# RESUME EXPLORER

Created by
Team :

Deep Vira
Vidhi Patel
Pratiksha Javanjal

# OBJECTIVE

The project aims to analyze the resume data to measure and classify the categories of the resumes of the candidates and parse the resume to predict the field domain of the resume which would help the companies select the best fit candidate for a particular position.

# PROBLEM STATEMENT

- To extract text and data from the resume document.
- To classify the categories of the resume data.
- To build a model for classification of categories of the resume.
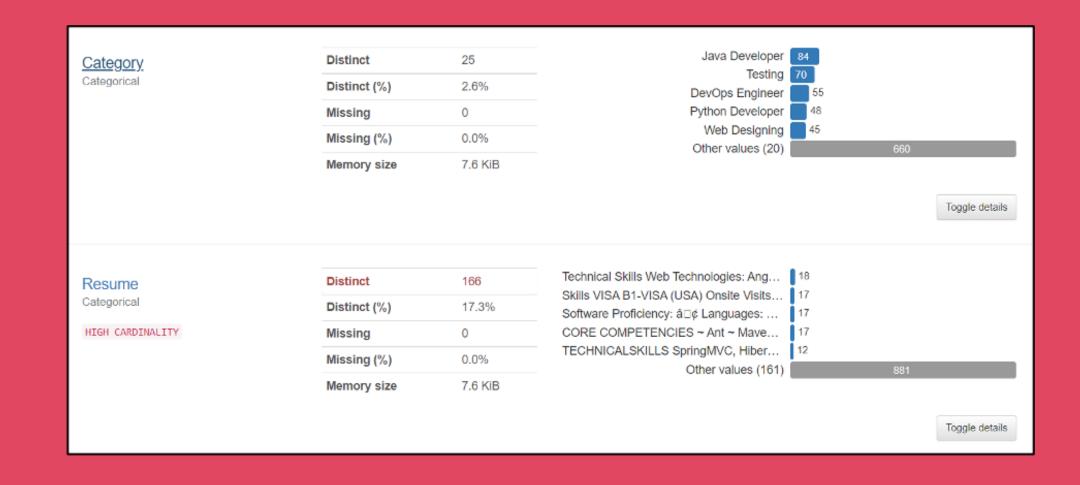- To gauge the best fit model which accurately predicts the data.

# DATA COLLECTION

- The dataset collected is from Kaggle which deals with the resume data consisting the categories and resume text.

- This dataset is used to train the model for classifying the categories of the field for each resume data.

- For the resume parser, a sample resume is collected which is in the PDF format for which the model is built.

- The model can then predict for any resume document that is passed to the model extracted from the PDF file.
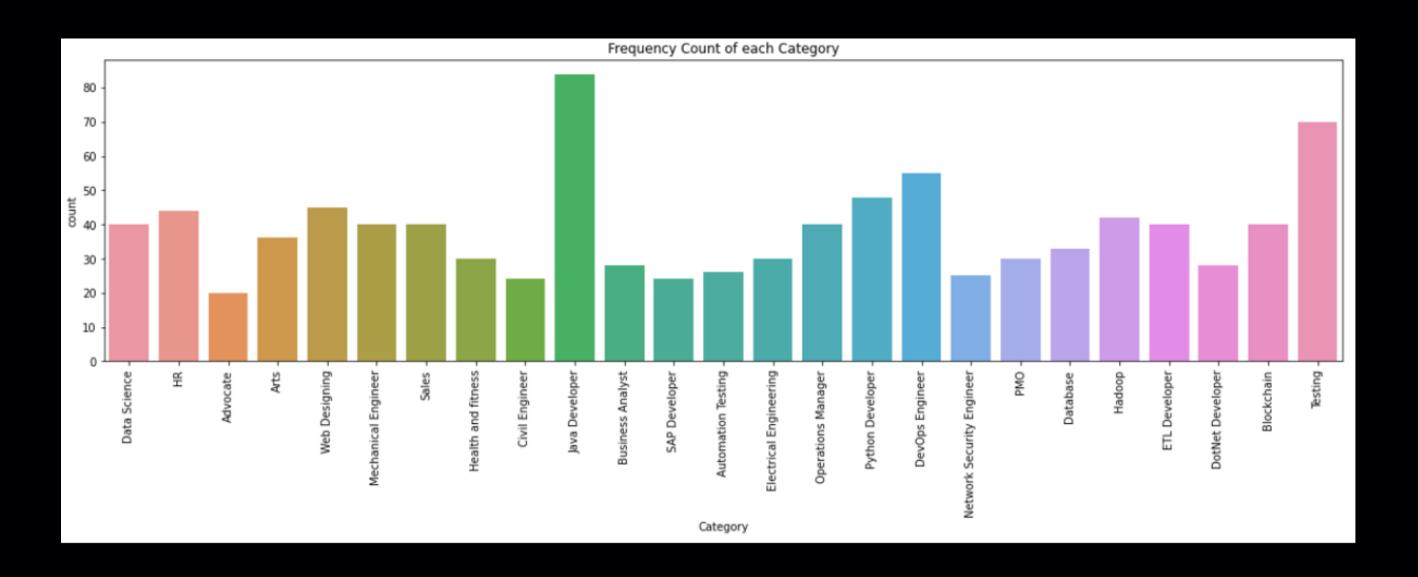
# DATA PROFILING RESULTS

