyboards were easily available people started using Devanagari keyboards for writing Hindi, but by that time so much English has come in day to day conversation that they felt it is uncomfortable to use Hindi words. So, they write like this.

“मैं इस फोन को यूज करना नहीं जानता”. Devanagari script with Hindi and English words. Over the period of time people started realizing it is becoming difficult to know which word is Hindi and which one is English therefore a word which come from English root should be written in Roman and word which are from Hindi root should be written in Devanagari. So, they started writing like this. “मैं इस phone को use करना नहीं जानता”. Devanagari & Roman mixed for Hindi and English words.

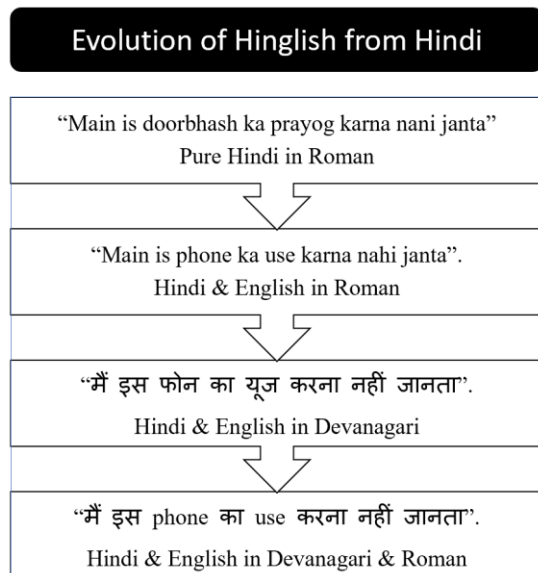


Figure 1: Evolution of Hinglish

Today if you read any Hindi speaker’s WhatsApp, twitter or Facebook message you will find they use words from different Indian languages like Urdu, Marathi, Bangla, Punjabi and write either in Devanagari or in Roman. “अमी मौजूलिका. अमी राजा को जरूर मारबो 😏”, but why you want to kill him?”. Here Hindi, Bangla, Urdu and English 4 languages used along with emoticon and written in two scripts Devanagari and Roman. This is Hinglish.

Today Hindi social media, Hindi comment boxes of product, Hindi news articles are full of this kind of language, Hinglish. Therefore, this work using Hinglish language is high value from the angle of practical usage.

1.3 What is Sarcasm?

Your friend come to you and speak something to you, from the tone of his language, his body language, choice of his words, time and situation he is speaking you realised that the real meaning of what he is saying is completely opposite. It may be easier for you to detect this opposite sense if you are aware about the complete context but if you are not aware about the context then even as intelligent human you may miss the real meaning of what is being said.

For example, you open the door for your friend, and he says wow! your looking handsome in this T-shirt. You know that this is an old T-shirt and many times your friend has seen this. But still not aware of full context, you hesitantly say thank and you invite him inside. After 15 minutes you check yourself in the mirror and realised that you are wearing T-shirt flip side. Now you are embarrassed for your “Thank you” response.

What your friend did was sarcastic remark on your dressing and you being unaware of the full context could not respond properly. In the absence of full context, understanding sarcasm is difficult task and most of the time we take literal meaning of the words or some other time get confused that why someone has made that remarks which was completely out of the context.

In English language this type of grammatical construct which has completely opposite meaning than what is said, it called sarcasm.

As per merriam-webster dictionary, sarcasm is²²

1: a sharp and often satirical or ironic utterance designed to cut or give pain

2a: a mode of satirical wit depending for its effect on bitter, caustic, and often ironic language that is usually directed against an individual

2b: the use or language of sarcasm

In Hindi it has several name and synonyms like कटाक्ष (Kataksha), तंज (Tanja), व्यंग/ व्यङ्ग (Vyanga), टोंट (Tonta)

Ten forms of humour are irony, satire, sarcasm, overstatement, self-deprecation, teasing, replies to rhetorical question, clever replies to serious statements, and transformations of frozen expressions. All these are functions of humour and found in the sitcom (situational comedy). What one finds hilarious or pun may be completely opposite to another person in another country or in other situation. Interpretation is filtered by cultural context. (Anggraini, 2014)

²² <https://www.merriam-webster.com/dictionary/sarcasm>

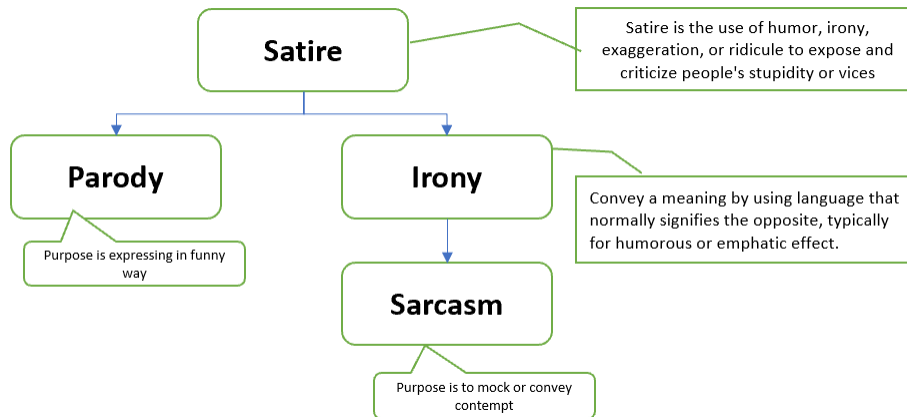


Figure 2: Sarcasm & Satire Relationship

In their work “The Differential Role of Ridicule in Sarcasm and Irony” (Lee et al., 2009) says sarcasm and irony are similar because they are both form of reminder yet they are different because sarcasm is about ridiculing a specific person however this is not required in case of irony. However, in our work we will ignore this specific aspect. There are two reasons for that 1- we are interested in predicting whether the statement is conveying real meaning or opposite meaning of what is being said 2- In Indian and specific to Hinglish context words like कटाक्ष or व्यंग doesn’t consider the aspect mentioned by the authors.

1.4 Why Sarcasm Detection is Critical?

If we do not understand the real intent of the speaker then we cannot respond him properly. Response can be physical action or verbally reply to the speaker or even no action.

Few examples where not understanding the real intent of the person can be catastrophic.

- In face to face communication with your customer when you miss his intent. Result is customer disengagement.
- In live program when you are listening a response or question from the audience in hall or live TV or Radio program or speaking over phone or video conferencing tool and you miss the intent. Result is dent on your reputation.
- In offline communication when you publish some content on blog, news, product selling page and receive some comment from the public. Someone expresses his opinion over your post or tweet and you are not able to understand that properly or not able to read. All other

people read that comment and think that either you are dumb or do not care or accept what is being said. Result you know very well.

