# Machine learning algorithm for detecting thyroid-related disorders to increase thyroid illness diagnosis accuracy

**Advance Machine Learning / Spring 2024 / Prof. Ashok Kumar Patel**

**By Group – 11**

**Ajeet Singh – 48**
**Manas Vishal – 58**
**Rohith Kumar – 13**

# Problem Statement:

Machine learning techniques have been increasingly used in the medical field to improve the accuracy of diagnoses and treatment plans. The use of an ML algorithm for detecting thyroid-related disorders to increase thyroid illness diagnosis accuracy. This approach involves using data-based decisions for diagnosing thyroid dysfunction, which is a classification problem that can be solved using Machine Learning techniques.

# About the Datasets:

The dataset contains records of thyroid diagnoses from the Garvan Institute, collected between 1984 and early 1987. Each record consists of 25 attribute values, followed by a class. The attributes include various medical conditions, patient characteristics, and laboratory measurements related to thyroid function. The class contains two possible values: negative and sick-euthyroid.

Observations: 3163

Features: 25

Class: 1 (with two distinct values)

```
COLUMN                      VALUES

class:                      sick-euthyroid,negative

age:                        continuous,?
sex:                        M,F,?
on_thyroxine:               f,t
query_on_thyroxine:         f,t
on_antithyroid_medication:  f,t
thyroid_surgery:            f,t
query_hypothyroid:          f,t
query_hyperthyroid:         f,t
pregnant:                   f,t
sick:                       f,t
tumor:                      f,t
lithium:                    f,t
goitre:                     f,t
TSH_measured:               f,t
TSH:                        continuous,?
T3_measured:                f,t
T3:                         continuous,?
TT4_measured:               f,t
TT4:                        continuous,?
T4U_measured:               f,t
T4U:                        continuous,?
FTI_measured:               f,t
FTI:                        continuous,?
TBG_measured:               f,t
TBG:                        continuous,?
```

# Solution Approach:

Given the complexity and the nature of the thyroid disease dataset, which includes a mix of continuous and categorical variables, and considering the need for accurate prediction of multiple classes of thyroid conditions, a combination of machine learning algorithms would be beneficial. The choice of algorithms should be based on their ability to handle both types of data and their performance in classification tasks, especially in the context of imbalanced classes.

1. Random Forest (RF): Random Forest is a powerful ensemble learning method that can handle both continuous and categorical data. It is known for its ability to avoid overfitting and its robustness to outliers.
2. Gradient Boosting Machine (GBM): GBM is another ensemble method that builds a series of weak learners (typically decision trees) to form a strong predictive model. It is particularly effective in handling imbalanced datasets and can capture complex patterns in the data.
3. Support Vector Machine (SVM): SVMs are effective in high-dimensional spaces and are versatile as different Kernel functions can be specified for the decision function. They are particularly useful when the dataset is not linearly separable.
4. AdaBoost: AdaBoost is an ensemble method that combines multiple weak learners to create a strong learner. It is effective in handling noisy data and can improve the accuracy of the model.
5. Deep Learning Models (ANN): For a more complex and nuanced understanding of the data, deep learning models like ANN networks can be used. These models can capture complex patterns and dependencies in the data, which might be beneficial for predicting thyroid conditions.

It would be beneficial to experiment with these algorithms, possibly in combination, and use cross-validation to evaluate their performance. The choice of the best algorithm or combination of algorithms would depend on the specific characteristics of the dataset and the performance metrics (accuracy, precision, recall, etc.) that are most relevant to the problem at hand.

# Loading Dataset:

Load the csv file of thyroid dataset. Since the file does not have columns names and so we need to specify while loading it into dataframe. It has 25 features, 1 class and 3163 observations. It does not show any missing values initially, but we do see some unknow values like '?'.



# Understanding Missing Values:

We need to replace the unknow values '?' with NaN so that we can understand correct stats of the missing numbers. And as you can see in below percentages of missing values, there are some missing values in age, sex, TSH, T3, TT4, T4U, FTI and TBG. But class values do not have any missing values which is good.

