

---

# Data Analytics I| Assignment 4

---

**Author:** Prateek Pandey

**ID:** 2022201035

October 30, 2023



INTERNATIONAL INSTITUTE OF  
INFORMATION TECHNOLOGY

---

H Y D E R A B A D

## Contents

- There are very less samples in the above dataset , How do you deal with that?
  - The dataset I'm working with contains a limited number of samples. In order to preprocess this data, I transformed textual labels into a one-hot encoding format. This conversion resulted in a sparse matrix with numerous columns, reflecting the categorical values as binary features.
  - Additionally, I standardized the numeric columns to follow a standard normal distribution. This scaling process ensures that all numerical features contribute uniformly to the analysis, preparing them for further processing.
  - Due to the high dimensionality caused by the sparsity from one-hot encoding, I employed Principal Component Analysis (PCA). PCA allowed me to address both the sparsity and the high number of dimensions. By reducing the dimensionality, I was able to obtain a denser representation of the data, condensing it down to 7 columns, while still retaining crucial information.
  - Thus only relevant data remained
  - While Evaluating Model we used Stratified K-Fold Cross Validation, It maximizes the use of the available data by iteratively using it for both training and validation, reducing the variance in model performance evaluation and providing a more reliable estimate of how the model will perform on unseen data.

```
[26] from sklearn.model_selection import StratifiedKFold
X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.3, random_state=17)
skf = StratifiedKFold(n_splits=10, shuffle=True, random_state=17)
lst_accu_stratified = []

for train_index, test_index in skf.split(X, Y):
    x_train_fold, x_test_fold = X[train_index], X[test_index]
    y_train_fold, y_test_fold = Y[train_index], Y[test_index]
    ovr=OneVsOne(x_train_fold,y_train_fold)
    ovr.fit()
    y_pred=ovr.predict(x_test_fold)
    lst_accu_stratified.append(accuracy_score(y_pred, y_test_fold))
```

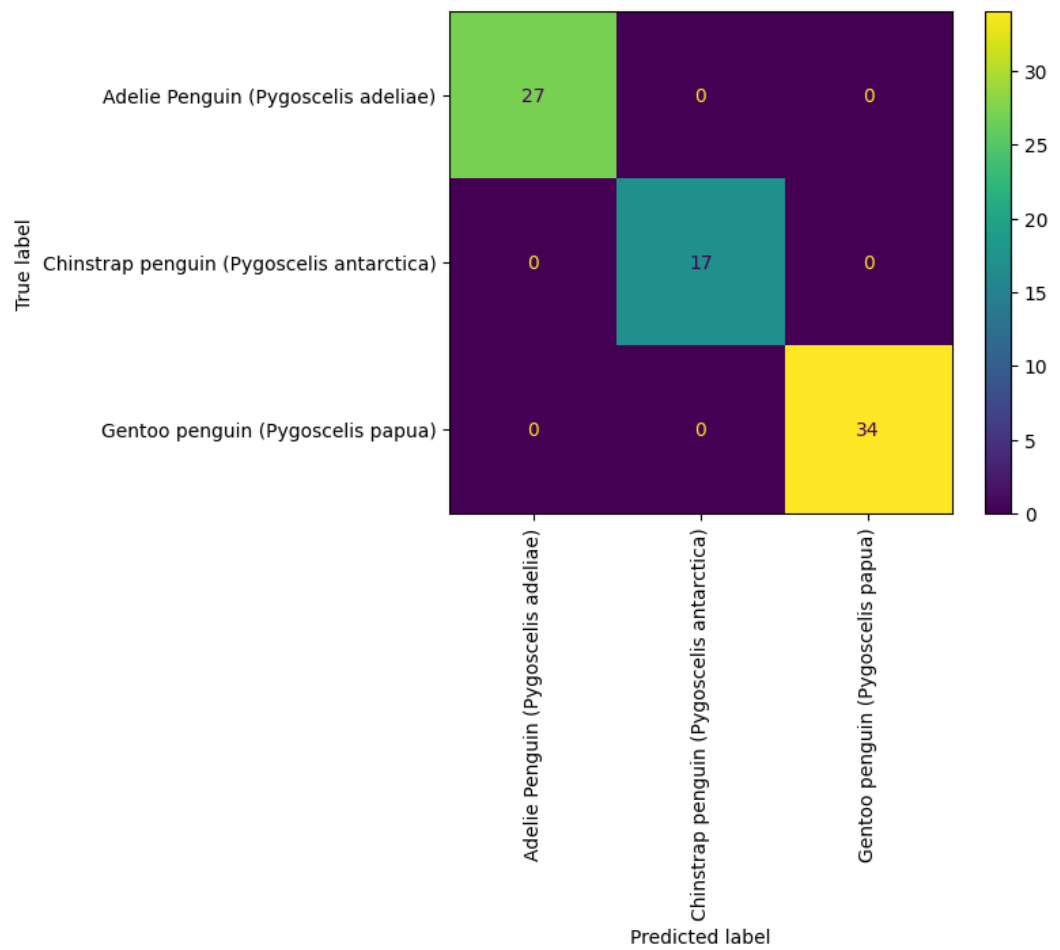
```
[26] from sklearn.model_selection import StratifiedKFold
X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.3, random_state=17)
skf = StratifiedKFold(n_splits=10, shuffle=True, random_state=17)
lst_accu_stratified = []

for train_index, test_index in skf.split(X, Y):
    x_train_fold, x_test_fold = X[train_index], X[test_index]
    y_train_fold, y_test_fold = Y[train_index], Y[test_index]
    ovr=OneVsOne(x_train_fold,y_train_fold)
    ovr.fit()
    y_pred=ovr.predict(x_test_fold)
    lst_accu_stratified.append(accuracy_score(y_pred, y_test_fold))
```