When you are dealing with your known people, friends, relatives and not responding properly in that situation, it will have lessor impact because they know your real nature and potential. But in public places, where you do not know the person to whom you need to respond, can cause huge dent on your image and brand.

1.5 Why Sarcasm Detection is Critical in Electronic Media?

With the advancement of online sales of products, social media and online blogs, new portals there is huge surge of online feedback. Post COVID19 pandemic there are clear trends of shifting in this direction. People prefer buying, reading, expressing, engaging online. This justifies the need of sophisticated real time sarcasm detection system.

1.6 Sarcasm Detection in Hindi

English is 3rd most spoken language in the world and many researchers across the world are working for sarcasm detection in English. But, Hindi is 4th most spoken language in the world and not much significant work is happening in sarcasm detection in Hindi. Due to this reason many of the feedback given on twitter, Facebook, product page, online news goes unnoticed.

Sarcasm is one kind of feedback and if we do not use this to improve our response then we prove ourselves foolish and customer shift to different product, service or platform. Similar things happen when people change their party or group. Therefore, we feel it is extremely important to detect the sarcastic feedback given by those people who write in Hindi.

1.7 Challenge in Processing Hinglish

A. Complexity due to English words in Hindi

Observe the variation of a sentence “I have purchased tickets”

मैंने (टिकिटें/ टिकटें/ टिकटे/ टिकिट) खरीद (ली/ लीं) (है/हैं). This simple sentence can be spoken in 16 different ways if written in Devanagari. If we mix Roman script in between then number of permutations goes beyond our normal imagination. Here we need to make note that

Ticket is English word, and people are making plural of that as they do with any Hindi word.

Let us see another sentence “She has boiled the rice”

उसने राइस बोइल कर दिया है

From the above Hinglish sentence you cannot figure out whether the doer is female or male. Secondly, राइस and बोइल are not words in Hindi dictionary. Sometime people will write letter in Roman like

उसने Rice बोइल कर दिया है / उसने Rice Boil कर दिया है / उसने राइस Boil कर दिया है /

उसने Rice बोयल कर दिया है

Like Guru, Karma are Hindi words and they are part of English dictionary. We do not have Hinglish dictionary which has word like यूज गुड नाइस क्वीन in that dictionary. Without transliterating words like Tickets, Boil into Devanagari and telling system that टिकिटें = टिकटें = टिकटे= टिकिट, बोइल= बोयल embedding will not give good results.

B. Mix Other Indian Language with Hindi

Observe the sentence below, Bangla written in Devanagari and clearly understandable by any Hindi speaking person. Most of the words in the sentence below are from Bangla language but written in Devanagari.

अमी मौंजुलिका.अमी राजा को मारबो दीदी ने केजरीवाल को भी पीछे छोड़ दिया. जि तो कमालई कर दओ ददू

India's business film Industry in Mumbai make film in Hindi. Rarely any film use as good Hindi as Hollywood uses English. Adoption of words from other language is not a problem. The problem is quantity of the words taken from other languages, availability of the updated vocabulary of the language. Many famous dialogues or songs from Hindi films which are taken different language or dialects. This increases complexity of sarcasm

detection in Hinglish. We do not have comprehensive dictionary which we can call Hinglish dictionary which has all the word being used by the Hinglish speakers.

Without telling system that अमी (Bangla word) = मैं, मारोबो (Bangla word) = मारुंगी = मारुंगा = मारना no embedding is going to help

C. Complexity of Synonyms in Hindi

For this let's understand what Synonyms is. A word or phrase that means exactly or nearly the same as another word or phrase in the same language²³, for example “shut” is a synonym of “close”. Few examples of synonyms

- The East = The Soviet Union (<https://www.lexico.com/en/definition/synonym>)
- Country of rising sun = Japan, Dragon Country = China,
- Fridge = Refrigerator
- Happy = Joyful, Cheerful, Contented, Jolly, Gleeful, Carefree

All the synonyms have different spelling, different pronunciation but almost same meaning and part of the same language. l'eau (French word for water) is not synonyms of water because they are two different languages.

Unlike other world languages, all Indian languages (except Tamil, this is debatable) heavily borrow words from Sanskrit.

Let's take English word “Water” and see how many words are available in sanskrit for “water” जल = पानी = तनि = नीरु = आपः = वाः = वारि = सलिलं = पयः = तोयं = मेघपुष्पं = घनरसः = पाणी. So all these words are synonyms of water in sanskrit.

Because all Indian languages have root in Sanskrit therefore most of the time, they take word from Sanskrit for communication. For example, Kannada uses नीरु, Bangla use पानी, Hindi uses पानी, सलिलं, मेघपुष्पं. If not regular, they are used in poetical or sometimes in sarcastic language. Because in sarcasm or poetry we often use loaded words.

²³ <https://www.lexico.com/en/definition/synonym>.

In Hindi language, can we say नीरू is synonym of पानी? No, because नीरू word is normally is used in Kannada and Sanskrit and not in Hindi. As per the definition of synonym another equal word should be from the same language and we know Hindi is not Kannada nor it is Sanksrit. The answer is yes also; because Sanskrit being mother of Hindi language, it borrows words freely from Sanskrit. Thus, we see synonym in Hinglish is not the way it is understood in the context of English.

Therefore, to be build a complete Hinglish dictionary we have take words from all other Indian languages and frequently used English words as well. Thus it should be like this.

जल = पानी = तनि = नीरू = आपः = वाः = वारि = सलिलं = पयः = तोयं = मेघपुष्पं = घनरसः = वाटर

D. Variation in Spelling of Same Word

In Hindi same word spoken and written with different spelling. Observe the spelling of the same word how they are varying. This kind of problem we do not have in English. As discussed earlier, synonym of Happy is Jolly. They both are not same, neither in spelling, nor in pronunciation, nor in full sense, but “happy” is close to “jolly”. That is why they are synonyms. But below all “=” signs are referring to the same thing.

विष्णु = बिश्णु = विश्णु = बिष्णु = विष्नु = बिष्नु,

दरसन= दर्शन= दर्सन = दरशन

करता = कर्ता,

यज्ञ = जग्य,

योग = जोग,

हरि=हरी,

We need to keep in mind Hindi is not Devanagari, nor Hindi is Avadhi or Marathi. Hindi is written in Devanagari script but it is heavily inflected by other languages like Awadhi, Bhojpuri, Rajasthani, Urdu etc.

