# PIRAMAL FINANCE HACKATHON -
# IIT GUWAHATI

1. Understanding the Data

2. Data Cleaning

3. Feature Selection

4. Feature Transformation

5. Data Splitting

6. Documentation

# UNDERSTANDING THE DATA

1. I initiated a comprehensive review of the dataset to grasp its intricacies.

2. Recognizing the dataset's mixed data types—ranging from strings to ranges—presented a notable challenge.

3. Additionally, the prevalence of null and empty values raised concerns regarding the dataset's reliability.

4. To streamline the analysis process, I utilized a CSV AI tool to compare unique values within each column.

5. Furthermore, as the CV files were initially in DOC format, I converted them to PDF for easier data extraction, facilitating smoother handling and analysis.
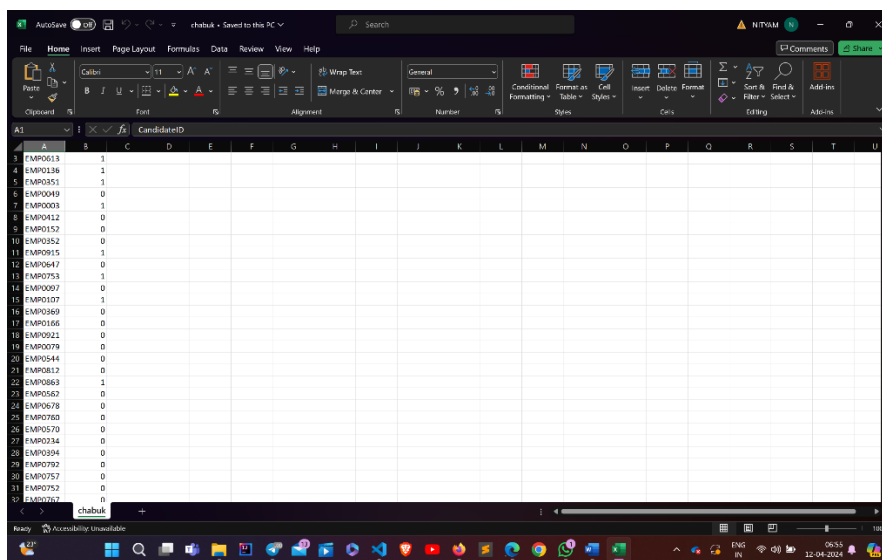
# DATA CLEANING

• Utilized tools like ChatCSV for efficient data cleaning, effectively addressing null values by employing statistical methods such as mean, median, or 0 where appropriate.

• Converted string data types into simpler numerical representations conducive to model processing, enhancing data compatibility for analysis.

• Leveraged the pandas library to develop Python scripts tailored for data cleaning, implementing specific parameters to ensure data integrity and quality.

• Executed Python scripts within Visual Studio Code (VSCode) and Google Colab, optimizing workflow and leveraging their capabilities for streamlined data processing.

• Employed various techniques for visualizing data, including tables, charts, and other graphical representations, facilitating insightful analysis and interpretation.

• Automated the conversion of document files, such as CVs in test and train sets, into PDF format using Python

scripts, enhancing data accessibility and compatibility.

• Implemented data refinement strategies, such as removing common data attributes like designations, to streamline analysis and focus on relevant information.

• Encoding String Data: For string-type categorical variables, we can encode them into numerical values using techniques like one-hot encoding or label encoding.

• One-Hot Encoding: This technique converts each category into a binary vector where each category is represented by a binary digit (0 or 1). This is particularly useful when the categories don't have a natural ordering or hierarchy.

• Label Encoding: Label encoding assigns a unique integer to each category. This is suitable when the categories have a natural order or hierarchy.

• Integer Encoding for Range Data: Range data, which consists of numerical intervals, can be encoded into integer values based on predefined ranges or bins.