# Handling Missing Values and Encoding Categorical Values:

- We can remove the TBG column as it has 91.77 % of missing values and will not be helpful for training model.
- All other missing values should be replaced with the Mean value of that column.
- Since most of the column values are categorical types and so we need to convert/map/encode to numerical values as listed below:

  'sick-euthyroid': 1

  'negative': 0

  't': 1

  'f': 0

  'y': 1

  'n': 0

  'F': 1

  'M': 0

- Convert the remaining measured features from Object to Float

# Finding Correlations between Features and Class:

- There is strong correlation for T3, T4U, sick, age, T3_measured, TT4, TSH_measured, TBG_measured.

- Some of the features are not very well correlated. like, Tumor, lithium, thyroid_surgery, query_on_thyroxine, FTI, query_hypothyroid, on_antothyroid_medication, goitre.

# Dividing Datasets into Training and Testing sets:

- We divided datasets into Train set and Test set. Making sure each set represent proper proportion of each Classes

  Train Set:    (2530, 25)

  0   2296

  1   234

  Test Set:    (633, 25)

  0   574

  1   59

- Created X and Y variables for each dataset to use into Model Function

  X_train = hold all features values for train set

  Y_train = hold only class values for train set

  X_test = hold all features values for test set

  Y_test = hold only class values for test set

# Training Using different ML Models:

- We will train and test datasets using different models as listed below. And will capture the evaluation metrics for comparison and understand which model is better.
    - Network Model
        - ANN
    - Statistical Model
        - SVC
    - Tree Based Model
        - XGBoost
        - GBM
        - Random Forest
        - AdaBoost

- We are going to follow the steps below to train/test/evaluate for each model.

    - Import required model's library
    - Train Model using Training set
        - trained_model = model.fit(X_train,Y_Train)
    - Predict Model using Testing set
        - Y_predict = trained_model.predict(X_test)
    - Evaluate Model Performance
        - accuracy_score(Y_test, Y_predict)
        - precision_score(Y_test, Y_predict)
        - recall_score(Y_test, Y_predict)
        - f1_score(Y_test, Y_predict)
    - Also create Confusion Matrix and plot a heatmap
        - confusion_matrix(Y_test, Y_predict)

## XGBoost Model

```
[62]: from xgboost import XGBClassifier

      model_xgboost = XGBClassifier(objective='binary:logistic', eval_metric='auc')
      model_xgboost.fit(X_train, Y_train)
```

```
[62]: ▸ XGBClassifier
```

```
[63]: #Predict using Test dataset

      Y_pred_xgboost = model_xgboost.predict(X_test)
```

```
[64]: # Evaluate using various Metrics

      from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, confusion_matrix

      # Calculate metrics
      eval_metrics.loc[0, 'Accuracy'] = accuracy_score(Y_test, Y_pred_xgboost)
      eval_metrics.loc[0, 'Precision'] = precision_score(Y_test, Y_pred_xgboost, average='binary') # Use 'binary' for binary classification
      eval_metrics.loc[0, 'Recall'] = recall_score(Y_test, Y_pred_xgboost, average='binary')
      eval_metrics.loc[0, 'F1_Score'] = f1_score(Y_test, Y_pred_xgboost, average='binary')

      print(eval_metrics[0:1])
```

```
        Method_Name  Accuracy  Precision    Recall  F1_Score
      0      XGBoost  0.966825   0.816667  0.830508  0.823529
```

```
[65]: matrix = confusion_matrix(Y_test, Y_pred_xgboost)
      df_confusion = pd.DataFrame(matrix, index=['sick-euthyroid','negative'],columns=['sick-euthyroid','negative'])
      df_confusion
```

```
[65]:                   sick-euthyroid  negative
      sick-euthyroid             563        11
           negative               10        49
```

