# Project Report on MBTI Personality Prediction

## Introduction to the Problem

The Myers-Briggs Type Indicator (MBTI) is a widely-used personality assessment tool that categorizes individuals into one of 16 distinct personality types based on four dichotomies. Understanding an individual's MBTI type can be beneficial in various domains such as career counseling, team building, personal development, and interpersonal relationships. Automating the prediction of MBTI types from textual data, such as social media posts or personal writings, can offer scalable insights into user personalities without the need for extensive surveys.

## Our Approach

## Data Preprocessing

To prepare the text data for modeling, we employed several preprocessing steps to clean and normalize the data. The `clear_text` function implemented these steps:

> **Lowercasing**: Standardizing the text to lowercase to ensure uniformity.
> **URL Removal**: Stripping away hyperlinks that do not contribute to personality type inference.
> **Symbol Removal**: Cleaning non-alphanumeric characters to focus on textual content.

# Handling Data Imbalance

Given the potential class imbalance in the dataset, we opted to use the `class_weight='balanced'` parameter in our models. This approach helps to ensure that each class is given appropriate importance during training, despite differences in their frequency. Additionally, during the train-test split, we utilized `stratify=y` to maintain a consistent distribution of MBTI types in both training and testing datasets, enhancing the reliability of our model evaluations.

# Machine Learning Models

We experimented with several machine learning models known for their effectiveness in text classification tasks:

- **Logistic Regression**
- **Support Vector Machine (SVM)**
- **Multinomial Naive Bayes**
- **Random Forest**
- **BERT**

These models were integrated into pipelines with two types of text vectorizers: TF-IDF and Count Vectorizer. This approach allowed us to convert raw text into a format suitable for model training while evaluating the impact of different text representation techniques.

# Model Evaluation and Selection

We conducted 5-fold cross-validation for each model and pipeline combination to robustly assess their performance. The evaluation metrics included average accuracy, F1-score, recall, precision, and a confusion matrix. The best-performing combinations were:

- **TF-IDF Vectorizer with Logistic Regression**:
    - ☐ Average Accuracy: 0.691
    - ☐ Average Precision: 0.697
    - ☐ Average Recall: 0.691
    - ☐ Average F1 Score: 0.691

Confusion Matrix: TF-IDF Vectorizer with Logistic Regression

| True\Predicted | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 22 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 6 | 2 | 2 | 0 | 0 | 2 | 0 | 2 |
| 1 | 3 | 82 | 5 | 5 | 0 | 3 | 0 | 0 | 6 | 13 | 9 | 5 | 1 | 0 | 1 | 2 |
| 2 | 1 | 0 | 31 | 2 | 0 | 0 | 0 | 1 | 2 | 2 | 4 | 3 | 0 | 0 | 0 | 1 |
| 3 | 1 | 9 | 1 | 98 | 0 | 0 | 0 | 0 | 3 | 4 | 6 | 8 | 2 | 2 | 2 | 1 |
| 4 | 0 | 0 | 0 | 1 | 5 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| 5 | 0 | 0 | 1 | 2 | 0 | 3 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 2 | 0 | 0 |
| 6 | 0 | 0 | 1 | 0 | 0 | 0 | 3 | 0 | 0 | 2 | 0 | 0 | 0 | 1 | 0 | 1 |
| 7 | 0 | 0 | 1 | 2 | 1 | 0 | 0 | 8 | 0 | 0 | 0 | 2 | 0 | 1 | 0 | 3 |
| 8 | 9 | 6 | 4 | 8 | 2 | 2 | 0 | 1 | 198 | 28 | 13 | 11 | 1 | 6 | 1 | 4 |
| 9 | 5 | 10 | 2 | 4 | 2 | 2 | 0 | 1 | 19 | 270 | 11 | 20 | 3 | 11 | 2 | 4 |
| 10 | 3 | 6 | 7 | 2 | 1 | 1 | 0 | 1 | 9 | 3 | 151 | 21 | 0 | 4 | 2 | 7 |
| 11 | 0 | 2 | 6 | 12 | 3 | 0 | 0 | 1 | 8 | 12 | 14 | 189 | 0 | 6 | 0 | 8 |
| 12 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 1 | 25 | 1 | 0 | 1 |
| 13 | 0 | 1 | 0 | 2 | 0 | 1 | 0 | 0 | 2 | 10 | 1 | 1 | 1 | 35 | 0 | 0 |
| 14 | 0 | 2 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 4 | 0 | 3 | 1 | 1 | 27 | 0 |
| 15 | 0 | 1 | 0 | 2 | 0 | 0 | 0 | 1 | 0 | 2 | 1 | 2 | 4 | 1 | 2 | 52 |