The data profiling report gives an overview of the dataset which includes reports of the statistics of each of the variables in the dataset, the dataset statistics, the missing data plot, and the correlation plot of the parameters.



| Dataset statistics | |
|---|---|
| Number of variables | 2 |
| Number of observations | 962 |
| Missing cells | 0 |
| Missing cells (%) | 0.0% |
| Duplicate rows | 796 |
| Duplicate rows (%) | 82.7% |
| Total size in memory | 5.7 MiB |
| Average record size in memory | 6.0 KiB |

**Category**
Categorical

| | |
|---|---|
| Distinct | 25 |
| Distinct (%) | 2.6% |
| Missing | 0 |
| Missing (%) | 0.0% |
| Memory size | 7.6 KiB |

Java Developer 84
Testing 70
DevOps Engineer 55
Python Developer 48
Web Designing 45
Other values (20) 660

Toggle details

**Resume**
Categorical

HIGH CARDINALITY

| | |
|---|---|
| Distinct | 166 |
| Distinct (%) | 17.3% |
| Missing | 0 |
| Missing (%) | 0.0% |
| Memory size | 7.6 KiB |

Technical Skills Web Technologies: Ang... 18
Skills VISA B1-VISA (USA) Onsite Visits... 17
Software Proficiency: â□¢ Languages: ... 17
CORE COMPETENCIES ~ Ant ~ Mave... 17
TECHNICALSKILLS SpringMVC, Hiber... 12
Other values (161) 881

Toggle details

# DATA VISUALIZATION

- The bar graph represents the frequency count of each category present in the resume dataset.
- This helps in understanding the total count of categories that are present for each resume data that is present in the data.
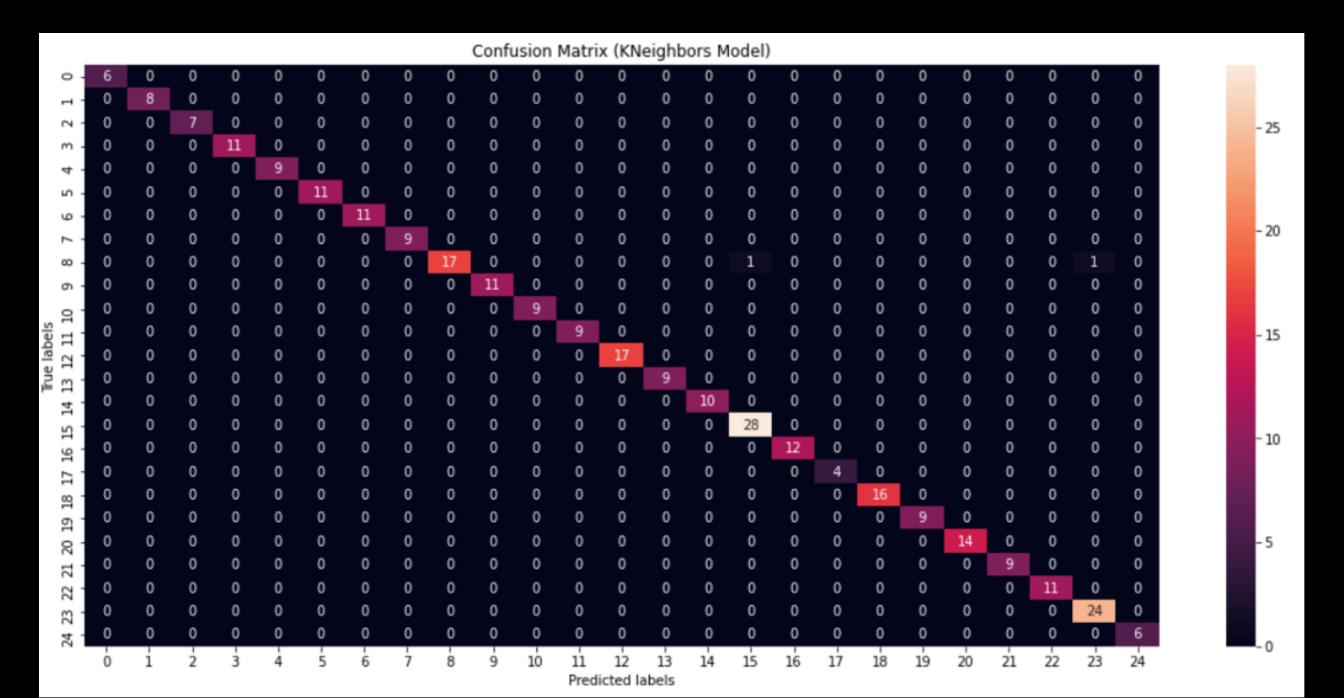


