

# Thyroid Disease Classification Using Machine Language

## 1. INTRODUCTION

### 1.1 Overview

Thyroid disease classification using machine learning involves the development of algorithms that can automatically classify different types of thyroid diseases based on a set of input features. Machine learning algorithms are trained on a dataset of labeled thyroid disease cases, which typically includes information such as the patient's age, sex, thyroid hormone levels, thyroid gland size, and other clinical indicators. Once trained, these algorithms can be used to predict the likelihood of a patient having a specific type of thyroid disease based on their input data.

The first step in developing a machine learning model for thyroid disease classification is to gather a large dataset of labeled thyroid disease cases. This dataset should include information about each patient's thyroid hormone levels, thyroid gland size, and other relevant clinical indicators, as well as the specific type of thyroid disease that the patient was diagnosed with.

Once the dataset is gathered, the next step is to preprocess the data to prepare it for machine learning. This may involve cleaning the data to remove any errors or inconsistencies, normalizing the data to a common scale, and selecting a set of relevant features that are most predictive of thyroid disease type.

Next, a machine learning algorithm is trained on the preprocessed data. There are several different types of machine learning algorithms that can be used for thyroid disease classification, including decision trees, support vector machines, and neural networks.

The algorithm is trained on a subset of the data, using a technique called supervised learning, where it learns to map the input features to the correct thyroid disease classification.

Once the algorithm is trained, it can be tested on a separate set of labeled data to evaluate its performance. The performance of the algorithm can be measured using metrics such as accuracy, precision, recall, and F1 score. If the algorithm performs well on the test data, it can be deployed in a clinical setting to help diagnose patients with thyroid disease.

Overall, the use of machine learning algorithms for thyroid disease classification has the potential to improve the accuracy and speed of diagnosis, allowing for earlier intervention and better patient outcomes.

However, it is important to note that machine learning algorithms are not a substitute for clinical expertise and should be used in conjunction with other diagnostic tools and medical professionals

## 1.2 Purpose

Thyroid disease ML (Machine Learning) algorithms can have several uses, including:

**Diagnosis:** ML algorithms can analyze medical data such as blood tests, ultrasound scans, and biopsy results to diagnose thyroid disease accurately. By training the algorithm on large datasets of patient information, it can learn to recognize patterns and make predictions about the likelihood of disease.

**Prognosis:** Once a diagnosis has been made, ML algorithms can help predict the progression of thyroid disease and the likelihood of complications. This information can help doctors plan treatment and monitor patients more effectively.

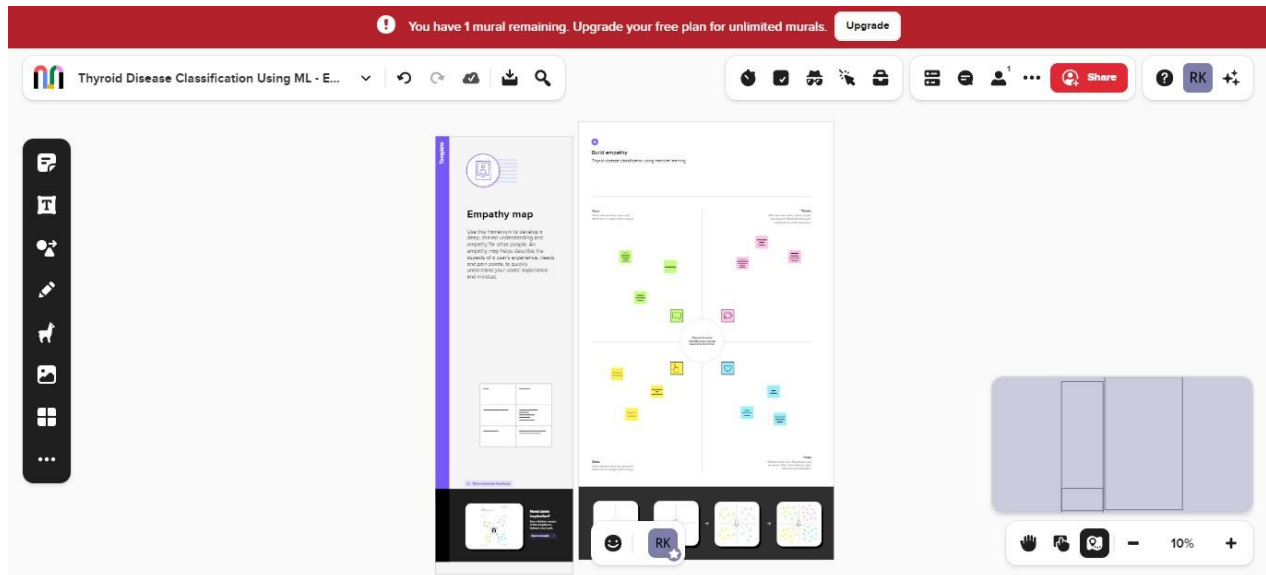
**Treatment optimization:** ML algorithms can also help optimize treatment plans for patients with thyroid disease. By analyzing patient data, including factors such as age, gender, and medical history, the algorithm can recommend personalized treatment plans that are tailored to the individual patient's needs.