Unless we have dictionary which tells विष्णु = बिश्णु = विश्णु = बिष्णु = विष्णु = बिष्णु, no embedding will help.

1.8 Common Challenges in Sarcasm Detection

Detecting Sarcasm is difficult if sentences are having following characteristics.

- E. **Idioms and Phrases:** Sarcasm detection become more difficult when people speak in idiomatic language. For example: What a wise man! what he did is nothing other than an axe to grind.
- F. **Speaking with Hint:** When people do not talk directly and use examples which are completely different than context. For example: You are behaving like Mir Jafar.
- G. **Culture:** Different languages have different degree of challenges in sarcasm detection. For example, English is spoken all over the world but the way American express their feeling is different than the way British express. The reason for that is the work and social culture of England and United States is hugely different. In English language what is call sarcasm in England may be considered a normal statement or abusive in US and vice versa.

1.9 Context Understanding a Challenge in Sarcasm Detection

Since the time human child take birth, baby has environment to learn from. Various types of formal or informal environment, social or business or cultural background forces human to think and learn. Either at physical or emotional or intellectual level if human fail to learn then his survival is challenged by the nature around. In this kind of environment, it is easy for any human to understand the context. If we are alert and interested in the topic then we need not to struggle hard to understand the context. But context understanding is extremely difficult in the case of Machine learning. Let us analyze one sarcastic tweet. “#JIO का सच नीता अंबानी ने मन्नत मांगी थी कि अनंत अम्बानी अपना वजन कम कर लेगा तो गरीबों में 3 महीने Net or call का भंडारा करवाऊँगी”

People living in India can understand that this is sarcasm. Because we know the full context. That

- Mukesh Ambani is owner of #Jio
- Neeta Ambani is Wife of Mukesh Ambani
- Anant Ambani is son of Neeta Ambani
- Anant Ambani has 200+ Kg body weight
- Normal body weight of human is around 70 kg
- Anand Ambani is overweight as per the normal standard

- Neeta Ambani desired that her son should have normal weight
- #Jio has launched 3 Month free internet package
- There is no direct connection between Anand Ambani weight reduction and 3-month free internet package

(Joshi et al., 2018) in their work “Investigation on Computational Sarcasm” says there are three type of context, Author Specific context, Conversational Context, Topical Context

We need to understand that keeping all the facts in mind we can say a statement is sarcasm and not normal statement. Even a human, who does not have all this information will fail to classify a statement as sarcasm. It is not easy to give all this information to a system to make a classification decision

1.10 Challenges in Sarcasm Detection in Hinglish

- A. 70% of the world population uses 26 letters of Roman script to write their language. The Roman alphabet is also used as the basis for the International Phonetic Alphabet, which is used to express the phonetics of all languages.²⁴ Due to this reason when people are writing different language like English, French, Indonesian, Tagalog, German, Turkish they need not to change much around the letters, so most of the cases script remain Roman. This advantage is not available to Devanagari script and Hindi language.
- B. An average westerner knows and speaks one language so written and verbal expression most of the time is that one language. An average Indian speaks minimum 2 languages, one is language of his state, plus national language, or English. In southern part of India, it is not uncommon when you find a taxi or truck driver who can speaker 3 or 4 languages but they cannot speak in English. This, one language- one script, advantage is not available for any Indian and they communicate in multiple language without realising that they have shifted language and borrowing words from different language.
- C. While typing feedback people write @account_name. Most of the time @account_name are proper name and written in Roman like @harithapliyal, @eating_point, @banarasi. Similarly, hashtag, which helps us understanding the context of the feedback, is also

²⁴ <https://www.worldatlas.com/articles/the-world-s-most-popular-writing-scripts.html> Accessed on 23-Jun-20

written in Roman script #Election2019 #COVID19 #Philosophy #Motivation
#NarendraModi.

D. Numerals: Many times, people use non English numerals like १, २, ३, ४, ५.

1.11 Degree of sarcasm

Although how a person perceive & responds to a sarcasm it also depends upon him, yet we need to know all sarcastic statements are not equally intense or powerful to generate pain to the listener or reader. Here are few examples of different degree of sarcasm.

- E. ओ भाई कचोरी समोसे की दुकानें खुल तो गयी है लेकिन ध्यान रखे कचोरी समोसे के चक्कर में आप की ही पूड़ी सब्जी न बट जाये #Covid_Unlock (Least Intense)
- F. NDTV की हैडलाइन एक बेजुबान अल्पसंख्यक भैंस को डूबा कर मारने की कोशिश करती बहुसंख्यक चिड़िया (Lessor intensity)
- G. करोना का दवा न होना यह एक साइंस है, और दवा न होते हुए भी बिल लाखों में आना ये एक आर्ट है !! (Moderate Intensity)
- H. ये शुक्र है जंगल में आरक्षण नहीं, बहोत नहीं तो जंगल का राजा शेर नहीं गधा होता. आरक्षण खत्म करो 70 साल हो गये यार #आरक्षण_भीख_है (Sharp Intensity)

1.12 Positive Side of Hinglish

Although India is big country with 1.35 billion people with different culture, religion, tradition but there is some common aspect in India culture and this does not change no matter where a Indian is living on the earth. That common culture helps us understanding the context and intent easily. Although there are many languages in India but because of one overarching culture it is easier to understand the meaning, a simple translation is good enough. Unlike English where Australian struggle to understand what American gentlemen want to say in English.

2. Background and Related Work

(Bharti et al., 2017) Sarcasm detection is one of the most complex work in Hindi Language and the reason for that is words in Hindi language are rich in morphology. This paper discusses a system to sarcasm detection in Hindi tweet but for that it is taking help of online news related to the tweet. This work demonstrates accuracy of 70.4%

Let us take one English verb “do”, in Hindi, it can be used like कर्ता (noun) , करता (verb with male), करती (verb with female), करूंगा (future tense with male), करूंगी (future tense with female), किया (done), करो (must do) करें (please do) etc. these all are with different gender, mood and tenses. However, in English we have inflection like do, does, did, done.

Now, let’s take another example but this time we take noun “Ram”. राम का, राम ने, राम को, राम द्वारा, राम में, राम पर, राम के लिए राम पर and many times you will see letters are written together. We never see any word like “ByRam” in English but in Hindi रामने and राम ने both have same meaning.