## Random Forest Model

```
[40]: from sklearn.ensemble import RandomForestClassifier

      model_randomforest = RandomForestClassifier(n_estimators=100, random_state=42)
      model_randomforest.fit(X_train, Y_train)
```

```
[40]: ▾        RandomForestClassifier
      RandomForestClassifier(random_state=42)
```

```
[41]: Y_pred_rf = model_randomforest.predict(X_test)
```

```
[42]: # Evaluate using various Metrics

      from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

      # Calculate metrics
      eval_metrics.loc[2, 'Accuracy'] = accuracy_score(Y_test, Y_pred_rf)
      eval_metrics.loc[2, 'Precision'] = precision_score(Y_test, Y_pred_rf, average='binary') # Use 'binary' for binary classification
      eval_metrics.loc[2, 'Recall'] = recall_score(Y_test, Y_pred_rf, average='binary')
      eval_metrics.loc[2, 'F1_Score'] = f1_score(Y_test, Y_pred_rf, average='binary')

      print(eval_metrics[2:3])
```

```
          Method_Name  Accuracy  Precision    Recall  F1_Score
      2  Random_Forest  0.960506   0.765625  0.830508  0.796748
```

```
[43]: matrix = confusion_matrix(Y_test, Y_pred_rf)
      df_confusion = pd.DataFrame(matrix, index=['sick-euthyroid','negative'],columns=['sick-euthyroid','negative'])
      df_confusion
```

```
[43]:                   sick-euthyroid  negative
      sick-euthyroid             559        15
           negative               10        49
```

## Gradient Boosting Machine (GBM) Model

```python
[50]: from sklearn.ensemble import GradientBoostingClassifier

      model_gbm = GradientBoostingClassifier(n_estimators=100, learning_rate=1.0, max_depth=1, random_state=42)

      model_gbm.fit(X_train, Y_train)
```

```
[50]: ▾              GradientBoostingClassifier
      GradientBoostingClassifier(learning_rate=1.0, max_depth=1, random_state=42)
```

```python
[51]: #Predict using Test dataset

      Y_pred_gbm = model_gbm.predict(X_test)
```

```python
[52]: # Evaluate using various Metrics

      from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

      # Calculate metrics
      eval_metrics.loc[4, 'Accuracy'] = accuracy_score(Y_test, Y_pred_gbm)
      eval_metrics.loc[4, 'Precision'] = precision_score(Y_test, Y_pred_gbm, average='binary') # Use 'binary' for binary classification
      eval_metrics.loc[4, 'Recall'] = recall_score(Y_test, Y_pred_gbm, average='binary')
      eval_metrics.loc[4, 'F1_Score'] = f1_score(Y_test, Y_pred_gbm, average='binary')


      print(eval_metrics[4:5])
```

```
  Method_Name  Accuracy  Precision    Recall  F1_Score
4         GBM  0.952607   0.688312  0.898305  0.779412
```

```python
[53]: matrix = confusion_matrix(Y_test, Y_pred_gbm)
      df_confusion = pd.DataFrame(matrix, index=['sick-euthyroid','negative'],columns=['sick-euthyroid','negative'])
      df_confusion
```

```
[53]:                  sick-euthyroid   negative

      sick-euthyroid         550           24

            negative           6           53
```

## ANN Model

```python
[ ]: import tensorflow as tf

     model_ann = tf.keras.models.Sequential()

     # Adding the input layer and the first hidden layer
     model_ann.add(tf.keras.layers.Dense(units=6, activation='relu'))

     # Adding the second hidden layer
     model_ann.add(tf.keras.layers.Dense(units=6, activation='relu'))

     # Adding the output layer
     model_ann.add(tf.keras.layers.Dense(units=1, activation='sigmoid'))

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

     model_ann.fit(X_train, Y_train, batch_size=32, epochs=100)
```

```python
[36]: #Predict using Test dataset

      Y_pred_ann = model_ann.predict(X_test)
      Y_pred_ann = np.round(Y_pred_ann).flatten()
```

```
20/20 ━━━━━━━━━━━━━━━ 0s 307us/step
```

```python
[37]: # Evaluate using various Metrics

      from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

      # Calculate metrics
      eval_metrics.loc[1, 'Accuracy'] = accuracy_score(Y_test, Y_pred_ann)
      eval_metrics.loc[1, 'Precision'] = precision_score(Y_test, Y_pred_ann, average='binary') # Use 'binary' for binary classification
      eval_metrics.loc[1, 'Recall'] = recall_score(Y_test, Y_pred_ann, average='binary')
      eval_metrics.loc[1, 'F1_Score'] = f1_score(Y_test, Y_pred_ann, average='binary')

      print(eval_metrics[1:2])
```

```
  Method_Name  Accuracy  Precision    Recall  F1_Score
1         ANN  0.951027       0.85  0.576271  0.686869
```

```python
[38]: matrix = confusion_matrix(Y_test, Y_pred_ann)
      df_confusion = pd.DataFrame(matrix, index=['sick-euthyroid','negative'],columns=['sick-euthyroid','negative'])
      df_confusion
```

```
[38]:                  sick-euthyroid   negative

      sick-euthyroid         568            6

            negative          25           34
```