Predicted labels / True labels

- **TF-IDF Vectorizer with SVM**:
  - ☐ Average Accuracy: 0.6280115273775215
  - ☐ Average Precision: 0.6463181833007485
  - ☐ Average Recall: 0.6280115273775215

Average F1 Score: 0.6138485255336215

Confusion Matrix: TF-IDF Vectorizer with Support Vector Machine



● **Count Vectorizer with SVM**:
   ☐ Average Accuracy: 0.657521613832853
   ☐ Average Precision: 0.6608670500584127
   ☐ Average Recall: 0.657521613832853
   ☐ Average F1 Score: 0.652151263402881

Confusion Matrix: Count Vectorizer with Support Vector Machine

# Encoding of MBTI Personality Types

To facilitate the machine learning processes, we have converted the Myers-Briggs Type Indicator (MBTI) personality types from textual representations to numerical encodings using LabelEncoder. This transformation is crucial for the effective application of algorithms that require numerical input. Below is the mapping used in our analysis:

- ★ **0**: ENFJ
- ★ **1**: ENFP
- ★ **2**: ENTJ
- ★ **3**: ENTP
- ★ **4**: ESFJ
- ★ **5**: ESFP
- ★ **6**: ESTJ
- ★ **7**: ESTP
- ★ **8**: INFJ
- ★ **9**: INFP
- ★ **10**: INTJ
- ★ **11**: INTP
- ★ **12**: ISFJ
- ★ **13**: ISFP
- ★ **14**: ISTJ
- ★ **15**: ISTP