Sarcasm is the major factor which can flip the meaning of a written or spoken phrase. To avoid the negativity people use positive words to communicate negative message. (Desai and Dave, 2016). They have used libsvm algorithm for multiclass classification. This paper uses 5 grades of sarcasm Non-Sarcastic, Mild Positive Sarcastic, Extreme Positive Sarcastic, Mild Negative Sarcastic and Extreme Negative Sarcastic. This work demonstrate the accuracy between 60% to 84% depending upon, whether sentence has any clue like emoticon, tag etc of sarcasm. This work suggests usage of lexical, pragmatic, and linguistic features along with emoticons, hashtag, punctuation marks to detect the sarcasm.

(Liebrecht et al., 2013) developed a sarcasm detection system. This was system was developed for tweets in Dutch language. They used 78,000 sarcastic tweets, along with normal tweets dataset, while adding normal tweet ensured that none of the normal tweet is part of sarcastic dataset. Split the sarcastic tweets into train-test and added with normal tweet into train dataset to train the model. Then test the model using test dataset which has only sarcastic tweets. There experiments leads to AUC of .79. This paper gives an overall approach of building sarcasm detection system in other than English language. But it does not address the problem which Hinglish language has. There test train split and model training approach looks good for non-English language.

(Asghar et al., 2014) developed a system to detect negative, positive, and neutral sentiments for English language tweets. As claimed by the authors their system can detect and score the slang

used in the tweet. This system has Accuracy of 92% for binary classification and 87% for multinomial classification. An approach to get tweets clean text is discussed for English language tweets. However, we need to look what extra we need to do for Hinglish language tweets.

(Turney, 2002) presents an unsupervised learning-based algorithm for classification of review in English language. Semantic Orientation (SO) is used to perform this work. SO of a phrase is calculated using adverbs and adjectives used in the phrase. The experiments were done for text of various domains like automobiles, banks, movie review and travels. The results of this experiment vary from domain to domain between 66 to 84%. The power of this SO in Hinglish language sarcasm detection can be used and verified.

Lot of work has been in English language sarcasm detection and authors mentioned different challenges in sarcasm detection, although results are not that great as for any other classification problems. Challenges exist because of context understanding, missing context, domain, culture, different words, or expression used by people to flip the meaning etc. There is not much work done in Hinglish Language Sarcasm detection. Hinglish language has a separate set of challenges like mixing script, mixing language, highly morphological words, using same morphology on English language words, meagre size of corpus etc.

3. Research Questions

- E. How sarcasm detection is done by other researchers for English and any other Indian languages?
- F. How sarcasm detection system should be designed when words from more than one scripts are used for communication. For example: “मेरा work पूरा हो गया है”, it has 2 scripts.
- G. How sarcasm detection system works when more than language are used for communicating idea. For example: “मेरा वर्क पूरा हो गया है”, it has 2 languages.
- H. Unlike English, Hindi is highly morphological language how does it influence overall approach? For example, करता, करती, करते all are equivalent to English “do” but depends upon the gender.
- I. Unlike Roman where we write words using consonants and vowels in Devanagari there is an extra concept called Maatra, this is not available in Roman script. For example, word

“Experience” in Roman is written using 5 vowels, 5 consonants. But in Devanagari it is written as “एक्सपिरिएन्स” 2 vowels (ए ऐ), 6 consonants(क् स् प् र् न् स्) , 5 Maatra (ा, ि, ि, े, ा).

How does Maatra of Devanagari influences text processing?

- J. How to do transliteration from Roman to Devanagari? Many options are available for reverse translation. For example “एकीकरण” => “Ekikaran” is easy lot many option there but “Ekikaran” => “एकीकरण” is not easy. Because Hindi speaking population is not aware about IAST²⁵ and nor they use it for transliteration. So confusion is “ra”=> र or र्, n=> न or न् or ण or ण् or ञ or ज् or ड or ड्, ki=> कि or की or क्ि or क्ी or क्इ or क्ई
- K. What kind of feature engineering need to be done when text is in multiple scripts?
- L. When some English word is written in Devanagari i.e. राइस, कुक then how to handle these words because they are not part of normal Hindi dictionary?
- M. Can we use the same approach for other Indian language words in Devanagari i.e. बरमंड (Garhwali word for “brain”), तुस्सी (Punjabi word for “you”), खाबो (Bangala word for “eating”)?
- N. If we create a feature using hashtag and say feature name is “context” then is it good enough to explain the context and produce better result?
- O. Can NER based features help in sarcasm detection?

4. Aim and Objectives

The aim of this research is to propose a model, which can predict sarcasm in a given Hinglish language sentence with highest possible accuracy.

Based on the above primary goal, objectives of this research are as following.

- D. To analyze the existing dataset of 300+ statements, clean it and labels each sentence.
- E. To expand existing dataset, minimum 600%, which can be used for training and testing a sarcasm detection model of Hinglish Language
- F. To determine which embedding technique best suits for Hinglish dataset
- G. To develop a preprocessing pipeline which can handle Hinglish language sentences.

²⁵ https://en.wikipedia.org/wiki/International_Alphabet_of_Sanskrit_Transliteration

- H. To develop models using different algorithm like Naive Bayesian, SVM, Logistic Regression, Recurrent Neural Network, and other. Consider minimum 4 suitable algorithms.
- I. To develop a prediction model using ensemble of different model and check whether performance improves.
- J. To evaluate different models and identify the best model.

To address issue related to the small dataset set we will use cross validation technique. Because we are going to develop this dataset therefore, we will try to create a balance dataset and hence no oversampling technique will be required. But, if we realize that results are not encouraging and we need to expand our dataset then in the interest of time we will put more non-sarcasm sentences and use oversampling technique to balance the dataset.

5. Significance of the Study

We didn't find one place which has done research and can say with conviction that approximately these are the number of Hindi speaker in the world. Different sources reveal different numbers. As per a [lingoda.com](https://blog.lingoda.com/en/most-spoken-languages-in-the-world-in-2020)²⁶ and [babbel.com](https://www.babbel.com/en/magazine/the-10-most-spoken-languages-in-the-world)²⁷ after English and Mandarin Hindi is 3rd most spoken language on earth. It is spoken by 615mn people. As per Wikipedia 176 million people speak Urdu.²⁸

Culture of Hindi speaking population and Urdu speaking population resembles a lot. While speaking or writing Hinglish many words of Urdu are spoken or written unknowingly. Therefore, any sarcasm analysis system in Hinglish will benefit Urdu speaking community as well.

With current trend of increasing online content in Hindi, it is practically not possible to read each and every review, even if you try it is very expensive and not worth work. We know, even one negative feedback or abuse which goes unnoticed can cause huge problem for the brand of the

²⁶ <https://blog.lingoda.com/en/most-spoken-languages-in-the-world-in-2020> Accessed on 22-Jun-20

²⁷ <https://www.babbel.com/en/magazine/the-10-most-spoken-languages-in-the-world> Accessed on 22-Jun-20

²⁸ <https://en.wikipedia.org/wiki/Urdu> Accessed on 22-Jun-20

company, product, or person. Therefore, performing sentiment analysis on every feedback makes a perfect sense and it can be done automatically almost in real time.

