

# 1: firstly, we import what we need, and observe the data

```
In [1]: import numpy as np
import pandas as pd
import xgboost as xgb
from xgboost import XGBRegressor as XGBR
from sklearn.ensemble import RandomForestRegressor as RFR
from sklearn.linear_model import LinearRegression as LinearR
from sklearn.datasets import load_boston
from sklearn.model_selection import KFold, cross_val_score as CVS, train_test_split
from sklearn.metrics import mean_squared_error as MSE
import matplotlib.pyplot as plt
from time import time
import datetime
```

```
In [2]: data = pd.read_csv("/Users/guangxinsu/Desktop/penguins_raw.csv")
```

```
In [3]: data.head()
```

Out[3]:

	studyName	Sample Number	Species	Region	Island	Stage	Individual ID	Clutch Completion	Date Egg	Culmen Length (mm)
0	PAL0708	1	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A1	Yes	2007- 11-11	39
1	PAL0708	2	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A2	Yes	2007- 11-11	39
2	PAL0708	3	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A1	Yes	2007- 11-16	40
			Adelie Penguin			Adult,			2007	

## 2: Data preprocess : drop useless features, Turn letters, genders and other words into numbers, Fill the nan values with proper values like mean number.

```
1. data.drop(["Comments", "Stage", "studyName", "Sample Number", "Region", "Individual ID",  
    "Date Egg"], inplace=True, axis=1)  
2. data["Culmen Length (mm)"] = data["Culmen Length (mm)"].fillna(data["Culmen Length  
    (mm)"].mean())  
3. data["Culmen Depth (mm)"] = data["Culmen Depth (mm)"].fillna(data["Culmen Depth (mm)  
    (mm)"].mean())  
4. data["Flipper Length (mm)"] = data["Flipper Length (mm)"].fillna(data["Flipper Leng  
    th (mm)"].mean())  
5. data["Body Mass (g)"] = data["Body Mass (g)"].fillna(data["Body Mass (g)"].mean())  
6. data["Delta 15 N (o/oo)"] = data["Delta 15 N (o/oo)"].fillna(data["Delta 15 N (o/oo  
    )"].mean())  
7. data["Delta 13 C (o/oo)"] = data["Body Mass (g)"].fillna(data["Body Mass (g)"].mean  
    ())  
8. data.loc[:, "Sex"] = (data["Sex"] == "male").astype("int")  
9. data.loc[:, "Clutch Completion"] = (data["Clutch Completion"] == "Yes").astype("int")  
10. #Check how many values are turned into lists  
11. labels = data["Island"].unique().tolist()  
12. #Convert letters to number  
13. data["Island"] = data["Island"].apply(lambda x: labels.index(x))  
14. labels = data["Island"].unique().tolist()  
15. #Check how many values are turned into lists  
16. labels = data["Species"].unique().tolist()  
17. print(labels)  
18. #Convert letters to number  
19. data["Species"] = data["Species"].apply(lambda x: labels.index(x))  
20. labels = data["Species"].unique().tolist()  
21. data.head(300)
```

```
data.head(300)
```

```
['Adelie Penguin (Pygoscelis adeliae)', 'Gentoo penguin (Pygoscelis papua)', 'Chinstrap penguin (Pygoscelis antarctica)']
```

```
:
```

species	Island	Clutch Completion	Culmen Length (mm)	Culmen Depth (mm)	Flipper Length (mm)	Body Mass (g)	Sex	Delta 15 N (o/oo)	Delta 13 C (o/oo)
0	0	1	39.10000	18.70000	181.000000	3750.000000	0	8.733382	3750.000000
0	0	1	39.50000	17.40000	186.000000	3800.000000	0	8.949560	3800.000000
0	0	1	40.30000	18.00000	195.000000	3250.000000	0	8.368210	3250.000000
0	0	1	43.92193	17.15117	200.915205	4201.754386	0	8.733382	4201.754386
0	0	1	36.70000	19.30000	193.000000	3450.000000	0	8.766510	3450.000000
...	...	...	...	...	...	...	...	...	...
2	2	1	49.20000	18.20000	195.000000	4400.000000	0	9.271580	4400.000000
2	2	1	42.40000	17.30000	181.000000	3600.000000	0	9.351380	3600.000000
2	2	1	48.50000	17.50000	191.000000	3400.000000	0	9.426660	3400.000000
2	2	0	43.20000	16.60000	187.000000	2900.000000	0	9.354160	2900.000000
2	2	0	50.60000	19.40000	193.000000	3800.000000	0	9.281530	3800.000000

vs x 10 columns

```
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 344 entries, 0 to 343
```

```
Data columns (total 10 columns):
```