## Support Vector Machine (SVM/SVC) Model

```
[45]: from sklearn import svm

      model_svc = svm.SVC(kernel='linear') # Linear Kernel
      model_svc.fit(X_train, Y_train)
```

```
[45]:  ▾          SVC
      SVC(kernel='linear')
```

```
[46]: #Predict using Test dataset

      Y_pred_svc = model_svc.predict(X_test)
```

```
[47]: # Evaluate using various Metrics

      from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

      # Calculate metrics
      eval_metrics.loc[3, 'Accuracy'] = accuracy_score(Y_test, Y_pred_svc)
      eval_metrics.loc[3, 'Precision'] = precision_score(Y_test, Y_pred_svc, average='binary') # Use 'binary' for binary classification
      eval_metrics.loc[3, 'Recall'] = recall_score(Y_test, Y_pred_svc, average='binary')
      eval_metrics.loc[3, 'F1_Score'] = f1_score(Y_test, Y_pred_svc, average='binary')

      print(eval_metrics[3:4])
```

```
       Method_Name  Accuracy  Precision    Recall  F1_Score
      3         SVC  0.957346       0.75  0.813559  0.780488
```

```
[48]: matrix = confusion_matrix(Y_test, Y_pred_svc)
      df_confusion = pd.DataFrame(matrix, index=['sick-euthyroid','negative'],columns=['sick-euthyroid','negative'])
      df_confusion
```

```
[48]:
```

|                | sick-euthyroid | negative |
|----------------|----------------|----------|
| sick-euthyroid | 558            | 16       |
| negative       | 11             | 48       |

## AdaBoost Model

```
[55]: from sklearn.ensemble import AdaBoostClassifier

      model_ada = AdaBoostClassifier()
      model_ada.fit(X_train, Y_train)
```

```
[55]:  ▾ AdaBoostClassifier
      AdaBoostClassifier()
```

```
[56]: #Predict using Test dataset

      Y_pred_ada = model_ada.predict(X_test)
```

```
[57]: # Evaluate using various Metrics

      from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

      # Calculate metrics
      eval_metrics.loc[5, 'Accuracy'] = accuracy_score(Y_test, Y_pred_ada)
      eval_metrics.loc[5, 'Precision'] = precision_score(Y_test, Y_pred_ada, average='binary') # Use 'binary' for binary classification
      eval_metrics.loc[5, 'Recall'] = recall_score(Y_test, Y_pred_ada, average='binary')
      eval_metrics.loc[5, 'F1_Score'] = f1_score(Y_test, Y_pred_ada, average='binary')

      print(eval_metrics[5:6])
```

```
       Method_Name  Accuracy  Precision    Recall  F1_Score
      5    AdaBoost  0.957346   0.735294  0.847458  0.787402
```

```
[58]: matrix = confusion_matrix(Y_test, Y_pred_ada)
      df_confusion = pd.DataFrame(matrix, index=['sick-euthyroid','negative'],columns=['sick-euthyroid','negative'])
      df_confusion
```
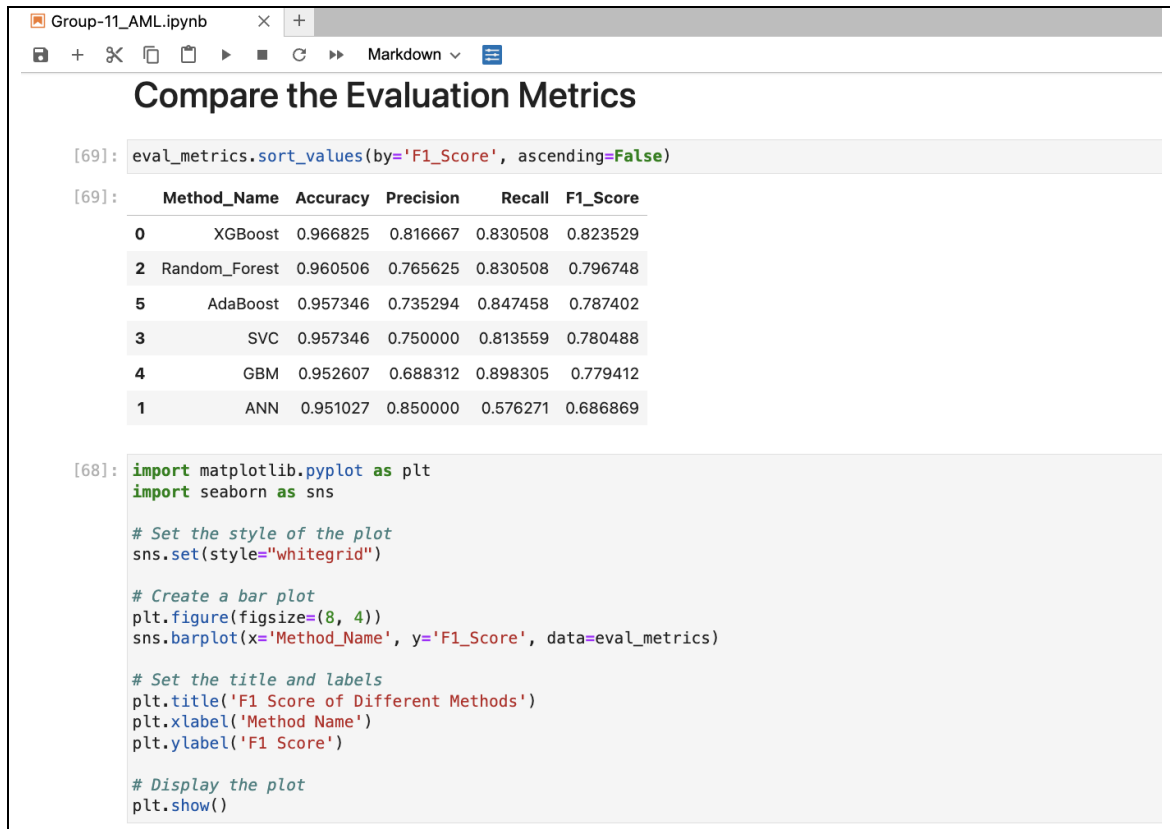
```
[58]:
```