Sarcasm is one type of sentiment and we are trying to discuss overall benefits of sentiment analysis keeping Sarcasm at the centre of discussion.

- E. Sentiment analysis has a broad range of applications like understanding whether a feedback is Sarcasm, Warning, Love Emotion, Hate Emotion, Advertisement of some other product, Contradicting statement, Pun, Abuse, Inspiring Quote, Sensational Revelation, Pleasant Surprise, Allegation, Poetry/Dohe/Chands etc.
- F. Government, NGO, religious leaders, product sellers are able to perform the sarcasm analysis against some product, political party, ideology, religion, company etc then they will be able to control the situation in much better way with minimum damage.
- G. Sarcasm analysis can be used to analyse the feedback on airlines service, travel service, bus or taxi service, telecom, health, government service, new articles, personal blog, food delivery, insurance service, personality page, book page are good places where sentiment analysis plays a critical role.
- H. In multinational companies it becomes exceedingly difficult to use humour to communicate the idea, crack joke or sarcasm, even if all the team member can speak English. The reason for that is different cultural background and different level of comprehension of English by non-native speakers. But when Hindi speaking people connect over video, telephonic or chat conversation it is easy for them to use idioms, joke, sarcasm and ensure that idea is understood. There is different kind of joy of working in lesser formal and light-hearted environment. When India people are speaking to each other using Hinglish we can perform sarcasm analysis to know the feeling of the group.

We are writing the examples of motivation in English language so that we can explain how sarcasm detection can help proper response from chatbot, but common use-case remain same.

Motivation in Travel Domain

Passenger: #ac_not_working. I love to get roasted in heat.

Chatbot: Sorry for the inconvenience. Our service engineer will call you.

Motivation in Hospital Business

Attendant: #expensive_treatment. We come to your hospital for this expensive treatment so that we can talk to your cute nurses.

Chatbot: We understand your concern about treatment cost. Our billing manager will call you.

Motivation in Restaurant Business

Customer: Last time, your food was so good that since last 2 days I am taking rest.

Chabot: I am sorry to hear that.

Motivation in Learning Portal

Learner: What a great content. I am still trying to understand the head and tell of that 30 min video.

Chatbot: Sorry, can you please share with us what difficulty you faced ?

Motivation in News Portal

Reader: What a great story! Did you read it after writing?

Chatbot: We are sorry that you didn't like this story.

Motivation in Airlines Business

Traveler: First time in my life I got such a wonderful service from any airlines. I reached to the destination one day before my check-in baggage.

Chatbot: We are sorry to hear that. We hope your baggage reached safe to you.

Motivation in Dialogue Analysis Work

A dialogue from a Hindi Film “Sholey”²⁹

मौसी मेरा दोस्त इतना अच्छा है कि वह शराब को कभी न नहीं बोल पाता। पीने के बाद जुआ खेलना उसकी खूबी है इसमें उसका कोई दोष थोड़ी है मौसी। बस हारने के बाद थोड़ा मारपीट

²⁹ <https://en.wikipedia.org/wiki/Sholay>

करता है और घर में आ के मेरे को गाली देता है। पर मेरा दोस्त दिल का बहुत अच्छा है मौसी आप अपनी बेटी की शादी मेरे दोस्त से पक्की कर दो

This is a pure sarcasm paragraph. These kind of dialogues makes movie interesting.

6. Scope of the Study

- G. This research is not related to any specific domain like philosophy, politics, history, current affair new etc. Rather it is trying to detect sarcasm in day to day informal conversation.
- H. Sarcasm in our communication can be expressed and experienced at Visual (facial express, body language), Vocal (tone, pace of speech, emphasis on certain word) and text (book, newspaper, articles, social media tweets, comments and feedback box on internet. Visual sarcasm is more universal than vocal. Because voice uses language and there are 7000+ languages on the earth so there is no universal vocal language of expressing sarcasm. But pause, pitch, pace, modulation between words, while speaking, are more universal like Visual. In this paper we are deal only with text-based sarcasm.
- I. Only Roman and Devanagari scripts are considered.
- J. Only Hindi and English language words are considered. If heavily used words from other languages which are part of day to social communication, then we will include that in our Hindi vocabulary.
- K. No analysis of degree of sarcasm.
- L. We know to understand the context datetime plays a critical role. And most of the text in the dataset is coming from tweet. Our base dataset does not have datetime. We could have included datetime. But we avoided that intentionally because in future when we are expanding the dataset further, we will extract information from different books and other sources and that time datetime will not be available. We wanted to develop a generic system which can understand the context using hashtag. Hashtag is part of the tweet. And we will be extracting it as a separate feature. We do not want that our system should be depending upon time to understand the context.

7. Research Methodology

In this section we are going to discuss a high-level approach to accomplish the research goal. The flow of discussion in the section is as following 7.1. About Dataset, 7.2. Dataset Structure, 7.3. Handling Small Dataset size, 7.4. Building Dataset, 7.5. Cleaning Text, 7.6. Labelling, 7.7. Transliteration, 7.8. Context Creation, 7.9. Emoticon Handling, 7.10. Embedding, 7.11. Feature Engineering, 7.12. Algorithm, 7.13. Prediction, 7.14. Result Comparison

7.1 About Dataset

Keeping the duration of project in mind, it was recommended that we should use an already existing dataset. We used Hindi tweet dataset.³⁰ This excel file had total 442. During the project planning phase we realised that these 442 are not sarcastic tweet but mix of normal and sarcastic and to determine the sarcastic-ness of a sentence developer of this dataset is using news context and there is not explicit labelling available in the given dataset.

Based on the feedback from research guide we decided to expand the dataset which should have minimum 1000 sarcastic sentences and 1000 normal sentences. To develop a dataset with minimum 2000 sentences we had adopted following approach.

1. Clean the base file and label the tweets as sarcastic and normal.
2. Updated dataset will also have non-tweet sentences
3. These 2000 sentences will be marked as sarcasm or normal by a team of minimum 3 people
4. Finally, whatever is the maximum vote will be the label of the sentence

7.2 Dataset Structure

4. Dataset will have 3 columns “Sentence”, “Context”, “Label”
5. Sentence: Sentence is text of the tweet or any normal sentence.
6. Context: This will be written in the hashtag format (one word). Those tweets which has hashtag it can be extracted from the text and for non-hash tagged sentences and non-tweet sentences context, it will be created manually. Many times, sentence will not have any context. For example “हां मुझे गाली सुनना बहुत पसन्द है” “Yes, I love to hear abuses” This is a sarcastic sentence and there is no context required.
7. Label: This column will have 0 for normal sentence and 1 for sarcastic sentence.