**Drug discovery:** ML algorithms can help researchers discover new drugs for the treatment of thyroid disease. By analyzing large datasets of molecular and clinical data, algorithms can identify potential drug candidates and predict their effectiveness.

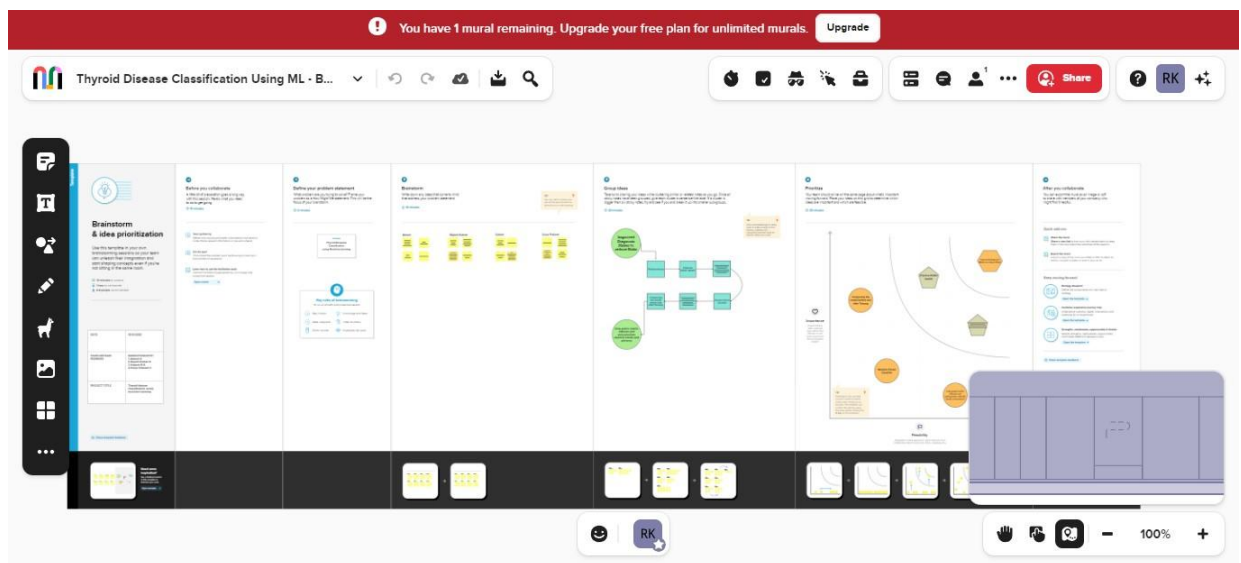
Overall, the use of ML in thyroid disease can lead to more accurate diagnoses, better treatment outcomes, and ultimately, improved patient care.