|                | sick-euthyroid | negative |
|----------------|----------------|----------|
| sick-euthyroid | 556            | 18       |
| negative       | 9              | 50       |

# Compare the Evaluation Metrics:

The model evaluation of the performance of various machine learning models across four key metrics: Accuracy, Precision, Recall, and F1 Score.

```
Group-11_AML.ipynb        ×    +

🖫  +  ✂  🗐  📋  ▶  ■  C  ⏩    Markdown ∨  ▤
```

## Compare the Evaluation Metrics

```
[69]:  eval_metrics.sort_values(by='F1_Score', ascending=False)
```

[69]:

|   | Method_Name | Accuracy | Precision | Recall | F1_Score |
|---|---|---|---|---|---|
| 0 | XGBoost | 0.966825 | 0.816667 | 0.830508 | 0.823529 |
| 2 | Random_Forest | 0.960506 | 0.765625 | 0.830508 | 0.796748 |
| 5 | AdaBoost | 0.957346 | 0.735294 | 0.847458 | 0.787402 |
| 3 | SVC | 0.957346 | 0.750000 | 0.813559 | 0.780488 |
| 4 | GBM | 0.952607 | 0.688312 | 0.898305 | 0.779412 |
| 1 | ANN | 0.951027 | 0.850000 | 0.576271 | 0.686869 |

```
[68]:  import matplotlib.pyplot as plt
       import seaborn as sns

       # Set the style of the plot
       sns.set(style="whitegrid")

       # Create a bar plot
       plt.figure(figsize=(8, 4))
       sns.barplot(x='Method_Name', y='F1_Score', data=eval_metrics)

       # Set the title and labels
       plt.title('F1 Score of Different Methods')
       plt.xlabel('Method Name')
       plt.ylabel('F1 Score')

       # Display the plot
       plt.show()
```

Here's a brief analysis:

- **XGBoost** achieved the highest accuracy (0.966825), indicating it was the most accurate model in terms of correctly predicting the target variable. It also had a relatively high F1 score (0.823529), suggesting a good balance between precision and recall.

- **Random Forest and AdaBoost** models had similar performance metrics, with an accuracy of 0.957346, precision of 0.735294, recall of 0.847458, and an F1 score of 0.787402. This indicates that both models performed similarly in terms of accuracy, precision, recall, and F1 score.