³⁰ <https://github.com/rkp768/hindi-pos-tagger/tree/master/News%20and%20tweets> (Accessed on 26-Jun-20)

7.3 Handling Small Dataset size

We will build our dataset which has 1000 sarcastic statements and 1000 non-sarcastic statements. Because dataset is not large enough therefore, we will use cross validation of 5 folds. For developing neural network-based model we will use 10 folds oversampling.

7.4 Building Dataset

Identify some twitter accounts, hashtags which posts sarcastic text. Write some code in python using tweepy to extract the text from these hashtags and accounts. Extract text from some blogs which write sarcastic articles. Extract each sentence of the blog as a record. Save all these tweets and sentences from the blog into a csv file.

7.5 Cleaning Text

We know that tweet text is unclean because it has text from different languages, in different scripts, extra space, emoticons, non-text sign like "~" ":", "<" etc, flag sign, line break, over used words like ".....", "???????", "beau.....tiful", "!!!!!!". Blog text may also have this kind of text but chances of that is extremely less. We will write a python script to clean all records. Now onwards we will not refer this as tweet or blog text but as sentences. Save all the clean sentence text in a new csv file.

7.6 Labelling

Identify 3 or 5 good Hindi reader who can read the text and identify which sentence is sarcasm and which not. Every manual labeller will label the sentence independently. After getting input from all the people majority of vote will decide whether a sentence is sarcastic or not.

7.7 Transliteration

We know Roman typing is much easy compare to typing in Devanagari therefore many time people use Roman letters in between the sentence. This is true especially if it is name of politician, film actor, place name, (#AmitShah, #Modi, #Salman, #Khan, #India #Bollywood, #Delhi, #Karnataka #Yogi) etc. Because same word will be written in Devanagari and other times in Roman and this is not good for text analysis. So, we will transliterate all the Roman words into Devanagari.

7.8 Context Creation

Whether a sentence is sarcastic or normal sentence, it also depends upon context. For example, "Thank you so much for your help" is normal sentence. But if context is "BJP said to Rahul Gandhi after winning election" then earlier sentence is sarcastic. We will use hashtag of the tweets to extract the context. If tweet has more than one hashtags then we will combine them using "_". If there is no hashtag, which will be true if text is taken from blog, in that we will manually write context. Context will not be sentence but one or two words connected with "_". We want to understand if context is given as hashtag and not as full sentence then how does it impact sarcasm detection.

7.9 Emoticon Handling

We will create another feature called "Emotions" using emoticons found in tweet. We will use corresponding English language word for creating this feature. Text taken from blog will not have any emoticons.

7.10 Embedding

(Sharma et al., 2014) in their work "A Sentiment Analyzer for Hindi Using Hindi Senti Lexicon" suggests using bootstrap approach to extract senti words from Hindi Wordnet. It has given encouraging results of 87% accuracy in sentiment analysis. We are going to test usefulness of this approach in sarcasm detection.

In their paper, Adaptive GloVe and FastText Model for Hindi Word Embeddings, (Gaikwad and Haribhakta, 2020) states that AGM gives better results than GloVe and FastTextWeb. They also mentioned that FastText embeddings which are trained on FastTextHin (Hindi Monolingual corpus) produce better results than FastTextWeb. We are planning to use FastTextHin corpus to check the performance.

Google research has introduced a multilingual BERT which is capable of working with more than 100 languages (Romano, n.d.). We will use this for our project and check how it can be used and how it performs for the task of sarcasm detection in Hinglish.

7.11. Feature Engineering

We will explore different methods for creating feature. For example:

- e) Based on number of Adjective or Adverbs
- f) Hashtag,
- g) Emoticon
- h) Bag of bowls using one word, two words, three words

In their work, (Joshi et al., 2018) have used 3 types of features POS, Named Entities, Unigram to predict the disagreement. However, the results are not encouraging but we would like to explore these features for sarcasm detection.

7.12. Algorithm

Depending upon time available and performance we can include more algorithm, but we will use following 4 algorithms to develop models.

- a. SVM, b. Logistic Regression, c. RNN/GRU/LSTM, d. Naïve Bayesian

7.13 Prediction

Models developed with different embedding and algorithm will be used to predict the result on test dataset. We will use train-test split of 50:50 and 80:20 and check which split helps training the model better.

7.14 Result Comparison

The result of prediction will be compared using Recall, Precision, Accuracy & F1-Score. Results of best models will be ensembled and best possible result with ensembled model will be discussed. This step will help us know that which embedding and algorithm works the best for Hinglish Language sarcasm detection.

