



MENTOR **NESS**

Machine learning internship

Disease prediction

Samuel MATIA

Project presentation

The objective of this project is to develop a predictive model that can accurately classify individuals into diseased or non-diseased categories based on their health attributes. By leveraging machine learning algorithms, we aim to create a reliable tool that healthcare providers can use to assist in disease diagnosis and prognosis.



Dataset overview

We have access to a dataset containing multiple health-related attributes such as cholesterol levels, blood cell counts, hormone levels, and other physiological measurements. The dataset also includes information on whether the individual has been diagnosed with a specific disease or not.

Data shape

Import data

```
1 train = pd.read_csv('Train_data.csv')
2 test= pd.read_csv('test_data.csv')
3
4 print(f'train data dims : {train.shape}')
5 print(f'test data dims : {test.shape}')
```

```
train data dims : (2351, 25)
test data dims : (486, 25)
```

Preprocessing

- Label encoding
- Standardisation

Label Encoding

```
1 disease = ['Diabetes', 'Thalasse', 'Anemia', 'Thromboc']  
2 non_disease = ['Healthy']
```

```
1 train['Disease'] = np.where(train['Disease'].isin(non_disease), 0,1)  
2 test['Disease'] = np.where(test['Disease'].isin(non_disease), 0,1)
```

Standardisation

```
1 y_train = train.pop('Disease')  
2 y_test = test.pop('Disease')
```

```
1 scaler = StandardScaler()  
2 x_train = scaler.fit_transform(train)  
3 x_test = scaler.fit_transform(test)
```

Model

Model

```
1 model = ExtraTreesClassifier(n_estimators=288, random_state = 0)
2 #the best n_estimators finded manually, GridSearchCV was inefficient in this case
```

```
1 model.fit(x_train, y_train)
2 model.score(x_train, y_train)
```

1.0

```
1 predictions = model.predict(x_test)
```


Metrics

Evaluation metrics

```
1 print(f'accuracy : {round(accuracy_score(predictions, y_test)*100, 2)} %')
2 print(f'precision : {round(precision_score(predictions, y_test)*100, 2)} %')
3 print(f'recall : {round(recall_score(predictions, y_test)*100, 2)} %')
4 print(f'f1_score : {round(f1_score(predictions, y_test)*100, 2)} %')
```

```
accuracy : 97.94 %
precision : 98.96 %
recall : 98.96 %
f1_score : 98.96 %
```



Machine learning internship

Samuel MATIA