• Binning or Discretization: This involves dividing the range of numerical values into discrete intervals or bins and assigning an integer value to each bin. This approach is useful for converting continuous data
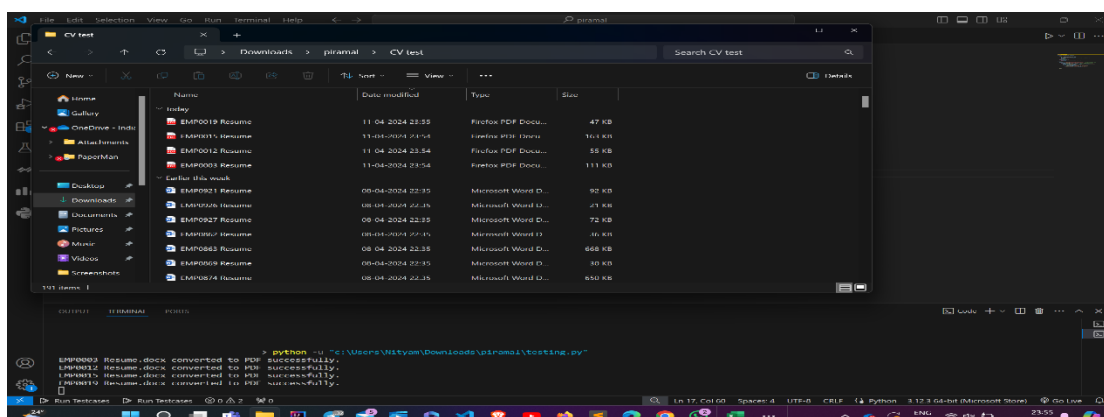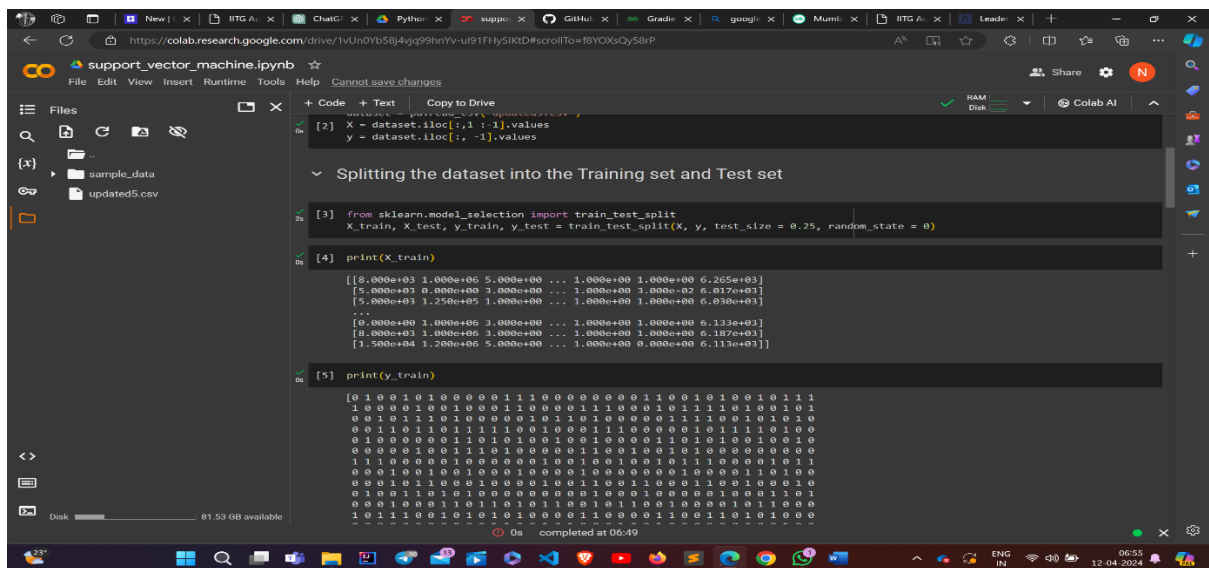
into categorical or ordinal data.

• Custom Encoding: Depending on the specific requirements and characteristics of the range data, custom encoding schemes can be devised. For example, assigning integer values based on predefined thresholds or business rules.