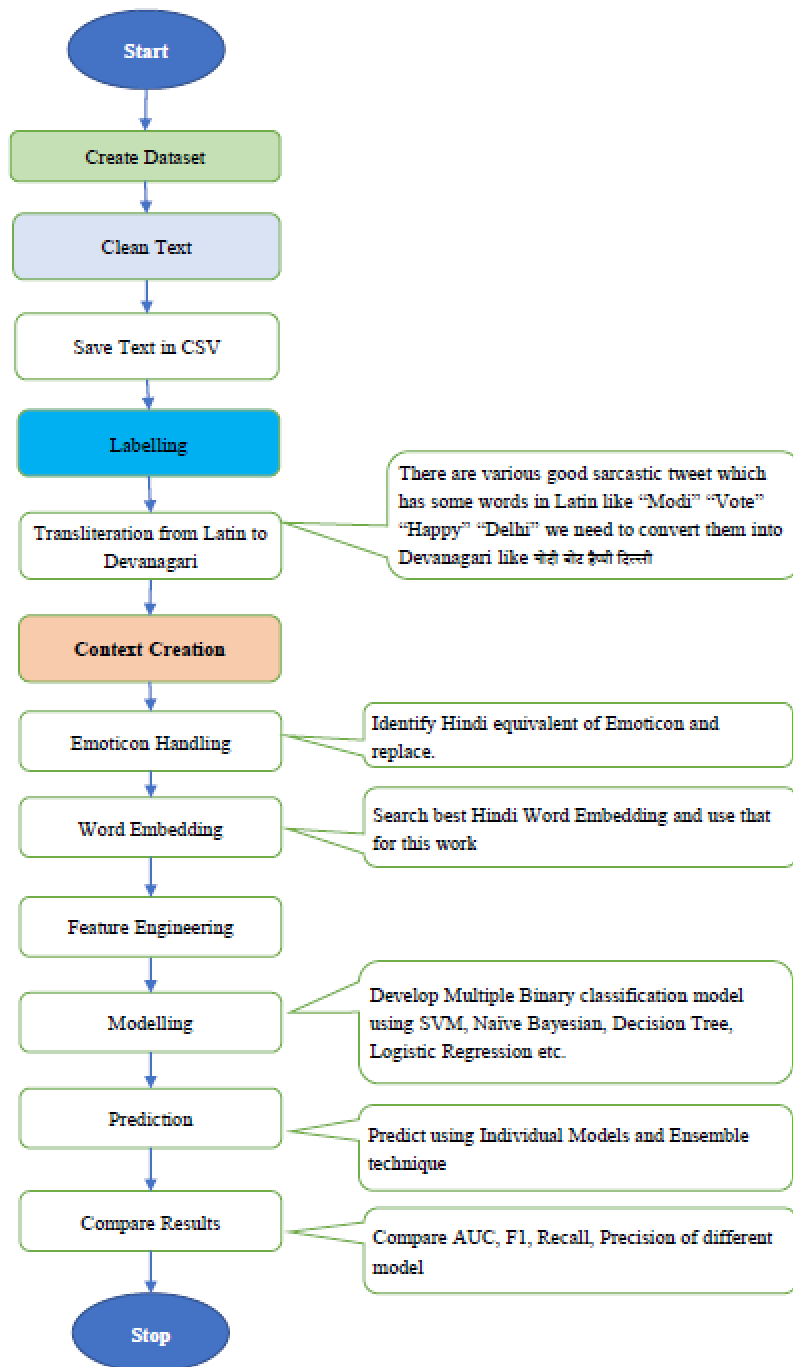
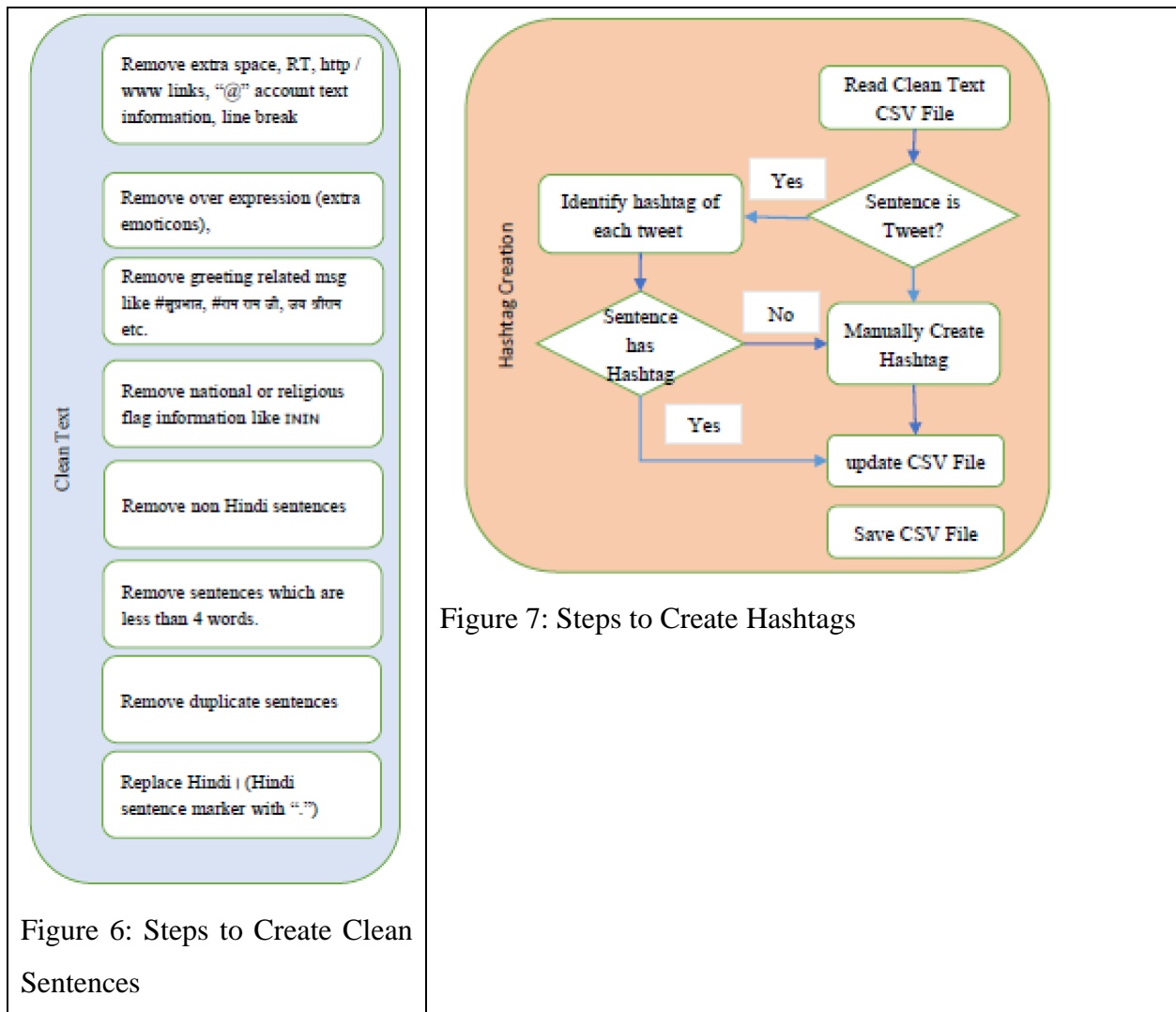
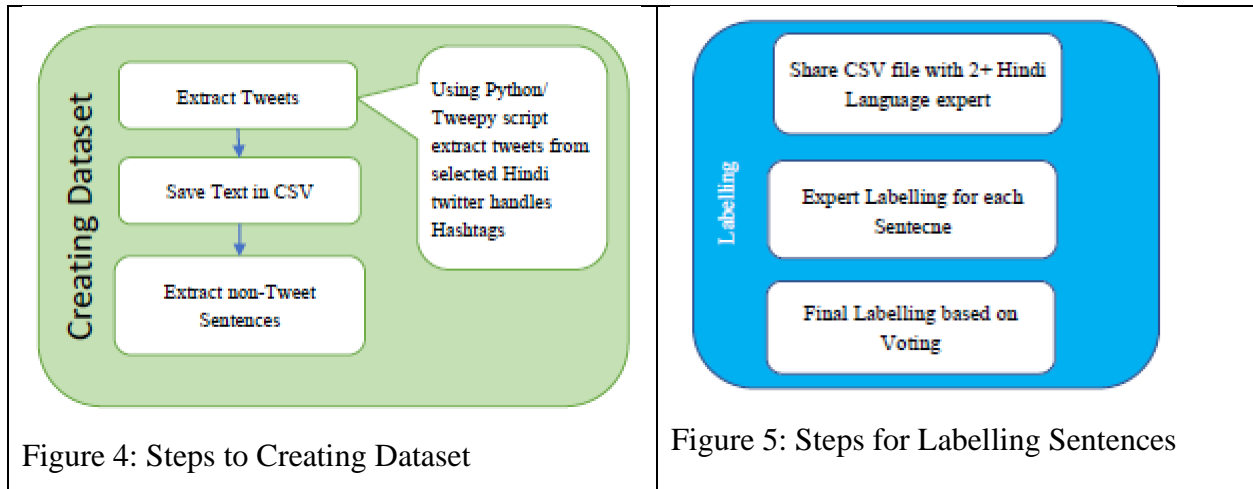