- Plot the confusion matrix, and include the precision, recall, f1-score metrics in the report .
  - One Vs All
    - Precision Recall Accuracy

	precision	recall	f1-score	support
0	1.00	1.00	1.00	27
1	1.00	1.00	1.00	17
2	1.00	1.00	1.00	34
accuracy			1.00	78
macro avg	1.00	1.00	1.00	78
weighted avg	1.00	1.00	1.00	78

## ■ Confusion Matrix

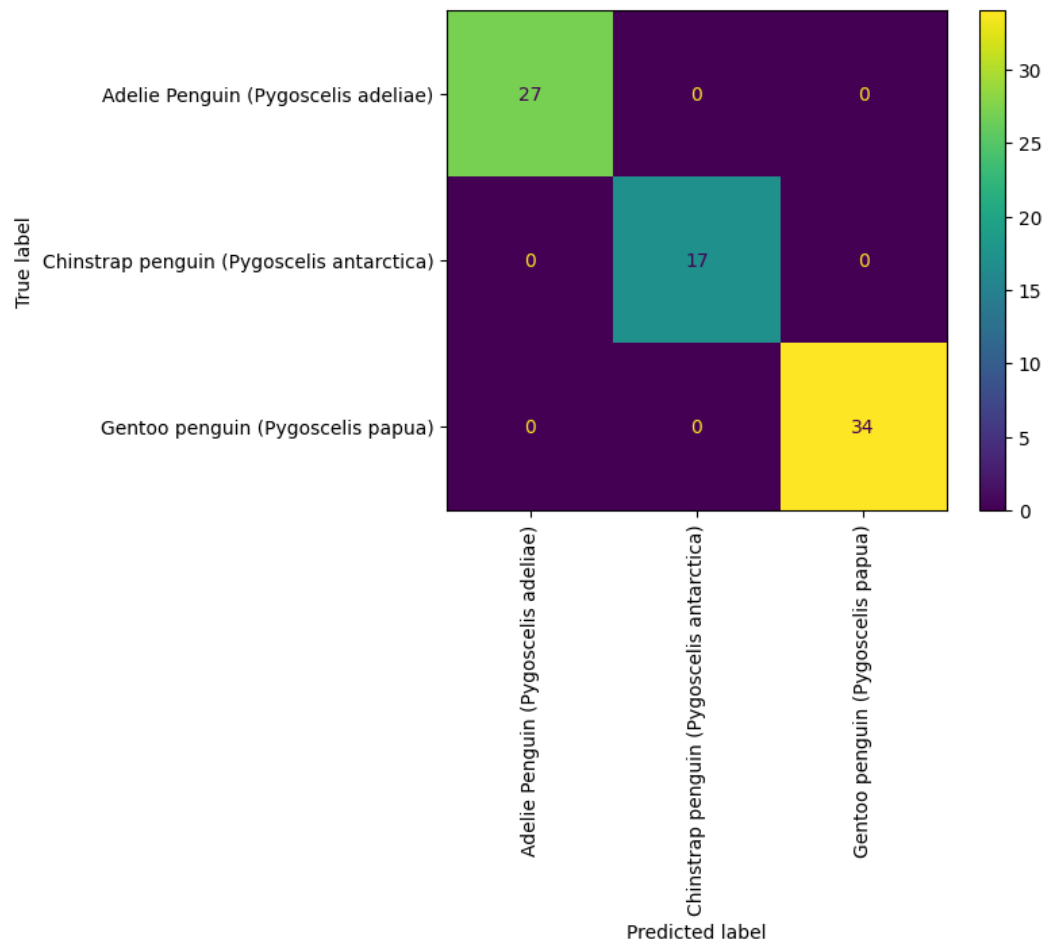


## ○ One vs One

### ■ Precision Recall Accuracy

	precision	recall	f1-score	support
0	1.00	1.00	1.00	27
1	1.00	1.00	1.00	17
2	1.00	1.00	1.00	34
accuracy			1.00	78
macro avg	1.00	1.00	1.00	78
weighted avg	1.00	1.00	1.00	78

## ■ Confusion Matrix



- Compare the results obtained for one-vs-one and one-vs-all(which according to you performs better for the above dataset)