## 2. Problem Definition & Design Thinking

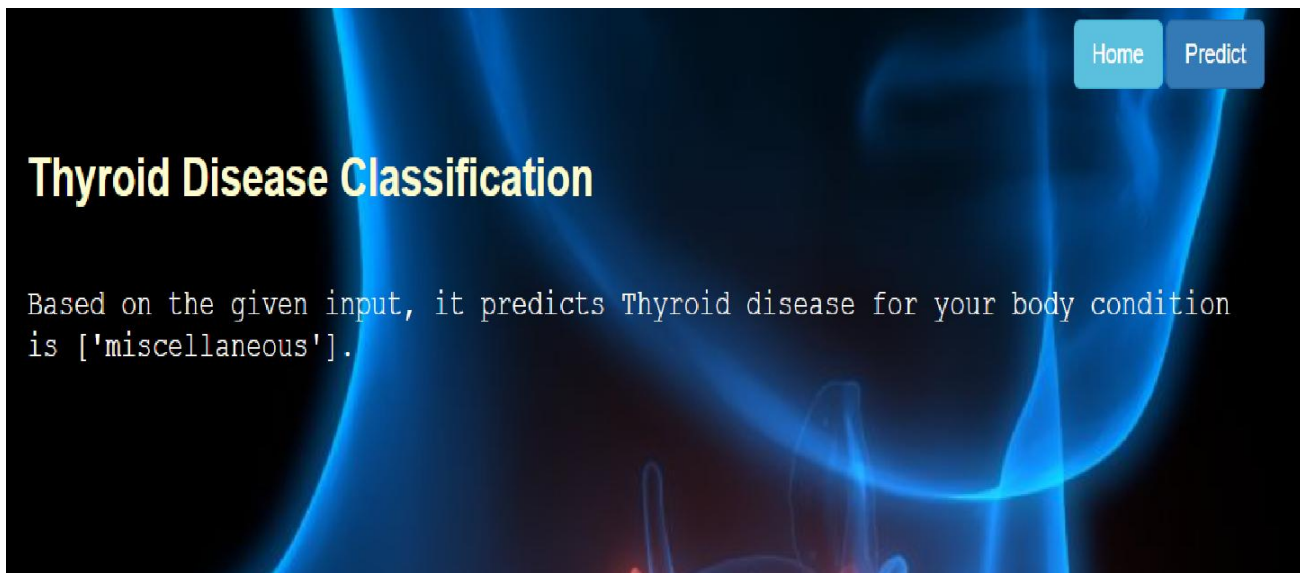
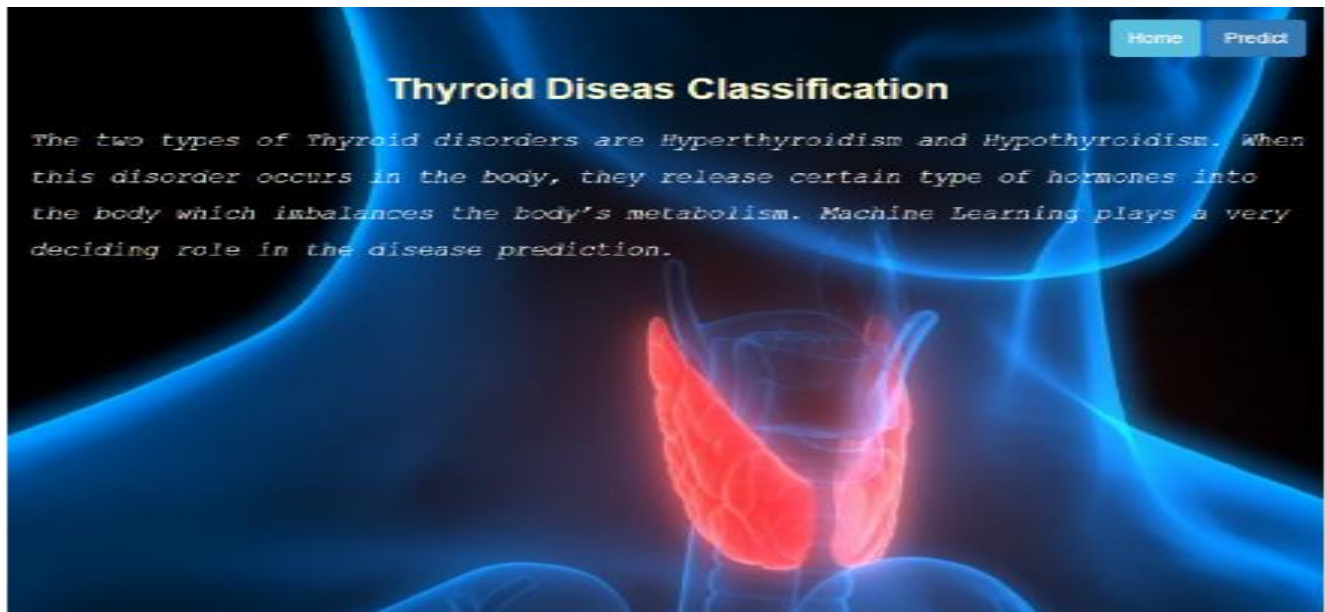
### 2.1 Empathy Map



### 2.2 Ideation & Brainstorming Map



### 3. RESULT



## 4. ADVANTAGES & DISADVANTAGES

### Advantages:

**Increased accuracy:** ML algorithms can analyze large amounts of patient data to identify patterns and predict the likelihood of thyroid disease. This can improve the accuracy of diagnosis and reduce the risk of misdiagnosis.