#	Column	Non-Null Count	Dtype
0	Species	344 non-null	int64
1	Island	344 non-null	int64
2	Clutch Completion	344 non-null	int64
3	Culmen Length (mm)	344 non-null	float64
4	Culmen Depth (mm)	344 non-null	float64
5	Flipper Length (mm)	344 non-null	float64
6	Body Mass (g)	344 non-null	float64
7	Sex	344 non-null	int64
8	Delta 15 N (o/oo)	344 non-null	float64
9	Delta 13 C (o/oo)	344 non-null	float64

```
dtypes: float64(6), int64(4)
```

```
memory usage: 27.0 KB
```

### 3: divide train set and test set

```
In [15]: #data and target
x = data.iloc[:,data.columns != "Species"]
y = data.iloc[:,data.columns == "Species"]
```

```
In [16]: #train and test set
Xtrain,Xtest,Ytrain,Ytest = TTS(x,y,test_size=0.3,random_state=420)
```

## 4: Xgboost Preliminary training & cross-validation & random forest & linear-regression

### Compare with each other

```
In [17]: reg = XGBR(n_estimators=100).fit(Xtrain,Ytrain)
         reg.predict(Xtest)
         reg.score(Xtest,Ytest)
```

```
Out[17]: 0.9175140971181192
```

```
In [19]: #One of the advantages of the tree model: being able to view the importance
         #the feature importances
         reg.feature_importances_
```

```
Out[19]: array([4.5520934e-01, 1.1017058e-03, 4.0963203e-01, 9.0594076e-02,
                4.1496810e-02, 1.8036369e-03, 0.0000000e+00, 1.6241733e-04,
                0.0000000e+00], dtype=float32)
```

```
In [20]: #Cross-validation
         reg = XGBR(n_estimators=100)
```

```
In [21]: CVS(reg,Xtrain,Ytrain,cv=5)
```

```
Out[21]: array([0.96334981, 0.93413385, 0.964922 , 0.94453331, 0.93087006])
```

```
In [22]: CVS(reg,Xtrain,Ytrain,cv=5,scoring='neg_mean_squared_error').mean()
```

```
Out[22]: -0.03033896242461126
```

```
] : #Let's take a look at all the model evaluation indicators in sklearn
import sklearn
sorted(sklearn.metrics.SCORERS.keys())
#random forest
rfr = RFR(n_estimators=100)
CVS(rfr,Xtrain,Ytrain,cv=5).mean()
```

```
estimator.fit(X_train, y_train, **fit_params)
```

```
Out[23]: 0.929550818059646
```

```
[25]: #linear regression
      lr = LinearR()
      CVS(lr,Xtrain,Ytrain,cv=5).mean()
```

```
[25]: 0.7479632602108952
```

```
[26]: CVS(lr,Xtrain,Ytrain,cv=5,scoring='neg_mean_squared_error').mean()
```

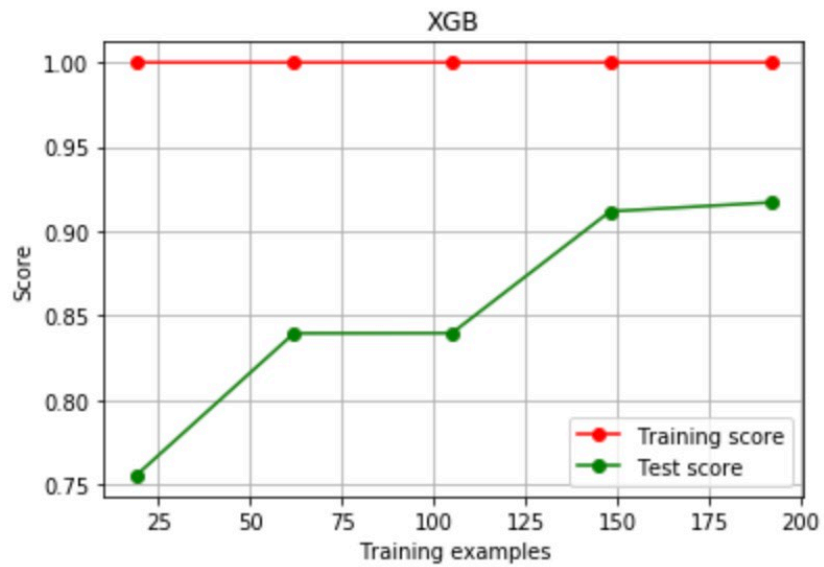
```
[26]: -0.14709472237947832
```

```
In [27]: reg = XGBR(n_estimators=10,silent=False)
         CVS(reg,Xtrain,Ytrain,cv=5,scoring='neg_mean_squared_error').mean()
```

