# Disease Prediction

**FINAL PROJECT**

**PRESENTED BY**

**MEHZABEEN GHEESAH**

# Project Objective

PREDICT A DISEASE BASED ON THE DESCRIBED SYMPTOMS

# Dataset Overview

Name: Symptom2Disease

Source: Kaggle

Rows: 1200

Columns: 3

No missing Values

47 Duplicate Rows

| Index | Text | Label |
|---|---|---|
| 1 – 1200 rows | 50 symptom descriptions per label | 24 diseases |

# Project Structure

Data Exploration

Pre-Processing

Modeling

Model Evaluation and Interpretation

Model Deployment

# Models and Results

| Models | Accuracy Score before Hyperparameter Tuning | | |
|---|---|---|---|
| | **CountVectorizer** | **Word2Vec** | **TF-IDF** |
| LinearSVC | 94 % | 9 % | 93 % |
| RandomForest Classifier | 98 % | 64 % | 96 % |

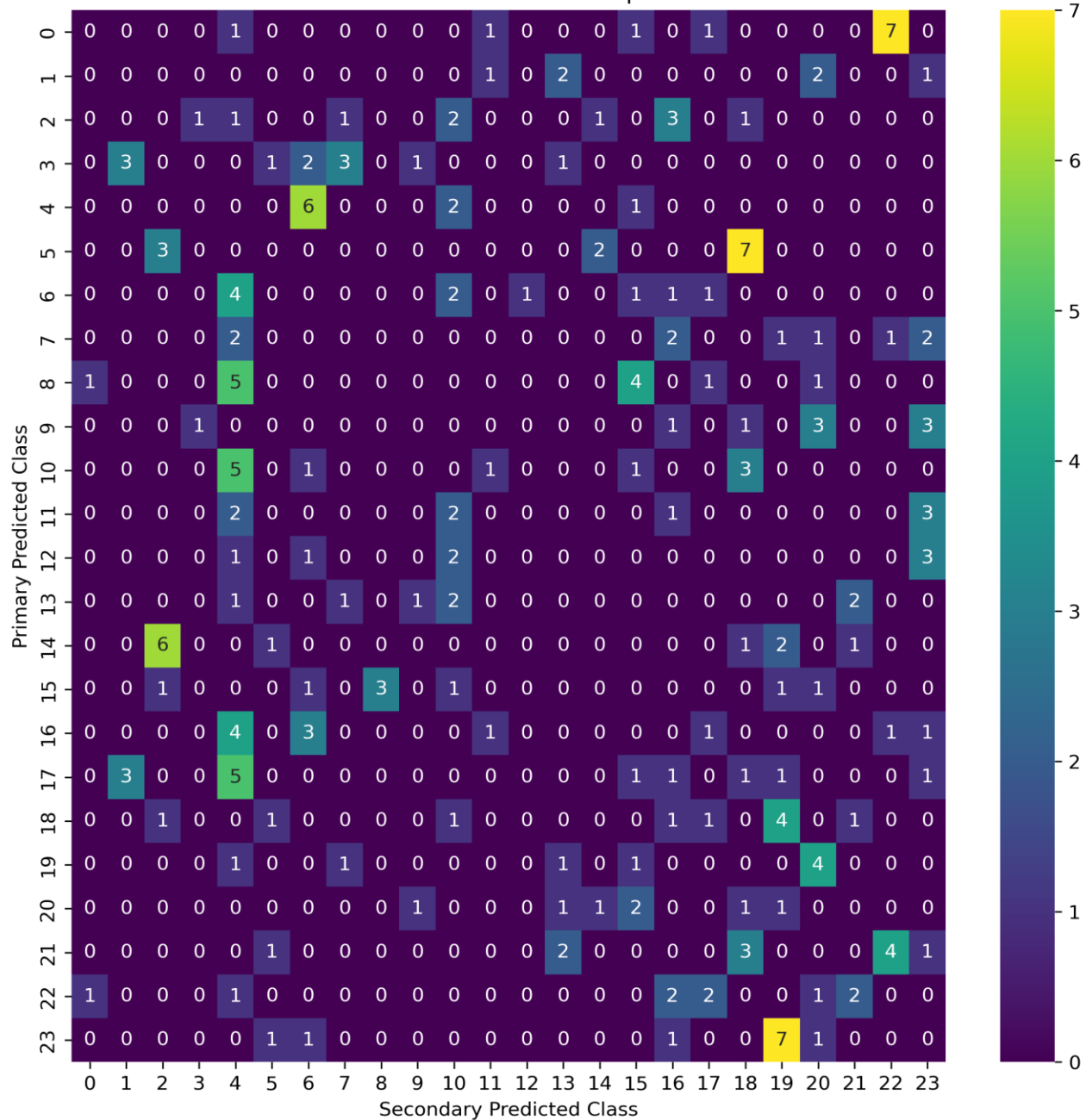| Model with CountVectorizer | Best Parameters | Accuracy Score after Hyperparameter tuning |
|---|---|---|
| LinearSVC | `{'C': 0.001, 'max_iter': 1000}` | 94 % |
| RandomForest Classifier | `{'max_depth': None, 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 100}` | 98 % |

