**Task:**

1. **Data Preprocessing:**
   - The data is checked for null values.
   - The description column is examined for the distribution of description length, which ranges between 2 to 3234 and plotted histplot. Explored the truncation function to use in Deep Learning models to truncate till 500 words.
   - The texts in description column are converted to lowercase, removed punctuations, removed stopwords, tokenized and lemmatized.
   - The labels are encoded with Label Encoder.



```
 1  df['description_length'].describe()

count    1759.000000
mean      289.876066
std       303.843952
min         2.000000
25%       108.000000
50%       203.000000
75%       348.500000
max      3234.000000
Name: description_length, dtype: float64
```

*Figure 1. Distribution of Description Length*

2. **Model Selection:**
   - GridSearch is employed for hyperparameter tuning and finding the better perfoming parameters in Logistic Regression and Random SearchCV is used for Random Forest Classifier.
   - Ensembling learning method, stacking classifier is used, where base models are Logistic Regression, Random Forest Classifier(50 estimators) and Gradient Boosting Classifier(50 estimators) and Logistic Regression as meta-learner.
   - MLFlow is employed for evaluating various model performances using the evaluation metrics accuracy, precision, recall, f1-score.
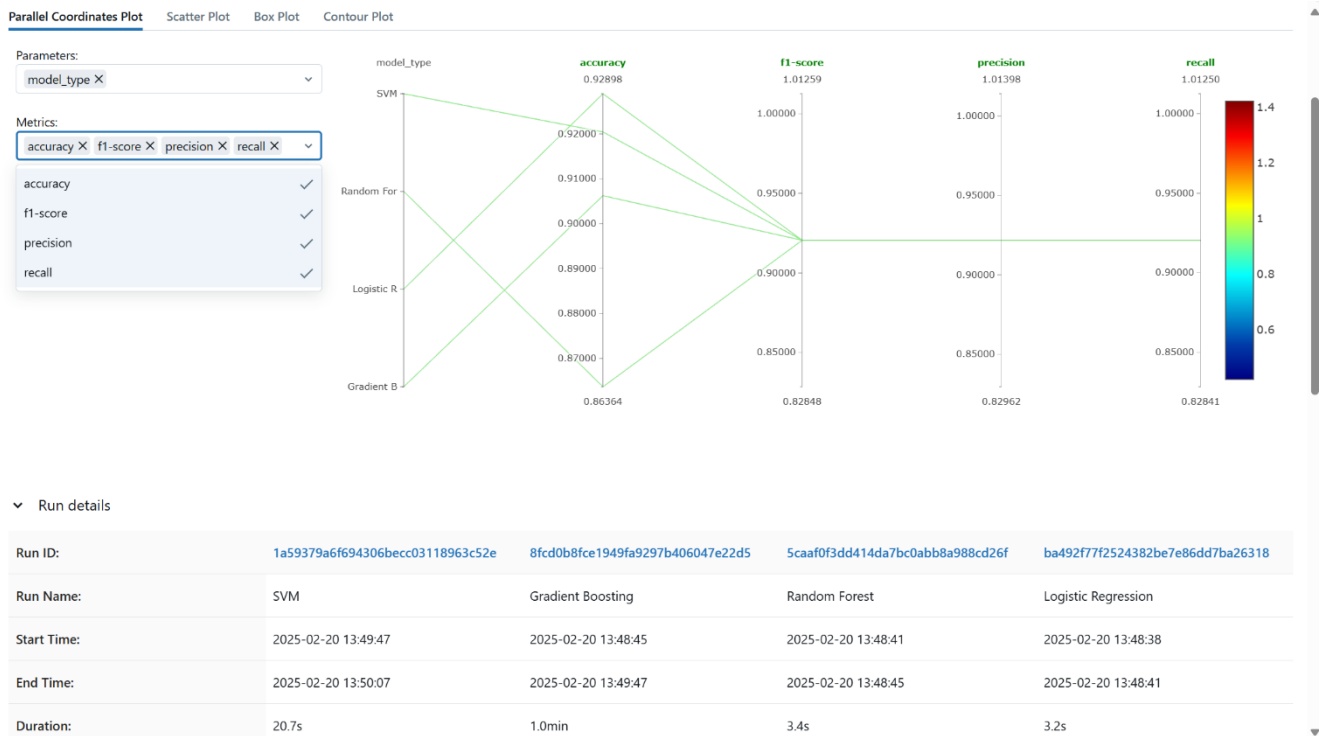