Figure 3: Overall Approach



Evaluation Metrics

ROC graphs are useful tool for visualizing and evaluating classifiers. ROC are able to provide a richer measure of performance than accuracy or error rate (Fawcett, 2004). However, for sake of illustration we will also use Accuracy, F1, Recall & Precision, because they have their relevance depending upon the domain where we use this for sarcasm detection. For example “Hospital administrators thinks I come to hospital because I have lot of money they have beautiful nurses to chat with” (writing sarcasm in English to make sure more readers understand the impact of choice of evaluation metrics). Healthcare domain, hospital administrators may be taking any sarcasm seriously and they do not want any sarcasm to be misclassified and they are ready for more False-True. To illustrate the choice of metrics, lets assume there are 1000 sentences in the dataset, 150 are sarcasm and 850 are normal sentences. Let’s say Model1 predicts 110 are sarcasm and 890 normal and Model2 predicts 140 sarcasm and 860 normal sentences. Accuracy of both the models is 90%. If we select Recall and F1 score then Model1 is better. If we select precision then Model2 is better. If we need to detect sarcasm in comment box of YouTube channel of some political party then we can go for Model1 which is giving recall of 73%. If we are dealing with some more serious product or service like healthcare, airlines service then we can go for Model2 which is giving Precision score of 63%.

| Model1 | | | | | Model2 | | | | |
|-----------------|-------|-----|----|-----|-----------------|-------|-----|----|-----|
| Observation | | | | | Observation | | | | |
| FALSE TRUE | | | | | FALSE TRUE | | | | |
| Actual | FALSE | 820 | 30 | 850 | Actual | FALSE | 805 | 45 | 850 |
| | TRUE | 70 | 80 | 150 | | TRUE | 55 | 95 | 150 |
| 890 110 1000 | | | | | 860 140 1000 | | | | |
| Accuracy 0.90 | | | | | Accuracy 0.90 | | | | |
| Recall 0.73 | | | | | Recall 0.68 | | | | |
| Precesion 0.53 | | | | | Precesion 0.63 | | | | |
| F1 Score 0.81 | | | | | F1 Score 0.77 | | | | |
| Error Rate 0.10 | | | | | Error Rate 0.10 | | | | |

Figure 8: Model Selection based on Evaluation Metrics

8. Expected Outcomes

- Tagged dataset of 2000 sentences
- A system to detect the sarcasm.

- c) Best practices for feature creation in Hinglish language NLP work

9. Requirements / resources

Hardware

- a) Laptop (already have)

Software/Packages

- a) Multilingual BERT
- b) Google Colab (available)
- c) NLTK (available)
- d) scikit-learn.org (available)
- e) seaborn (available)
- f) matplotlib (available)
- g) Google Sheet (for creating dataset)
- h) Microsoft Word (available)
- i) Mendeley (available)
- j) Hindi SentiWordnet
- k) Indic Translation

10. Research Plan

10.1 Risks or contingency plan

| Risk # | Risk Name & Response Plan |
|--------|--|
| 1 | Risk: Latin to Devanagari Transliteration May be more complex than planned Contingency Plan: If we are not able to find or build a suitable solution for translation then we will proceed without transliteration or perform manual transliteration. |
| 2 | Risk: Due to non-availability of any good corpus of Named Entities in Hindi we may not be able to perform NER tagging of sentences. Contingency Plan: We will drop NER experiment from this project. |
| 3 | Risk: If time is constrained and we may not able to write context of all then sentences Contingency Plan: We will develop two solution a- with only those sentences which context b- without context column. Whatever gives better results we will make conclusion based on that. |
| 4 | Risk (Positive Risk): If we have more time and primary goal is achieved. Contingency Plan: We will increase dataset size and perform experiments on the new dataset. |

10.2 Project Schedule

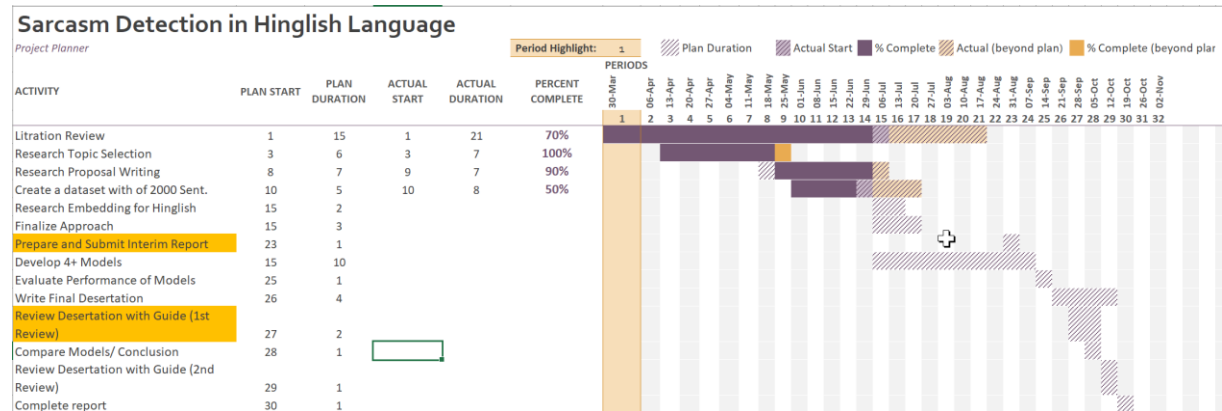


Figure 9: Project Schedule