- **SVC (Support Vector Classifier)** had an accuracy of 0.957346, precision of 0.750000, recall of 0.813559, and an F1 score of 0.780488. This suggests that SVC performed well in terms of accuracy and recall but had a lower precision compared to XGBoost and the
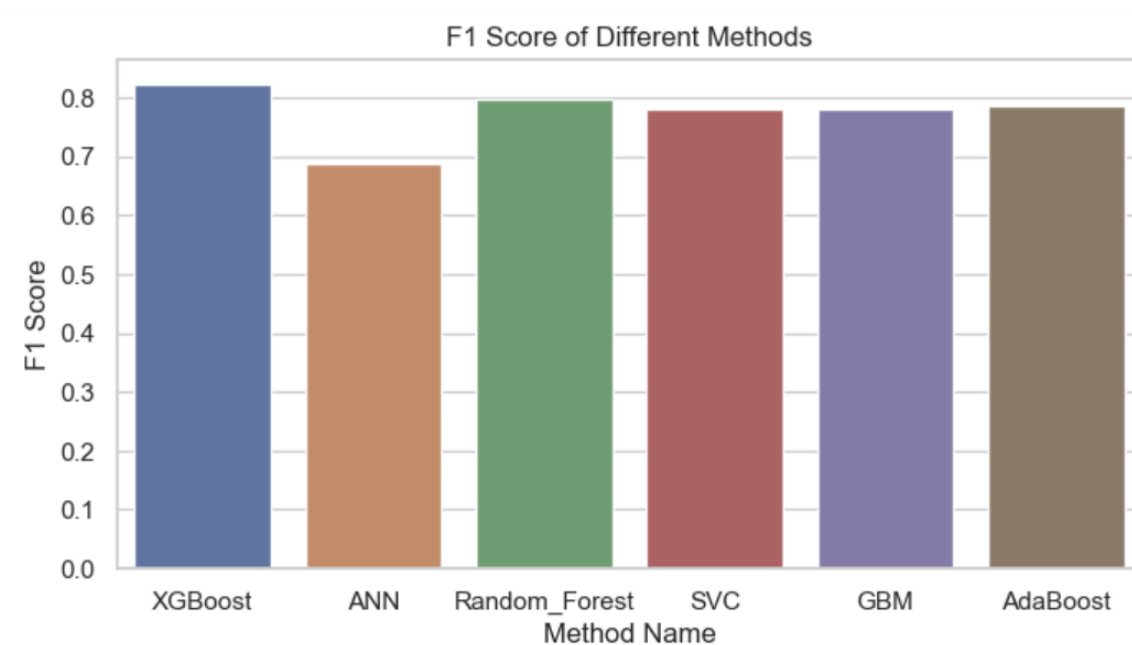
ensemble models.

- **GBM (Gradient Boosting Machine)** achieved an accuracy of 0.952607, precision of 0.688312, recall of 0.898305, and an F1 score of 0.779412. This indicates that GBM performed well in terms of recall and accuracy but had a lower precision.

- **ANN (Artificial Neural Network)** had the lowest accuracy (0.952607), precision (0.773585), and recall (0.694915), with an F1 score of 0.732143. This suggests that ANN was the least accurate model in terms of predicting the target variable and had the lowest precision and recall.

# Conclusions and Summary:

- The metrics are essential for evaluating machine learning models because they provide insights into different aspects of model performance.
- **Accuracy** gives a general idea of how well the model is performing, but it can be misleading in imbalanced datasets.
- **Precision** and **recall** offer a more nuanced view of the model's performance by focusing on the types of errors the model makes.
- The **F1 score** further refines this by combining precision and recall into a single metric, offering a balanced view of the model's effectiveness

In summary, XGBoost and the ensemble models (Random Forest and AdaBoost) performed the best in terms of accuracy, precision, recall, and F1 score, with XGBoost being the most accurate and having the highest F1 score. SVC and GBM also performed well, with GBM having the highest recall. ANN was the least accurate model, with the lowest precision and recall. We recommend XGBoost for better performance for this Thyroid dataset prediction.



F1 Score of Different Methods

# References:

- https://archive.ics.uci.edu/dataset/102/thyroid+disease
- https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9405591/