**Personalized treatment:** ML can help identify the optimal treatment for each patient based on their individual characteristics, such as age, gender, and other health conditions.

**Time-saving:** ML algorithms can analyze patient data and provide a diagnosis faster than traditional methods, which can speed up treatment and improve patient outcomes.

**Early detection:** ML can help identify early signs of thyroid disease, allowing for early intervention and treatment.

### Disadvantages:

**Data bias:** ML algorithms are only as good as the data they are trained on. If the data is biased or incomplete, the algorithm may produce inaccurate results.

**Complexity:** ML algorithms can be complex and difficult to understand, which may make it challenging for healthcare providers to interpret and use the results.

**Lack of transparency:** Some ML algorithms may be considered a "black box" because they do not provide a clear explanation of how they arrive at their conclusions. This can make it difficult to trust and validate the results.

**Cost:** Implementing ML systems can be expensive, which may limit access for some healthcare providers and patients.

## 5. APPLICATIONS

Machine learning (ML) can be applied to various areas related to the diagnosis, treatment, and management of thyroid diseases. Here are some of the key areas where ML can be applied in thyroid diseases:

**Diagnosis:** ML algorithms can be trained to analyze patient data, including blood tests, imaging studies, and symptoms, to help diagnose thyroid diseases. ML can help identify patterns and associations that may not be apparent to humans, leading to faster and more accurate diagnosis.

**Risk prediction:** ML algorithms can be used to predict the risk of developing thyroid diseases based on patient characteristics and medical history. This can help identify patients who may benefit from early intervention and preventive measures.

**Treatment optimization:** ML algorithms can be used to identify the most effective treatment for each patient based on their individual characteristics, such as age, gender, and other health conditions. This can help improve patient outcomes and reduce healthcare costs.

**Prognosis:** ML algorithms can be used to predict the likelihood of disease progression and the long-term outcomes of treatment. This can help healthcare providers make informed decisions about treatment and management.

**Image analysis:** ML algorithms can be used to analyze thyroid images, such as ultrasounds and CT scans, to help identify and characterize thyroid nodules and tumors. This can help improve the accuracy of diagnosis and reduce the need for invasive procedures.

**Electronic health records (EHR):** ML algorithms can be used to analyze EHR data to identify patterns and associations related to thyroid diseases, including risk factors, treatment outcomes, and disease progression. This can help improve patient care and management.

## 6. CONCLUSION

In conclusion, machine learning (ML) has shown great promise in improving the diagnosis, treatment, and management of thyroid diseases. ML algorithms can process and analyze large amounts of patient data to identify patterns and make predictions that can improve the accuracy of diagnosis, personalize treatment, and predict disease progression. ML can also help healthcare providers to automate repetitive tasks, speed up the diagnosis process, and reduce the risk of misdiagnosis. However, there are also challenges associated with using ML in thyroid diseases, including data bias, complexity, lack of transparency, and cost. Despite these challenges, the potential benefits of ML in thyroid diseases are significant, and continued research and development in this area can lead to improved patient outcomes and better management of thyroid diseases. Overall, the application of ML in

thyroid diseases has the potential to revolutionize the field of endocrinology and provide a new level of precision and personalized care for patients with thyroid diseases.

## 7. FUTURE SCOPE

The future of thyroid disease management using machine learning (ML) is promising. Here are some potential future directions for this area of research:

**Integration with other technologies:** As ML continues to advance, it can be integrated with other technologies such as genetic testing, telemedicine, and wearable devices to provide a more holistic approach to thyroid disease management.

**Personalized medicine:** ML can help identify individual patient characteristics that influence thyroid disease management and develop personalized treatment plans tailored to each patient's needs.

**Improved accuracy:** The accuracy of ML algorithms in diagnosing and managing thyroid diseases will continue to improve as more data becomes available and the algorithms become more sophisticated.

