Durgesh Samariya

# Tweelyzer

## An Approach to Sentiment Analysis of Tweets

**Anchor Academic Publishing**

*disseminate knowledge*

# ABSTRACT

The underlying trend of people using microblogging to express their thoughts on various topic has increased the need for developing computerised techniques for automatic sentiment analysis on texts that do not excessed 200 characters. Twitter is a "micro-blogging" social networking site that has a large and rapidly growing base on users. Twitter's tweets or messages are limited to 140 characters. Because of limitation, it is more difficult to express sentiment and the classification of the tweets difficult as well. The sentiment analysis can be done by two types: emotion and opinion. This research completely focus on sentiment analysis of opinions. These opinions can be divide in three different classes: positive, negative and neutral ( Between positive and negative).

The main goal of this project is to build a model that predict election movement and provide sentiment score from Twitter message (which can not exceed 140 characters). In this project, I apply the novel approach that classify sentiment and emotions of Twitter tweets automatically. After that message is categorised in classes (positive, negative, neutral). For the sentiment first of all, I retrieved tweets from twitter and the convert in to dataset. Therefore applied pre-processing (Data Cleaning) to dataset. After pre-processing applied proposed algorithm namely : TWEELYZER to dataset. At the end I measured performance of TWEELYZER in term of accuracy and recall.

In this project, all tweets of people regarding to movie, brand, actor, actress was collected form twitter and then cleaned and analysed according to proposed algorithm. These tweets were collected using R Studio software. Different process took place in pre-processing tweets, I clean tweets in order to make better for analysis. In process of pre-processing I add some stopping words, removing Hash(#) tag, removing punctuation marks, removing of URLs (ex: www.abcd.com/xyz/abc), etc. After pre-processing, using R Studio developed different insights.

***Keywords:*** Sentiment Analysis, Twitter, Data Analysis, Classification, Big Data, Social Media, R Studio

# ACKNOWLEDGEMENT

# TABLE OF CONTENTS

# LIST OF TABLES

# LIST OF FIGURES