# Classification of Data based on Sentiment of Text using Machine Learning & Deep Learning

## (High Integrity Systems, Master's Thesis)

**Presented by:** Rozana Alam

**Advisor:** Prof. Dr. rer.nat. Matthias F. Wagner

**Co-advisor:** Prof. Dr. Eicke Godehardt

Department of Computer Science and Engineering,

Frankfurt University of Applied Sciences, Frankfurt(Germany)

# Outline

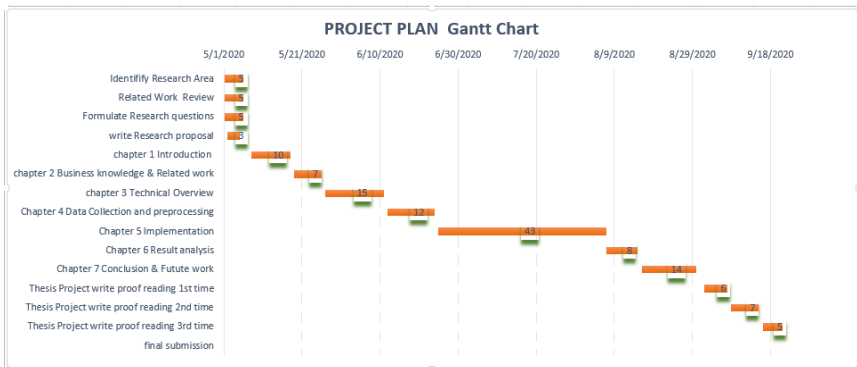# Introduction

- ▶ Overview

- ▶ Motivation

- ▶ Related Work

- ▶ Project plan

# Related work

| Result analysis of Related works | | |
|---|---|---|
| **Model Type** | **Model name** | **Accuracy** |
| Machine Learning | Naïve Bayes | 88% |
| | Support Vector machine | 81% |
| | Random Forest | 63% |
| Deep Learning | Deep Pyramid CNN (Amazon 2) | 94.68% |
| | BLSTM-(SST2) | 89.5% |
| Bert model | Bert base(IMBD) | 95.63% |
| Classification | eight category | 59% |

# Project Plan



PROJECT PLAN Gantt Chart

| | 5/1/2020 | 5/21/2020 | 6/10/2020 | 6/30/2020 | 7/20/2020 | 8/9/2020 | 8/29/2020 | 9/18/2020 |
|---|---|---|---|---|---|---|---|---|
| Identifify Research Area | 5 | | | | | | | |
| Related Work Review | 5 | | | | | | | |
| Formulate Research questions | 5 | | | | | | | |
| write Research proposal | 3 | | | | | | | |
| chapter 1 Introduction | 10 | | | | | | | |
| chapter 2 Business knowledge & Related work | | 7 | | | | | | |
| chapter 3 Technical Overview | | 15 | | | | | | |
| Chapter 4 Data Collection and preprocessing | | | 12 | | | | | |
| Chapter 5 Implementation | | | | 43 | | | | |
| Chapter 6 Result analysis | | | | | 8 | | | |
| Chapter 7 Conclusion & Futute work | | | | | | 14 | | |
| Thesis Project write proof reading 1st time | | | | | | | 6 | |
| Thesis Project write proof reading 2nd time | | | | | | | 7 | |
| Thesis Project write proof reading 3rd time | | | | | | | | 5 |
| final submission | | | | | | | | |

# Business Knowledge and proposed solution

- Natural Language processing(NLP)
- Multi-class Classification
- Sentiment Analysis
- Applications of Sentiment Analysis
- Problem Statement
- Proposed Solution
- Objectives
- Process model of sentiment analysis

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

# Process model of sentiment analysis

# AI, Machine Learning and Deep Learning



**Figure:** AI, ML, and Deep Learning relationship

# Selected Methods for the experiment on our Data

## Text classification with Machine Learning models

# Selected Methods for the experiment on our Data

**Text classification with Machine Learning models**

- ▶ Naïve Bayes classifier

- ▶ Logistic Regression

- ▶ Decision Tree

**Text classification with Deep Learning models**

- ▶ General Neural Network

- ▶ Convolutional Neural Network

- ▶ Long Sort-term Memory(LSTM)

- ▶ Bert_base_uncased_model

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

# Naïve Bayes Theorem ML model

**Naïve Bayes:**

$$P(A/B) = (P(B/A)P(A))/(P(B)) \tag{1}$$

Where P(A/B) and P(B/A) are conditional probabilities.