Confusion Matrix for RandomForest

# Model Analysis and Evaluation

- Diseases were accurately predicted for a maximum of 13 times

- Confusion matrix shows very few misclassification

Pairwise Predicted Probabilities Heatmap for RandomForest

# Model Analysis and Evaluation

- The heatmap highlights some potential confusions between diseases

- Maximum number of confusion 7 times

# Conclusion

- Best performing model is the RandomForest Classifier

- Highest accuracy score of 98 %

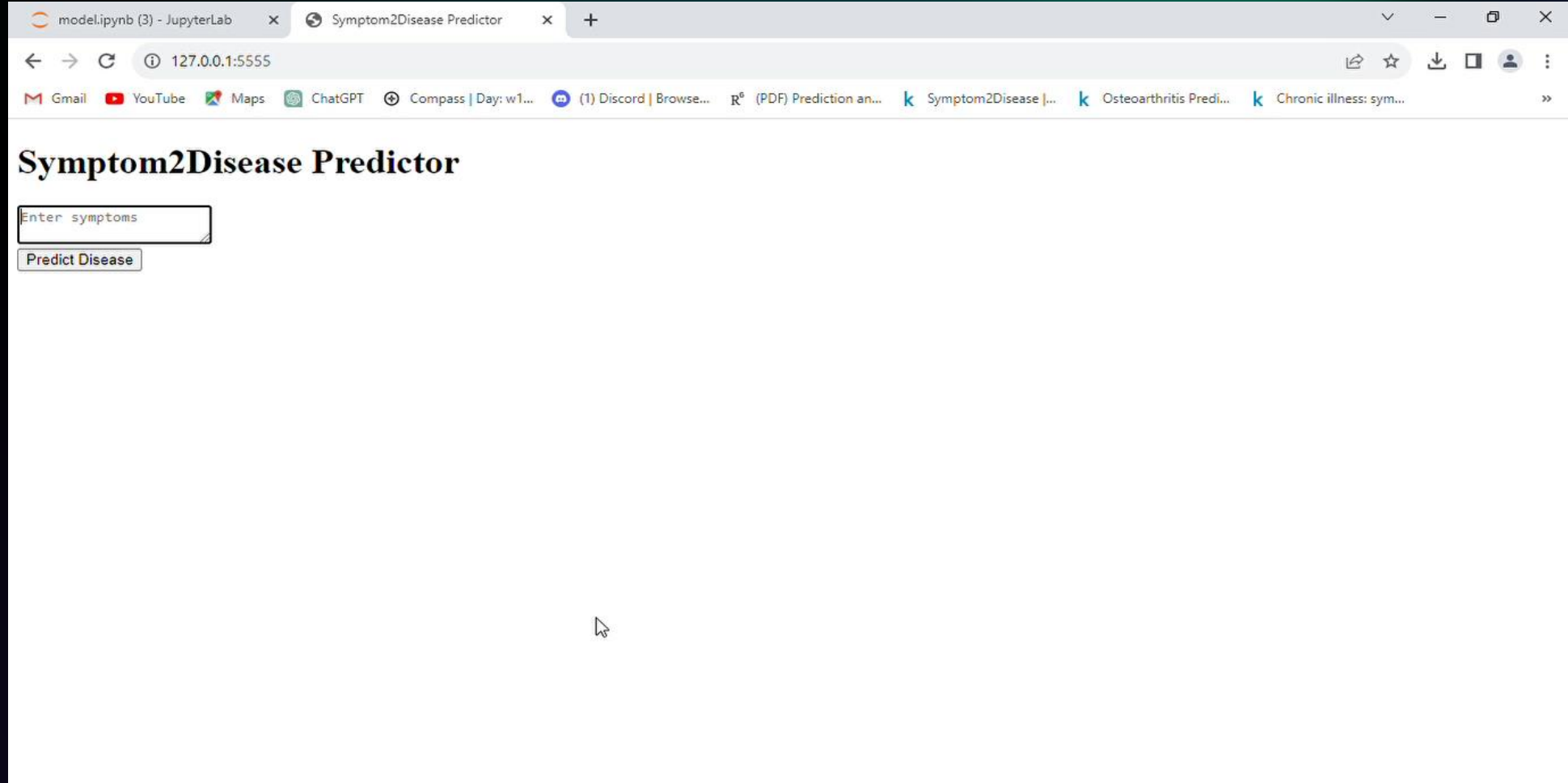- Low variability in its predictions

- Some misclassifications

# Challenges

- Dealing with text

- Determining the relevance of words

- Hyperparameter Tuning is time consuming

# Future Scope

- Refine model performance

- Learn how to effectively remove irrelevant words and noise

# Symptom2Disease App Demonstration

THANK YOU