# FEATURE SELECTION

• Statistical Analysis on Google Colab: Initially, I performed a detailed

statistical analysis using Google Colab, which allowed me to sift through

the multitude of parameters and discern their impact on performance.

This step was crucial in understanding which variables significantly

influenced the outcome.

• Linear Regression for Parameter Importance: Following the initial

analysis, I employed linear regression techniques in Python to further

investigate the importance of various parameters. By running regression

models, I could quantify the influence of each parameter and identify

the ones that were most critical for performance.

• Conversion of Parameters into Usable Form: Once the important

parameters were identified, I converted them into a usable format

suitable for further analysis and modeling. This step involved

transforming the data into a structured and manageable form for

subsequent processing.

• Optimization with Singular Value Decomposition (SVD): To enhance

the efficiency of the testing data, I utilized Singular Value

Decomposition (SVD) techniques. By decomposing the data matrix, SVD

helped optimize the test dataset, improving the accuracy and reliability

of subsequent analyses.

• Utilization of scikit-learn for ML Classification Models: I employed

scikit-learn, a powerful machine learning library in Python, to

experiment with various classification models. By employing different

algorithms, I aimed to ascertain which parameters were indispensable

for accurate predictions and classification tasks.

• PDF CV Parsing with PyPDF2: In order to extract relevant details such

as languages spoken and skills possessed by candidates, I utilized

PyPDF2's PdfReader module. This enabled the automated extraction of

information from PDF-based resumes, streamlining the data collection

process for further analysis.

# FEATURE TRANSFORMATION

• Employed encoding techniques to transform string-type categorical variables into numerical values, such as:

• Utilized one-hot encoding for non-ordinal categories, converting each category into a binary vector representation.

• Implemented label encoding for categories with a natural order or hierarchy, assigning unique integer values to each category.

• Applied integer encoding methods to handle range data, including:

• Utilizing binning or discretization to divide numerical intervals into predefined bins or categories.

• Assigned integer values to each bin, facilitating the conversion of continuous variables into categorical or ordinal ones.

• By employing these feature transformation techniques, ensured all data types, including strings and range data, were converted into numerical formats suitable for machine learning algorithms.

• Prepared the dataset for integration into various machine

learning models, enabling robust analysis and accurate predictions.

Integrated normalization for selected parameters to ensure consistent scaling across features:

• Identified specific parameters where normalization would be

beneficial, such as those with varying scales or units.

• Applied normalization techniques to rescale these parameters

to a standard range, typically between 0 and 1

• Experimented with a combination of classification models including SGBoost, CatBoost, Random Forest Classification, and K-Nearest Neighbors (KNN).

• Implemented ensemble techniques to leverage the strengths of each model and mitigate their individual weaknesses.

• Combined predictions from SGBoost, CatBoost, Random Forest, and KNN to generate a consolidated output, aiming to achieve higher accuracy and robustness.

• Fine-tuned the ensemble approach to optimize performance based on validation metrics and domain-specific requirements.

• Evaluated the effectiveness of the ensemble method compared to individual models through rigorous testing and analysis.

# DATA SPLITTING

• Data Preparation:

• Created separate training and testing files from the provided

train.csv and test.csv datasets.

• Preserved approximately 20% of the data from train.csv for

validation purposes, ensuring a sufficient amount for training the

model.

• Designated test.csv exclusively for prediction and final model

evaluation.

• Data Splitting Strategy:

• Utilized a 80-20 data splitting strategy, allocating 80% of train.csv

for model training and the remaining 20% for validation.

• Ensured that the splitting process preserved the integrity of the

dataset and maintained representative samples across the

subsets.

• Random Shuffling:

• Randomly shuffled the data to eliminate any inherent ordering or

biases that may exist.

• This step was crucial for enhancing model accuracy and

preventing overfitting by exposing the model to diverse samples

during training.

• Validation and Evaluation:

• Reserved the validation set for fine-tuning model

hyperparameters and assessing performance during training.