**Multinomial NB model:**

$$P(c/d) \; \alpha \prod_{1 \leq_k \leq_n d} P(c)P(tk/c) \tag{2}$$

# Logistic Regression



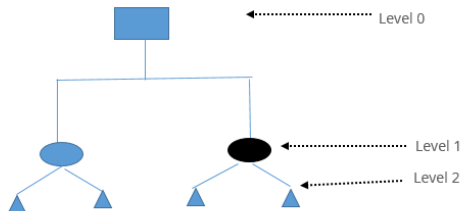**Figure:** Logistic Curve, The values of y cannot be less than 0 or greater than 1

# Decision Tree



Root node: is the beginning of the tree

Internal Node: Splits into further

Leaf node: is a node that no longer splits

Branches: is the link between nodes

Level 0

Level 1

Level 2

# Simple Neural Network

$$\sum_{i=1}^{3} x_i w_i = w_1 * x_1 + w_2 * x_2 + w_3 * x_3 \tag{3}$$
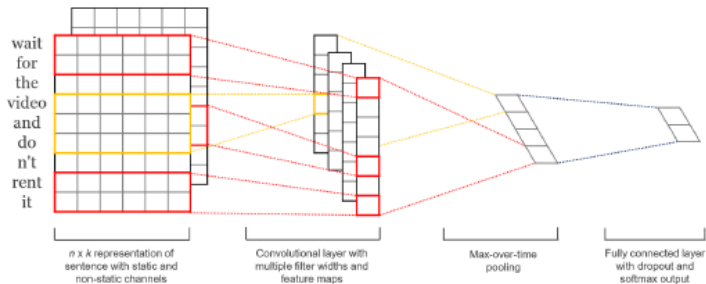


**Figure:** Perceptron

# Convolutional Neural Network



**Figure:** The architecture of a sample CNN model for text classification
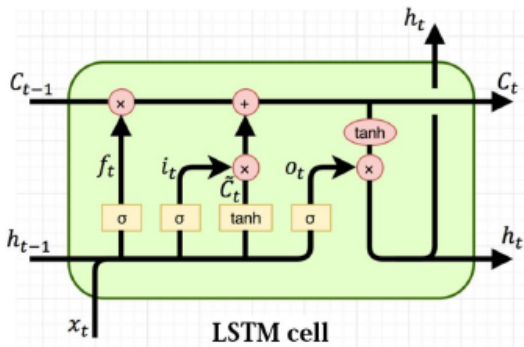
# Long Short-Term Memory(LSTM)