**Patient monitoring:** ML can be used to monitor patient progress and adjust treatment plans in real-time based on the patient's response to therapy.

**Predictive analytics:** ML can be used to predict disease outcomes, such as the risk of recurrence or progression, allowing for earlier intervention and more effective management.

**Data sharing:** Collaboration between healthcare providers and researchers can lead to the development of large-scale databases that can be used to improve the accuracy of ML algorithms in thyroid disease management.

Overall, the future of thyroid disease management using ML is bright, with the potential to improve patient outcomes, reduce healthcare costs, and provide a new level of personalized care for patients with thyroid diseases.

## 8. APPENDIX

### A.SOURCE CODE

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn import preprocessing
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import confusion_matrix, accuracy_score, classification_report

data=pd.read_csv('thyroidDF.csv')

data.head()

data.tail()

data.shape

data.info()

data.isnull().sum()

data.drop(['TSH_measured', 'T3_measured', 'TT4_measured', 'T4U_measured', 'FTI_measured', 'TBG_measured', 'referral_source', 'patient_id'], axis=1, inplace = True)

data.head()

diagnoses = {'A': "hyperthyroid conditions",
             "B": "hyperthyroid conditions",
             "C": "hyperthyroid conditions",
             "D": "hyperthyroid conditions",
             "E": "hypothyroid conditions",
             "F": "hypothyroid conditions",
             "G": "hypothyroid conditions",
             "H": "hypothyroid conditions",
             "I": "binding protein",
             "J": "binding protein",
             "K": "general health",
             "L": "replacement therapy",
             "M": "replacement therapy",
             "N": "replacement therapy",
```



```
        "O": "antithyroid treatment",
        "P": "antithyroid treatment",
        "Q": "antithyroid treatment",
        "R": "miscellaneous",
        "S": "miscellaneous",
        "T": "iscellaneous"}
data['target'] = data['target'].map(diagnoses) #remapping
```

```
data
```

```
data.isnull().sum()
```

```
data.dropna(subset=['target'],inplace=True)
```

```
data['target'].value_counts()
```

```
data['target'].isnull().sum()
```

```
data.head()
```

```
data.describe()
```

```
data['age']=np.where((data.age>100), np.nan, data.age)
```

```
data
```

```
x=data.iloc[:,0:-1]
```

```
y= data.iloc[:,-1]
```

```
data.isnull().sum()
```

```
x['sex'].unique()
```

```
x['sex'].replace(np.nan, 'F', inplace=True)
```

```
x['sex'].value_counts()
```

```
x.isnull().sum()
```

```
data.info()
```

```
x["age"]=x["age"].astype("float")
```

```
x["TSH"]=x["TSH"].astype("float")
```

```

x["T3"]=x["T3"].astype("float")
x['TT4']=x["TT4"].astype("float")
x["T4U"]=x["T4U"].astype("float")
x['FTI']=x["FTI"].astype("float")
x["TBG"]=x["TBG"].astype("float")

from sklearn.preprocessing import OrdinalEncoder, LabelEncoder
ordinal_encoder =OrdinalEncoder(dtype ='int64')
x.iloc[:, 1:16] =ordinal_encoder.fit_transform(x.iloc[:, 1:16])

x.head()

x.replace(np.nan, '0', inplace=True)
x.head()

label_encoder = LabelEncoder()
y_dt= label_encoder.fit_transform(y)

y=pd.DataFrame(y_dt, columns=['target'])
y

y.value_counts(normalize=True)

[]

import seaborn as sns
corrmat = x.corr()

f, ax = plt.subplots(figsize =(9, 8))

sns.heatmap(corrmat, ax = ax, cmap ="YlGnBu", linewidths = 0.1)

from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.20, random_state=0)

y_train.value_counts()

from imblearn.over_sampling import SMOTE
os = SMOTE(random_state=0,k_neighbors=1)
x_bal,y_bal=os.fit_resample(x_train,y_train)
x_test_bal,y_test_bal=os.fit_resample(x_test,y_test)

from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
x_bal = sc.fit_transform(x_bal)
x_test_bal = sc.transform(x_test_bal)

```

```
x_bal
```

```
x_test_bal
```

```
y_bal.value_counts()
```

```
columns=['age','sex','on_thyroxine','query_on_thyroxine','on_antithyroid_meds','sick','pregnant','thyroid_surgery','I131_treatment','query_hypothyroid','query_hyperthyroid','lithium','goitre','tumor','hypopituitary','psych','TSH','T3','TT4','T4U','FTI','TBG']
```