• Kept the test set separate for final evaluation, enabling unbiased
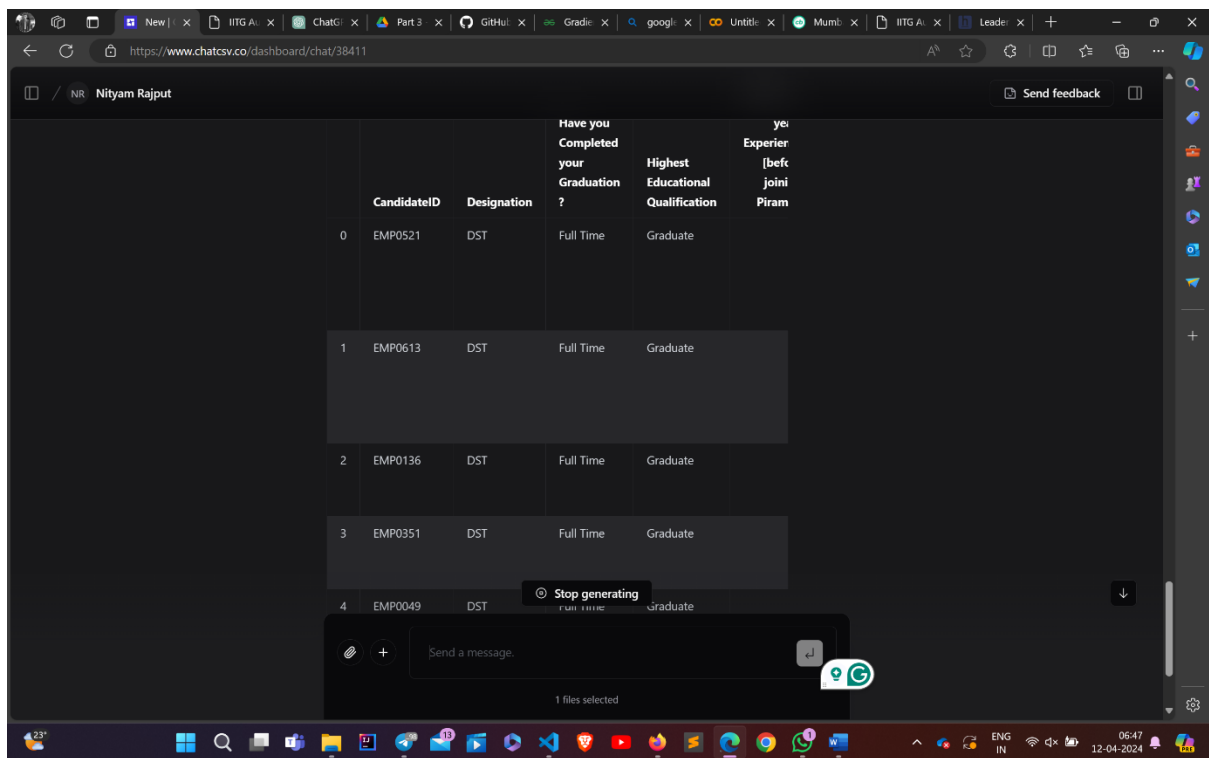
assessment of the model's generalization capabilities on unseen

data.

• Data Integrity Check:

• Ensured that both the output files and the original datasets

(train.csv and test.csv) had the correct number of rows and

columns.

