



# RUBIX

PRERIT JAIN

TARUN BHATTAR

# PROBLEM STATEMENT

- SENTIMENT ANALYSIS OF EMPLOYEES REVIEWS FROM “GLASSDOOR.COM” USING MACHINE LEARNING APPROACH.



OR



A hand holding a pen is visible in the background. Above the hand is a laptop with a tree-like structure growing from its screen. The tree is composed of orange lines and numerous circular icons representing various digital concepts like social media, communication, and technology. The word 'CHALLENGES' is written in large, bold, orange capital letters within a thin orange rectangular border.

# CHALLENGES

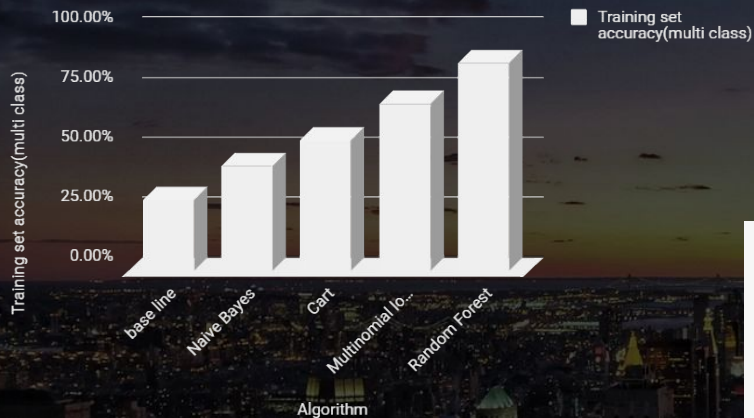
- DATA EXTRACTION
  - Scraping Glassdoor website using a google chrome extension.
- TRAINING DATASET CREATION
  - Manual Approach
  - Dictionary-Based Approach
- DECIDING APPROPRIATE MODELS
  - Trying Different Machine Learning Models.
  - Checking Accuracy.



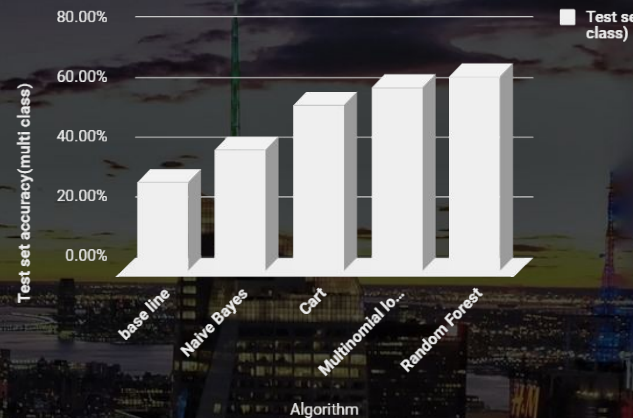
# APPROACH

- SPLITTING THE DATASET.
- CREATION OF CORPUS.
- PREPROCESSING STEPS:
  - Removing Stop Words.
  - Removing Punctuations.
  - Stemming.
- CREATION OF DOCUMENT TERM MATRIX
- REMOVING SPARSITY
- MACHINE LEARNING MODELS APPLICATION.
  - Classification Trees.
  - Random Forests.
  - Naive Bayes Algorithm.
  - Multinomial Logistic Regression.
  - Support Vector Machines
- APPLYING ON TEST SET AND CALCULATING ACCURACY.

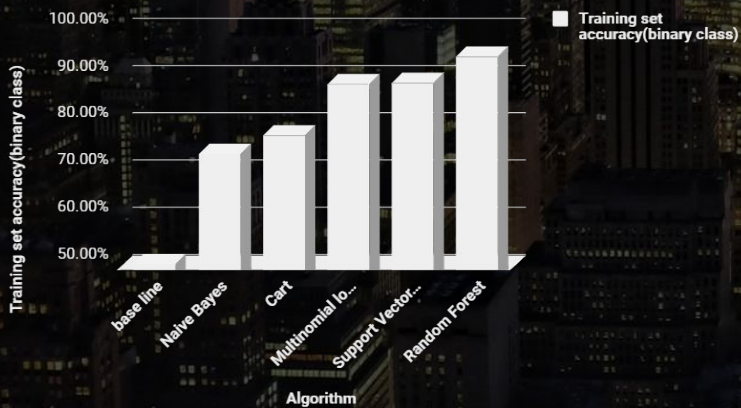
### Training set accuracy(multi class) vs. Algorithm



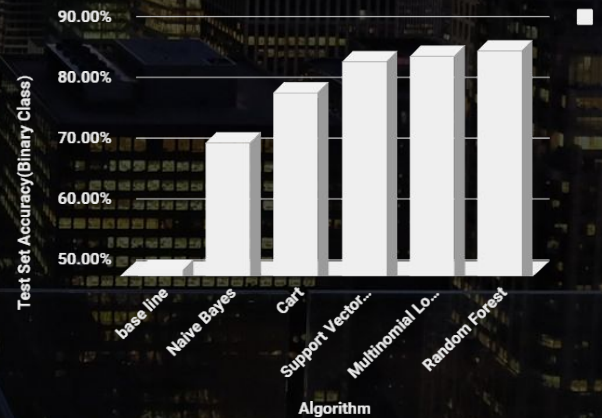
### Test set accuracy(multi class) vs. Algorithm



### Training set accuracy(binary class) vs. Algorithm



### Test Set Accuracy(Binary Class) vs. Algorithm





# OVERALL RATING OF A COMPANY

BENCHMARK SCORE 1 =3.3

BENCHMARK SCORE 2 =3.9

OVERALL RATING =  $5 \cdot \text{POS} / (\text{POS} + \text{NEG})$

POS = No. of Reviews classified as positive

NEG= No. of Reviews classified as negative

- Overall rating > Benchmark score 2

GOOD

- Benchmark score 1 < Overall rating < Benchmark score 2

AVERAGE

- Overall rating < benchmark score 1

BAD



# **LIMITATIONS**

- DATASET SIZE
- SARCASM HANDLING:
- STATEMENTS HAVING MULTIPLE SENTIMENTS.
- LIMITED DICTIONARY
- MULTICLASS CLASSIFICATION

# HOW OUR WORK IS RELATED TO RUBIX ?

## SENTIMENT ANALYSIS

Understanding Customer/Employer Perspective for a company through their sentiment expressed in reviews.

## OVERALL RATING OF A COMPANY

Overall rating by aggregating the individual review ratings.

## AN INDEPENDENT VARIABLE FOR CREDIT RATING

Key Independent variable for analysing the behaviour of a company towards loan/credit.