

Business Problem

Our goal is to create a model that can help predict a species of a penguin based on physical attributes, then we can use that model to help researchers classify penguins in the field, instead of needing an experienced biologist

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [2]: df = pd.read_csv("penguins_size.csv")
df.head()
```

Out[2]:

	species	island	culmen_length_mm	culmen_depth_mm	flipper_length_mm	body_mass_g
0	Adelie	Torgersen	39.1	18.7	181.0	3750.0
1	Adelie	Torgersen	39.5	17.4	186.0	3800.0
2	Adelie	Torgersen	40.3	18.0	195.0	3250.0
3	Adelie	Torgersen	NaN	NaN	NaN	NaN
4	Adelie	Torgersen	36.7	19.3	193.0	3450.0

"Palmer Penguins" dataset Summary:

- species: penguin species (Chinstrap, Adélie, or Gentoo)
 - culmen_length_mm: culmen length (mm)
 - culmen_depth_mm: culmen depth (mm)
 - flipper_length_mm: flipper length (mm)
 - body_mass_g: body mass (g)
 - island: island name (Dream, Torgersen, or Biscoe) in the Palmer Archipelago (Antarctica)
 - sex: penguin sex

```
In [3]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 344 entries, 0 to 343
Data columns (total 7 columns):
#   Column                Non-Null Count  Dtype
---  -
0   species               344 non-null   object
1   island                344 non-null   object
2   culmen_length_mm      342 non-null   float64
3   culmen_depth_mm       342 non-null   float64
4   flipper_length_mm     342 non-null   float64
5   body_mass_g           342 non-null   float64
6   sex                   334 non-null   object
dtypes: float64(4), object(3)
memory usage: 18.9+ KB
```

```
In [4]: df.isnull().sum()
```

```
Out[4]: species      0
island              0
culmen_length_mm    2
culmen_depth_mm     2
flipper_length_mm   2
body_mass_g         2
sex                 10
dtype: int64
```

Data Preprocessing

EDA

```
In [ ]:
```

```
In [ ]:
```

Feature Engineering

```
In [5]: #Dropping the missing values
df = df.dropna()
```

```
In [6]: #shape of data after dropping missing values
df.shape
```

```
Out[6]: (334, 7)
```

```
In [7]: df = df[df['sex']!='.']  
df.shape
```

```
Out[7]: (333, 7)
```

X & y

```
In [8]: X = pd.get_dummies(df.drop('species',axis=1),drop_first=True)  
y = df['species']
```

Train Test Split

```
In [9]: from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=101)
```

Modelling

Random Forest Classifier with default parameters

```
In [10]: # Import Classifier  
from sklearn.ensemble import RandomForestClassifier
```

```
In [11]: # Model with default parameters  
model = RandomForestClassifier()
```

```
In [12]: # Fit the model  
model.fit(X_train,y_train)
```

```
Out[12]: RandomForestClassifier()
```

```
In [13]: # Predict on X_train  
ypred_train = model.predict(X_train)
```

```
In [14]: # Predict on X_test  
ypred_test = model.predict(X_test)
```

Evaluation

```
In [15]: from sklearn.metrics import accuracy_score
print("Train accuracy:",accuracy_score(ypred_train,y_train))
print("Test accuracy:",accuracy_score(ypred_test,y_test))
```

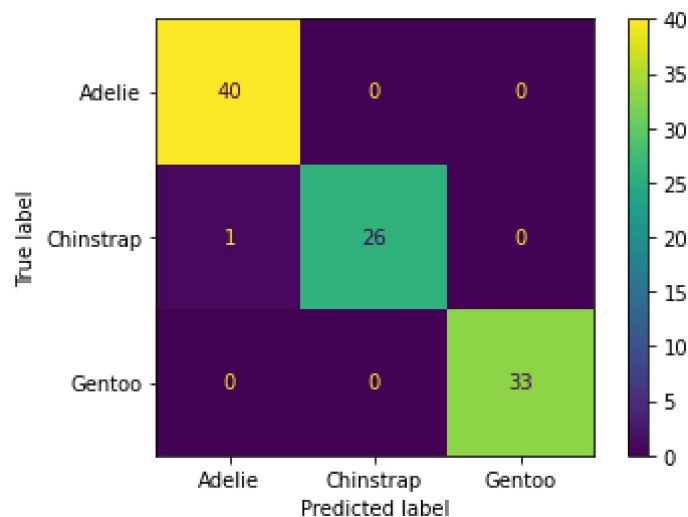
Train accuracy: 1.0
Test accuracy: 0.99

```
In [16]: # Cross validation
from sklearn.model_selection import cross_val_score
scores = cross_val_score(model,X,y,cv=5)
print(scores)
scores.mean()
```

[1. 0.98507463 0.98507463 0.98484848 1.]

Out[16]: 0.9909995477159657

```
In [17]: from sklearn.metrics import plot_confusion_matrix
plot_confusion_matrix(model,X_test,y_test)
plt.show()
```



```
In [18]: from sklearn.metrics import classification_report
print(classification_report(y_test,ypred_test))
```

	precision	recall	f1-score	support
Adelie	0.98	1.00	0.99	40
Chinstrap	1.00	0.96	0.98	27
Gentoo	1.00	1.00	1.00	33
accuracy			0.99	100
macro avg	0.99	0.99	0.99	100
weighted avg	0.99	0.99	0.99	100

Feature Importance

Very useful attribute of the trained model!

```
In [19]: model.feature_importances_
```

```
Out[19]: array([0.33223947, 0.14151084, 0.26307415, 0.11810937, 0.11197618,  
                0.02524413, 0.00784585])
```

HyperParameter Tuning

```
In [20]: from sklearn.model_selection import GridSearchCV
```

```
In [21]: # model  
estimator = RandomForestClassifier()  
  
# parameters (which you want to tune and identify the best)  
param_grid = {'n_estimators':list(range(1,101))}
```

```
In [22]: grid = GridSearchCV(estimator,param_grid, scoring="accuracy",cv=5)
```

```
In [23]: grid.fit(X_train,y_train)
```

```
Out[23]: GridSearchCV(cv=5, estimator=RandomForestClassifier(),  
                    param_grid={'n_estimators': [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11,  
12,  
                    13, 14, 15, 16, 17, 18, 19, 20, 21,  
                    22, 23, 24, 25, 26, 27, 28, 29, 30,  
                    ...]}},  
                    scoring='accuracy')
```

```
In [24]: grid.best_params_
```

```
Out[24]: {'n_estimators': 14}
```

Final Model with best parameters

```
In [25]: # Modelling  
final_model = RandomForestClassifier(n_estimators=14)  
final_model.fit(X_train,y_train)
```

```
Out[25]: RandomForestClassifier(n_estimators=8)
```

```
In [26]: # prediction  
predictions = final_model.predict(X_test)
```

```
In [27]: accuracy_score(predictions,y_test)
```

```
Out[27]: 0.99
```

```
In [28]: plot_confusion_matrix(final_model,X_test,y_test)  
plt.show()
```