• Maintained consistency in data dimensions to prevent errors

during training and evaluation processes.By meticulously

organizing and preparing the data in this manner, facilitated the

development of robust and generalizable machine learning

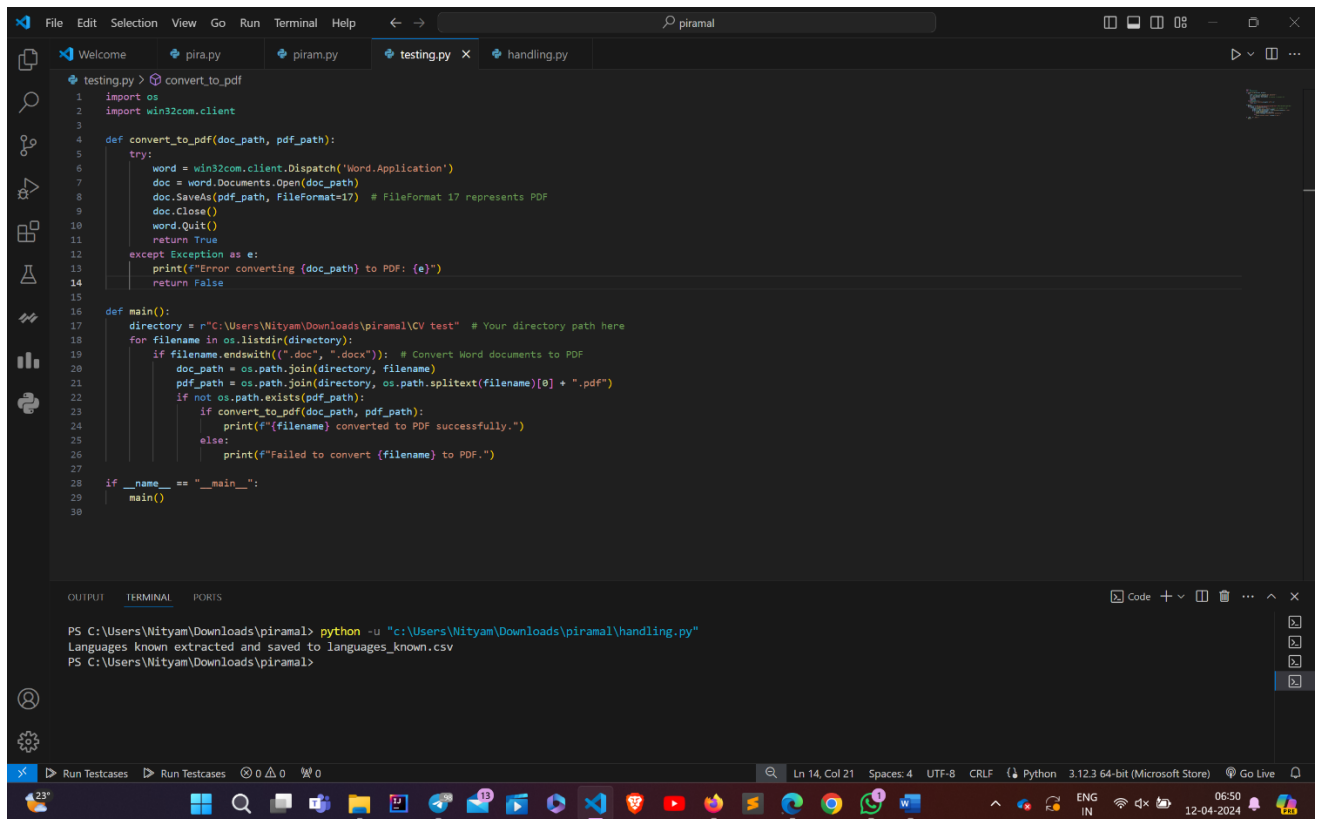models capable of delivering accurate predictions on new, unseen

data.



Chatcsv for getting details of data

# DOCUMENTATION

• Utilized cepstrum IITG'S classification ML models for initial classification tasks.

• Incorporated SGBoost and CatBoost software from GitHub repositories to enhance model performance.

• Referenced documentation from GeeksforGeeks (GFG) on SGBoost for implementation guidance and best practices.

• Leveraged ChatGPT to amalgamate results from multiple models, aiming to improve overall accuracy.

• Employed a combination of techniques from various sources to create a robust and accurate classification system.

CONVERTING THE DOC FILES TO PDF

EXTRACTING LANGUAGE DATA FROM CV IN TRAIN

NITYAM

ROLL NUMBER – 220107059

CONTACT – n.nityam@iitg.ac.in