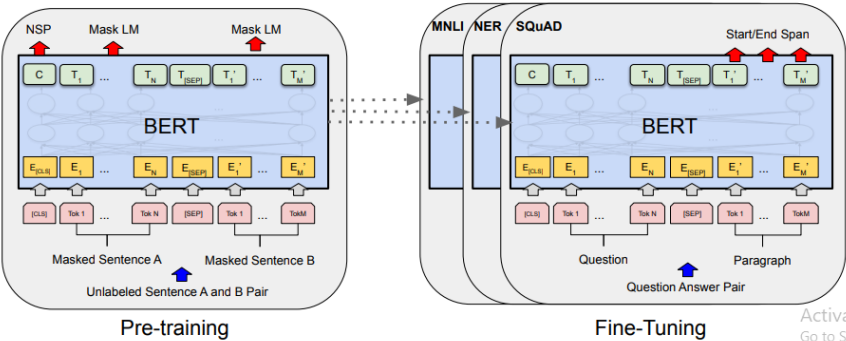**Figure:** Structure of the LSTM cell

# Bert_base_uncased model



**Figure:** BERT model
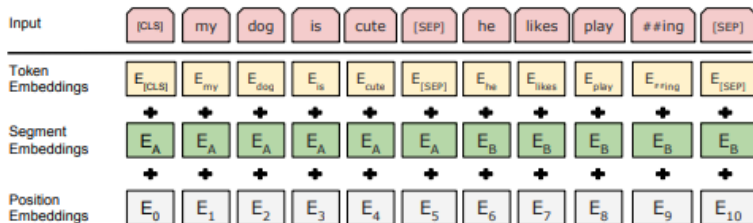
# BERT input representation



**Figure:** BERT input representation

# Data Collection And Data Investigation



```
S_df.head(10)  # first 5 series
```

| | textID | text | selected_text | sentiment |
|---|---|---|---|---|
| 0 | cb774db0d1 | I`d have responded, if I were going | I`d have responded, if I were going | neutral |
| 1 | 549e992a42 | Sooo SAD I will miss you here in San Diego!!! | Sooo SAD | negative |
| 2 | 088c60f138 | my boss is bullying me... | bullying me | negative |
| 3 | 9642c003ef | what interview! leave me alone | leave me alone | negative |
| 4 | 358bd9e861 | Sons of ****, why couldn`t they put them on t... | Sons of ****, | negative |
| 5 | 28b57f3990 | http://www.dothebouncy.com/smf - some shameles... | http://www.dothebouncy.com/smf - some shameles... | neutral |
| 6 | 6e0c6d75b1 | 2am feedings for the baby are fun when he is a... | fun | positive |
| 7 | 50e14c0bb8 | Soooo high | Soooo high | neutral |
| 8 | e050245fbd | Both of you | Both of you | neutral |
| 9 | fc2cbefa9d | Journey!? Wow... u just became cooler. hehe.... | Wow... u just became cooler. | positive |

**Figure:** head lines of S_df
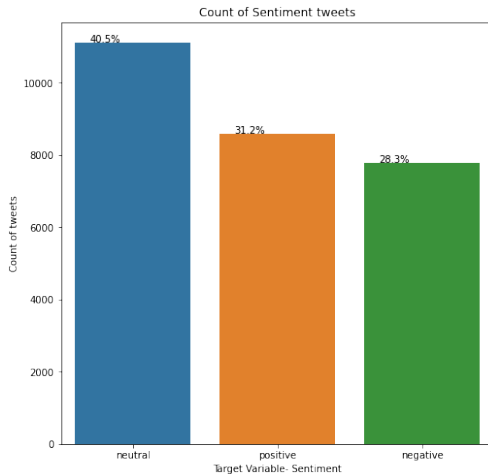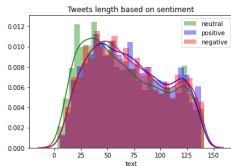
# Data visualization of 'S_df' DataFrame



**Figure:** S_df visualization

# 'Selected_text' column WordCloud Visualalization



**Figure:** 'Selected_text' column WordCloud Visualalization

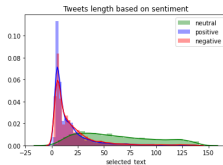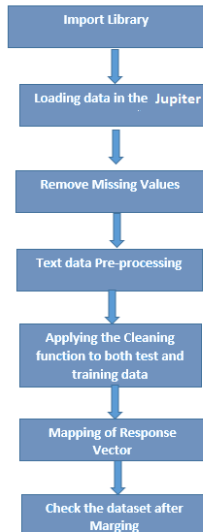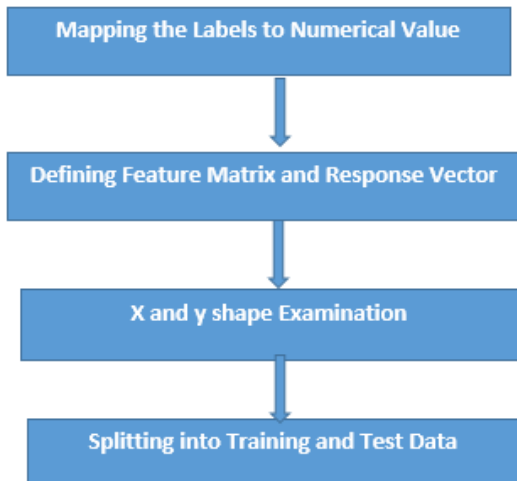# Comparison of 'selected_text' & 'text' column