| Run ID: | 1a59379a6f694306becc03118963c52e | 8fcd0b8fce1949fa9297b406047e22d5 | 5caaf0f3dd414da7bc0abb8a988cd26f | ba492f77f2524382be7e86dd7ba26318 |
|---|---|---|---|---|
| Run Name: | SVM | Gradient Boosting | Random Forest | Logistic Regression |
| Start Time: | 2025-02-20 13:49:47 | 2025-02-20 13:48:45 | 2025-02-20 13:48:41 | 2025-02-20 13:48:38 |
| End Time: | 2025-02-20 13:50:07 | 2025-02-20 13:49:47 | 2025-02-20 13:48:45 | 2025-02-20 13:48:41 |
| Duration: | 20.7s | 1.0min | 3.4s | 3.2s |

*Figure 2.Evaluation of Machine Learning Models using MLFlow with accuracy, precision, recall, f1-score*

- The models evaluated are Logistic Regression, Support Vector Machine, Random Forest, Gradient Boosting, in which Logistic Regression showed better performance with accuracy 92.89.
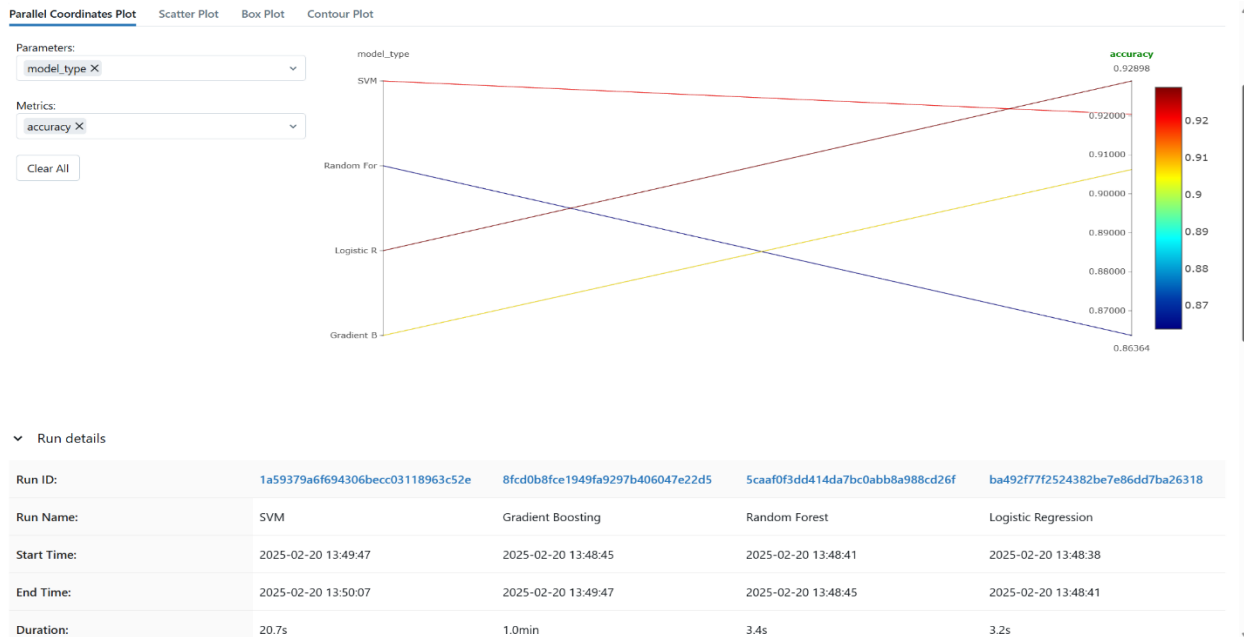


*Figure 3. Evaluation of Machine Learning Models using MLFlow with accuracy*

- Logistic Regression is model is selected and saved for inferencing and exhibiting excellent performance and would able to predict the text related to the 4 medical conditions.



*Figure 4. Model Predictions*