**Documentation for the Resume Categorization System**

I reviewed a collection of PDF resumes from different types of individuals. In the dataset, I found that there were around 100 resumes per group, on average. The goal is to develop a machine learning or deep learning model that predicts the category of people based on their resume and saves the results of the test dataset in a CSV file.

To accomplish this, I first extracted the text from each PDF and labeled them accordingly. I used regular expressions to remove stopwords and extra spaces. I applied TF-IDF to reduce the impact of common words while emphasizing rarer words that are more frequent in specific groups' resumes.

For classification, I chose to use the Support Vector Machine (SVM) model. Although I could have opted for other models like Random Forest, ANN, RNN, LSTM, or even heavier models like BERT, due to time constraints and the limited processing power of my personal computer, I decided to go with the simplest approach.

After completing the code, the model's evaluation metrics for the validation set were as follows: Accuracy: 0.6660, Precision: 0.6838, Recall: 0.6660, and F1-score: 0.6611. For the testing set, the metrics were Accuracy: 0.6177, Precision: 0.6313, Recall: 0.6177, and F1-score: 0.6046.

- **Accuracy** measures the proportion of correct predictions across the entire dataset.

$$Accuracy \ = \ \frac{TP + TN}{TP + FP + TN + FN}$$

- **Precision** indicates the reliability of positive predictions.
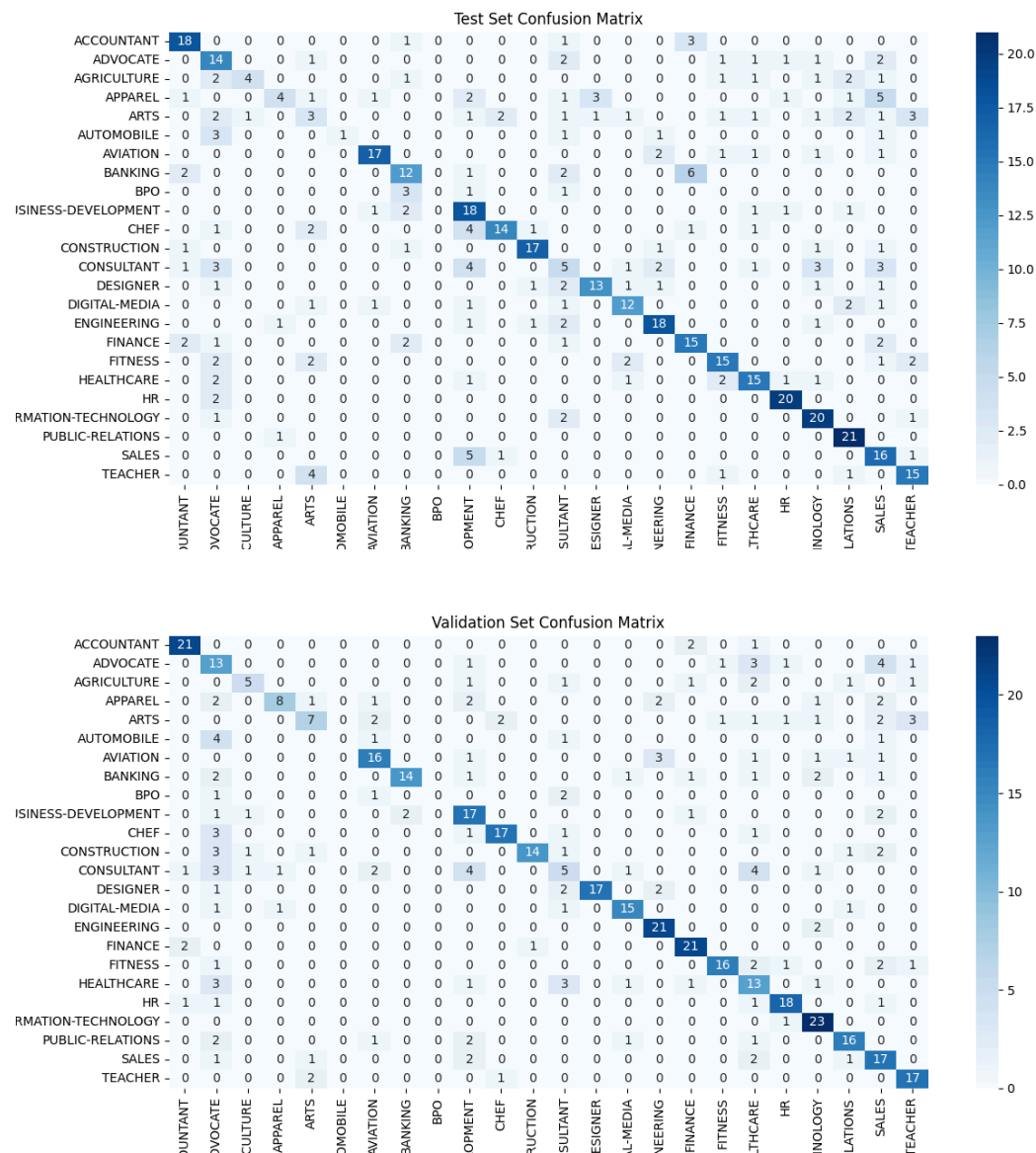
$$Precision \ = \ \frac{TP}{TP + FP}$$

- **Recall** measures the model's ability to correctly identify actual positive instances.

$$Recall \ = \ \frac{TP}{TP + FN}$$

- **F1-score** seeks to balance Precision and Recall.

$$F1 - score \ = \ \frac{Precision \times Recall}{Precision + Recall}$$

Here the photo of confusion matrices:



Test Set Confusion Matrix



Validation Set Confusion Matrix

**Improvements**: There are several ways to improve these results. In the dataset, the highest number of resumes belongs to Information Technology professionals (120), while the lowest number of resumes is from BPO professionals (22). By applying random oversampling to the smaller datasets, the results could potentially be improved.