Frequency Count of each Category

# MODEL BUILDING

- Extracting text from a PDF file
- Extracting email address from resume
- Extracting names from resume
- Extracting phone number from resume
- Extracting skills from resume
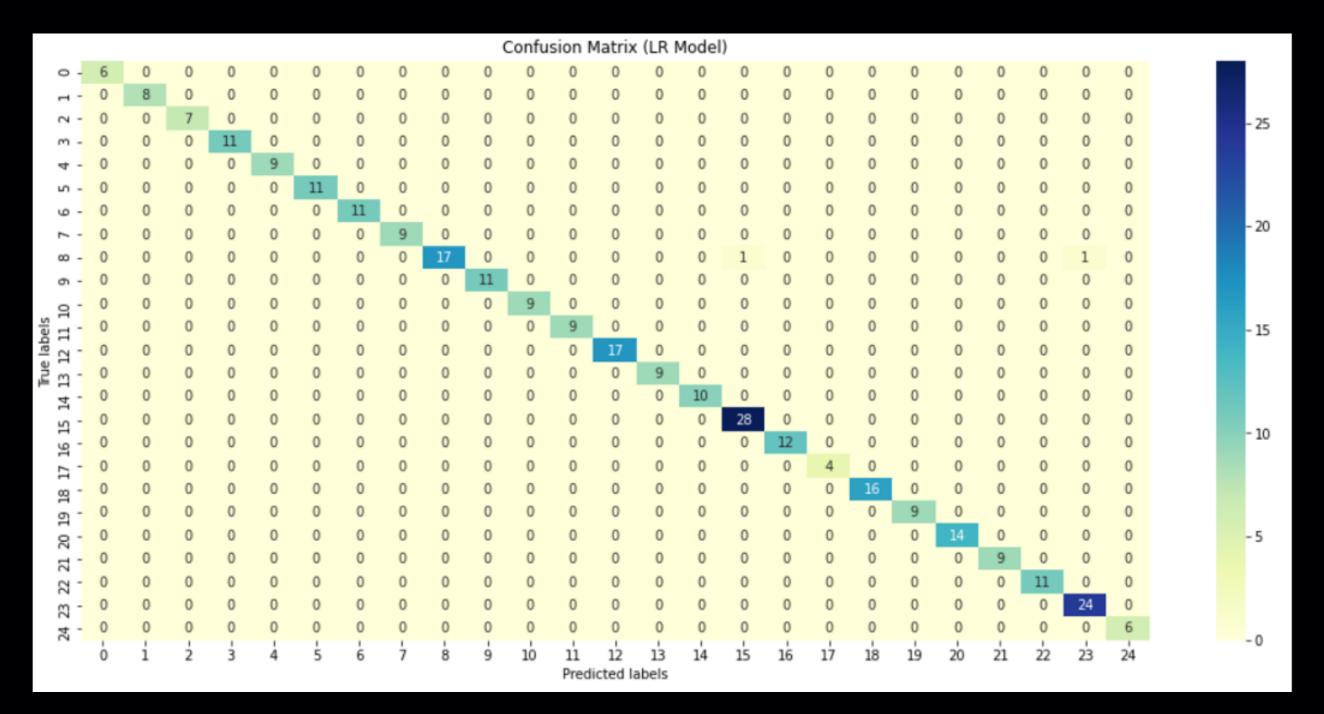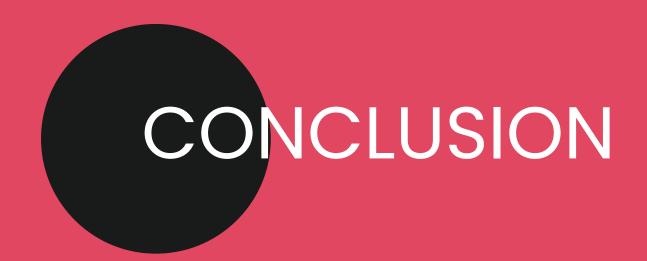- Extracting education and schools from resume

# MODEL BUILDING

Accuracy obtained for KNeighbors Classifier is 97% for the test dataset.



Confusion Matrix (KNeighbors Model)

# MODEL BUILDING

Accuracy obtained for Logistic Regression model is 99% for the test dataset.

# CONCLUSION

- The model built for classification and prediction of the resume uses the dataset containing the categories and the resume information extracted from the document.

- Data is extracted from the sample PDF document which can be further used in the model building phase for classification of the fields and categories.

- The model built takes document as input passed through the model where the various methods of natural language processing are applied and information from the resume document is extracted in order to send the it to the model for training.

- The accuracy obtained for the KNeighbors model was 97% whereas the accuracy for the Logistic Regression model was 99%.