# Hyperparameter Tuning

Based on the initial results, we selected the top three model configurations for further optimization through hyperparameter tuning. The tuning process also considered different parameters for the vectorizers to maximize performance. This step is crucial for refining our models to achieve the best possible predictions and is currently ongoing, with results pending at the time of this report's deadline.
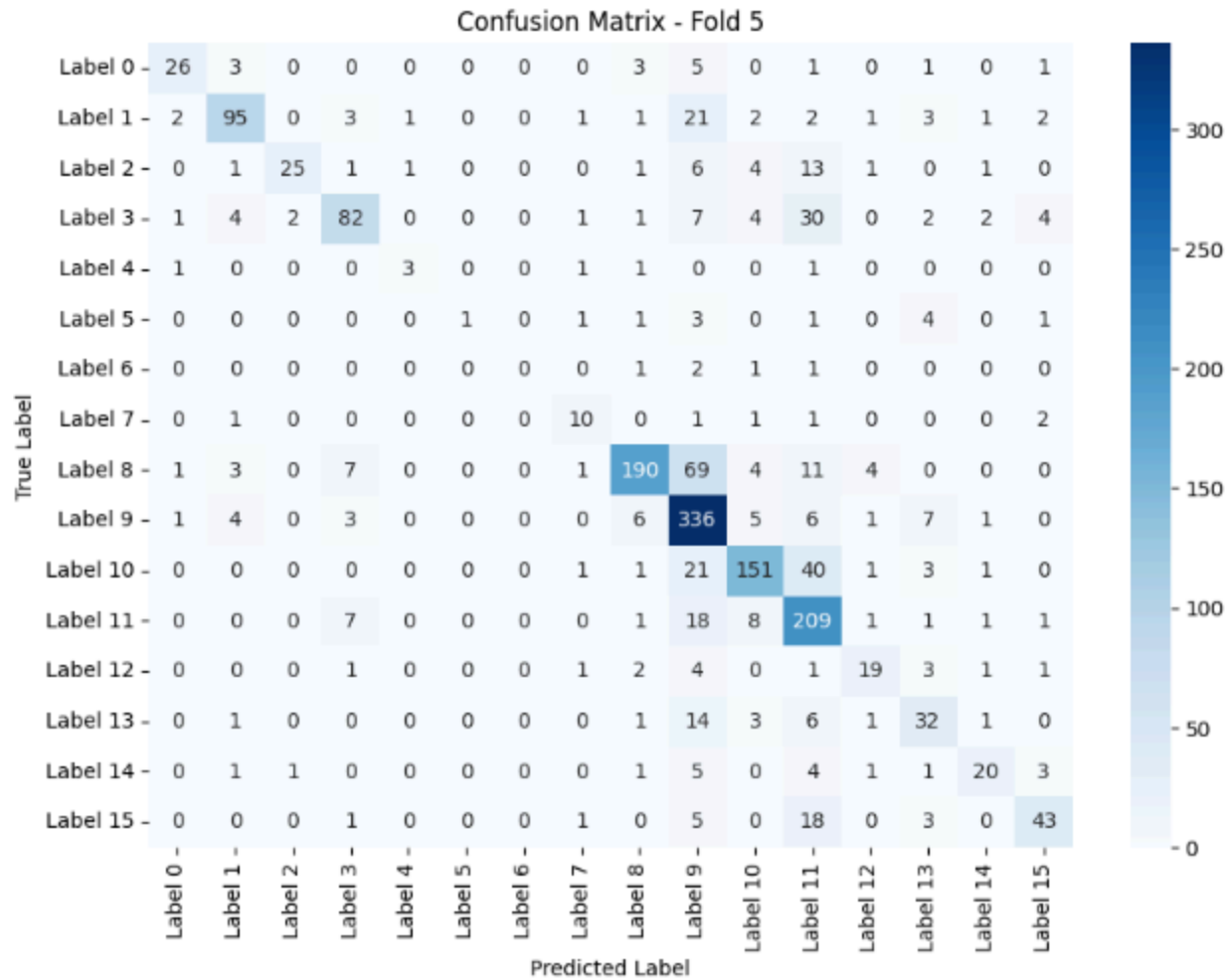
# Bert

With BERT, we can arguably see the best accuracy of up to 71% with only 5 epochs thus this offers significantly better performance. However, we observe that labels with very less data are hard to predict.

Final accuracy average:
{'accuracy': [0.715850144092219],
 'precision': [0.7441454649692727],
 'recall': [0.715850144092219],
 'f1': [0.7114172551038348]}

```
Running fold 1
Epoch 0: Train Loss: 2.1977952563542926, Val Loss: 1.8380326507827653, Val Accuracy: 0.44726224783861673
Running fold 2
Epoch 0: Train Loss: 1.6370930902419552, Val Loss: 1.4173097518457245, Val Accuracy: 0.579250720461095
Running fold 3
Epoch 0: Train Loss: 1.352871043307166, Val Loss: 1.222883027132755, Val Accuracy: 0.6461095100864553
Running fold 4
Epoch 0: Train Loss: 1.1919911193408175, Val Loss: 1.0116676601328058, Val Accuracy: 0.7060518731988472
```

Confusion Matrix - Fold 5

**Label Mappings:**

```
Label mappings:
0: ENFJ
1: ENFP
2: ENTJ
3: ENTP
4: ESFJ
5: ESFP
6: ESTJ
7: ESTP
```

```
8:  INFJ
9:  INFP
10:  INTJ
11:  INTP
12:  ISFJ
13:  ISFP
14:  ISTJ
15:  ISTP
```

## Conclusion

This project aims to demonstrate the feasibility of predicting MBTI personality types using machine learning techniques applied to text data. Through systematic preprocessing, model evaluation, and hyperparameter tuning, we strive to develop a robust predictive system. Future work will focus on completing the tuning process and potentially deploying the model in a real-world application to provide automated personality assessments.