**Figure:** Value counts for 'selected_text' and 'text' column
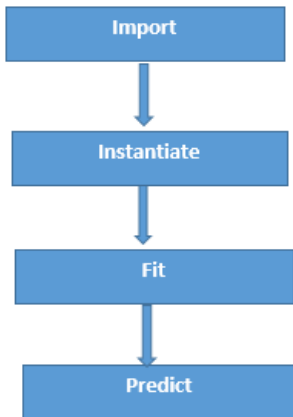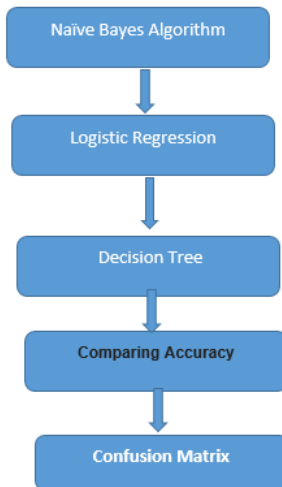
# Data pre-processing for ML models

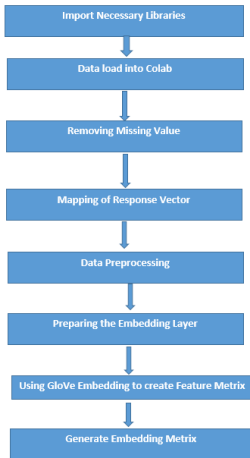# Defining Training and Test Data for ML models
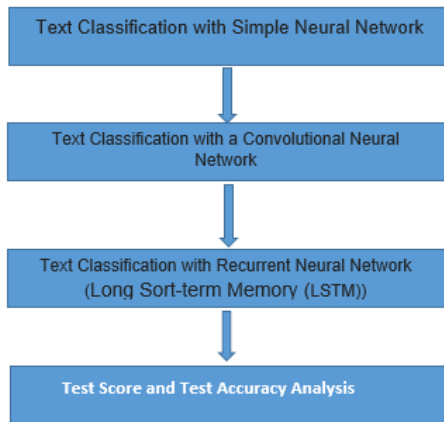
# Flow of Feature Engineering steps
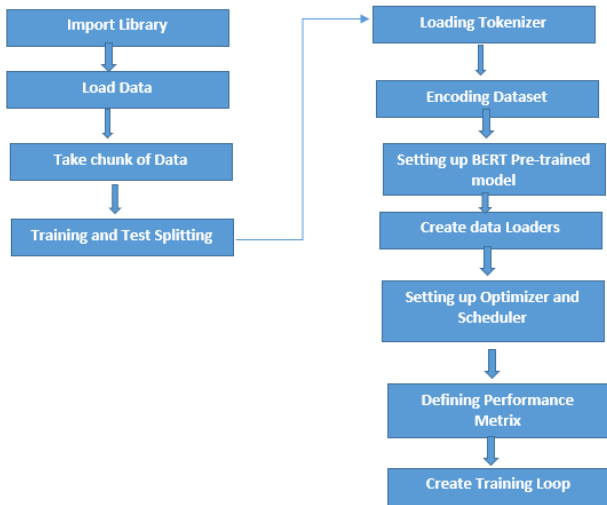
# Steps for best Machine Learning model

# Data pre-processing for Simple NN, CNN and RNN

# Flow of Deep Learning Best model

# Flowchart of BERT model

| Result analysis of all the model | | |
|---|---|---|
| **Model Type** | **Model name** | **Accuracy** |
| ML model | Naïve Bayes | 78.31% |
| | Logistic Regression | 78.85% |
| | Decision Tree | 75.96% |
| DL Model | Simple NN | 48.66% |
| | CNN | 49% |
| | LSTM | 32.50% |
| Bert model | Bert_Base_Fine_Tuning | 88.56% |
| Sentiment for Bert model | Neutral | 86.028% |
| | negative | 88.724% |
| | positive | 90.29% |

FRANKFURT UNIVERSITY OF APPLIED SCIENCES

# Conclusion

---

- Limitations

- Problem Faced

- Future Plan

# Thank You

Questions??