  - One Vs One



- One Vs Rest



- Since both the models are showing 1.0 accuracy at a split of 0.7-0.3 both the models are accurate

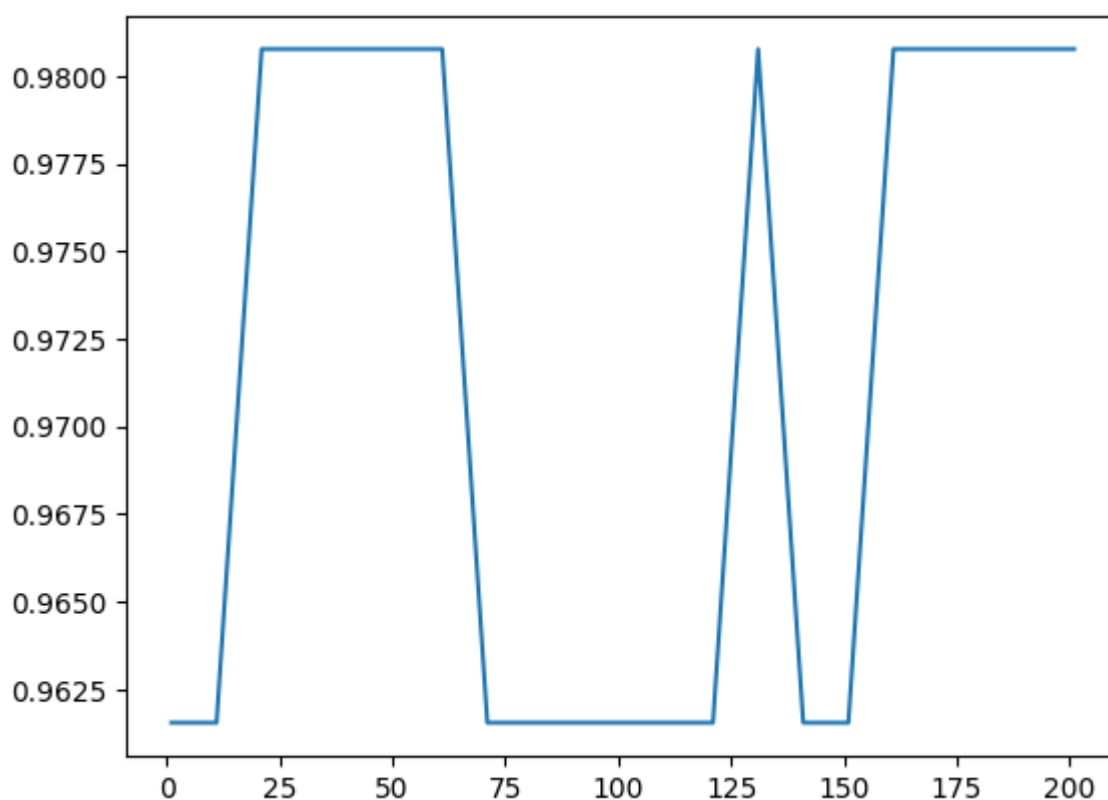
## Random Forest

- Predicted vs Actual Labels



	precision	recall	f1-score	support
Adelie Penguin ( <i>Pygoscelis adeliae</i> )	1.00	0.89	0.94	18
Chinstrap penguin ( <i>Pygoscelis antarctica</i> )	0.82	1.00	0.90	9
Gentoo penguin ( <i>Pygoscelis papua</i> )	1.00	1.00	1.00	25
accuracy			0.96	52
macro avg	0.94	0.96	0.95	52
weighted avg	0.97	0.96	0.96	52

- Accuracy Score over no. of trees 1-201 at interval of 10



## Confusion Matrix