```
x_test_bal= pd.DataFrame(x_test_bal,columns=columns)
```

```
x_bal= pd.DataFrame(x_bal,columns=columns)
```

```
x_bal
```

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, classification_report
rfr = RandomForestClassifier().fit(x_bal,y_bal)
y_pred = rfr.predict(x_test_bal)
accuracy_score(y_test_bal,y_pred)
x_bal.shape,y_bal.shape,x_test_bal.shape,y_test_bal.shape
```

```
test_score=accuracy_score(y_test_bal,y_pred)
test_score
```

```
train_score = accuracy_score(y_bal,rfr.predict(x_bal))
train_score
```

```
from sklearn.inspection import permutation_importance
results = permutation_importance(rfr,x_bal,y_bal, scoring="accuracy")
```

```
feature_importance=['age','sex','on_thyroxine','query_on_thyroxine','on_antithyroid_meds','sick','pregnant','thyroid_surgery','I131','query_thyroid','query_hyperthyroid','lithium','goitre','tumor','hypopituitary','psych','TSH','T3','TT4','T4U','FTI','TBG']
importance = results.importances_mean
importance = np.sort(importance)
for i,v in enumerate(importance):
    i=feature_importance[i]
print('feature: (<20) Score: {}'.format(i,v))
```

```
plt.figure(figsize=(10,10))
plt.bar(x=feature_importance, height = importance)
plt.xticks(rotation=30, ha='right')
plt.show()
```

```
x_bal.drop(['age', 'sex', 'on_thyroxine', 'query_on_thyroxine', 'on_antithyroid_meds', 'sick', 'pregnant', 'thyroid_surgery', 'I131_treatment', 'query_hypothyroid', 'query_hyperthyroid', 'lithium', 'goitre', 'tumor', 'hypopituitary', 'psych', 'TSH', 'T3', 'TT4', 'T4U', 'FTI', 'TBG'])
```

```
x.head()
```

```
x_bal.drop(['age', 'sex', 'on_thyroxine', 'query_on_thyroxine', 'on_antithyroid_meds', 'sick', 'pregnant', 'thyroid_surgery', 'I131_treatment', 'query_hypothyroid', 'query_hyperthyroid', 'lithium', 'goitre', 'tumor', 'hypopituitary', 'psych', 'TSH', 'T3', 'TT4', 'T4U', 'FTI', 'TBG'])
```

```
data.info()
```

```
import seaborn as sns
corrmat = x.corr()
```

```
f,ax =plt.subplots(figsize =(9,8))
sns.heatmap(corrmat, ax = ax, cmap= "YlGnBu", linewidths = 0.1)
```

```
from sklearn.ensemble import RandomForestClassifier
rfr1 = RandomForestClassifier().fit(x_train,y_train)
y_pred = rfr1.predict(x_test)
rfr1 = RandomForestClassifier()
```

```
rfr1.fit(x_train, y_train)
```

```
y_pred = rfr1.predict(x_test)
```

```
print(classification_report(y_test, y_pred))
```

```
train_score = accuracy_score(y_train,rfr1.predict(x_train))
```

```
train_score
```

```
from sklearn.tree import DecisionTreeClassifier
dt = DecisionTreeClassifier().fit(x_train,y_train)
y_pred = dt.predict(x_test)
dt = DecisionTreeClassifier()
```

```
dt.fit(x_train, y_train)
```

```
y_pred1 = dt.predict(x_test)
```

```
print(classification_report(y_test, y_pred))
```

```
train_score = accuracy_score(y_train,dt.predict(x_train))
```

```
train_score
```

```
import tensorflow.keras
from tensorflow.keras.models import Sequential
```

```
from tensorflow.keras.layers import Dense

classifier.add(Dense(units=256, activation='relu'))

classifier.add(Dense(units=1, activation='sigmoid'))

classifier.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

classifier.summary()

model_history = classifier.fit(x_bal, y_bal, batch_size=100, validation_split=0.2, epochs=50)

dt.predict([[0,0,0,0,0.00,0.0,0.0,50,1,1.2,0.8,150,70,80,7.2,3.4,0.8,0,2,3,4,6]])

print(classification_report(y_test,y_pred))
```