## 5: plot learning\_curve figure of xgboost and analysis, it's overfitting now!

```
1. def plot_learning_curve(estimator, title, X, y,
2.                         ax=None, #Select subgraph
3.                         ylim=None, #Set the value range of the ordinate
4.                         cv=None, #cross validation
5.                         n_jobs=None #Set the thread to use
6.                         ):
7.     from sklearn.model_selection import learning_curve
8.     import matplotlib.pyplot as plt
9.     import numpy as np
10.    train_sizes, train_scores, test_scores = learning_curve(estimator, X, y
11.                                                            ,shuffle=True
12.                                                            ,cv = cv
13.                                                            ,random_state = 420
14.                                                            ,n_jobs = n_jobs)
15.    if ax == None:
16.        ax = plt.gca()
17.    else:
18.        ax = plt.figure()
19.    ax.set_title(title)
20.    if ylim is not None:
21.
22.        ax.set_ylim(*ylim)
23.    ax.set_xlabel("Training examples")
24.    ax.set_ylabel("Score")
25.    ax.grid() #grid
26.    ax.plot(train_sizes, np.mean(train_scores, axis=1), 'o-'
27.            , color="r",label="Training score")
28.    ax.plot(train_sizes, np.mean(test_scores, axis=1), 'o-'
29.            , color="g",label="Test score")
30.    ax.legend(loc="best")
31.    return ax
```

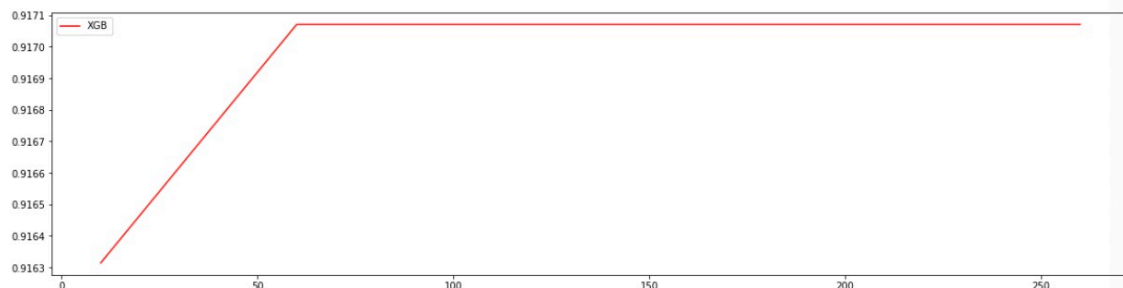
```
29]: cv = KFold(n_splits=5, shuffle = True, random_state=42)
plot_learning_curve(XGBR(n_estimators=500,random_state=420)
                    , "XGB",Xtrain,Ytrain,ax=None,cv=cv)
plt.show()
```



## 6: find the best n\_estimators/the number of trees

```
1. #Observe the influence of n_estimators on the model by using a parameter learning curve
2. axisx = range(10,310,50)
3. rs = []
4. for i in axisx:
5.     reg = XGBR(n_estimators=i,random_state=420)
6.     rs.append(CVS(reg,Xtrain,Ytrain,cv=cv).mean())
7. print(axisx[rs.index(max(rs))],max(rs))
8. plt.figure(figsize=(20,5))
9. plt.plot(axisx,rs,c="red",label="XGB")
10. plt.legend()
11. plt.show()
```

110 0.917070470188386



```
1. #Verify that the model effect has improved?
2. time0 = time()
3. print(XGBR(n_estimators=1,random_state=420).fit(Xtrain,Ytrain).score(Xtest,Ytest))
4. print(time()-time0)
5. time0 = time()
6. print(XGBR(n_estimators=10,random_state=420).fit(Xtrain,Ytrain).score(Xtest,Ytest))
7. print(time()-time0)
8.
9. time0 = time()
10. print(XGBR(n_estimators=40,random_state=420).fit(Xtrain,Ytrain).score(Xtest,Ytest))
11. print(time()-time0)
12. time0 = time()
13. print(XGBR(n_estimators=110,random_state=420).fit(Xtrain,Ytrain).score(Xtest,Ytest)
14. )
14. print(time()-time0)
15.
16. time0 = time()
17. print(XGBR(n_estimators=300,random_state=420).fit(Xtrain,Ytrain).score(Xtest,Ytest)
18. )
18. print(time()-time0)
```

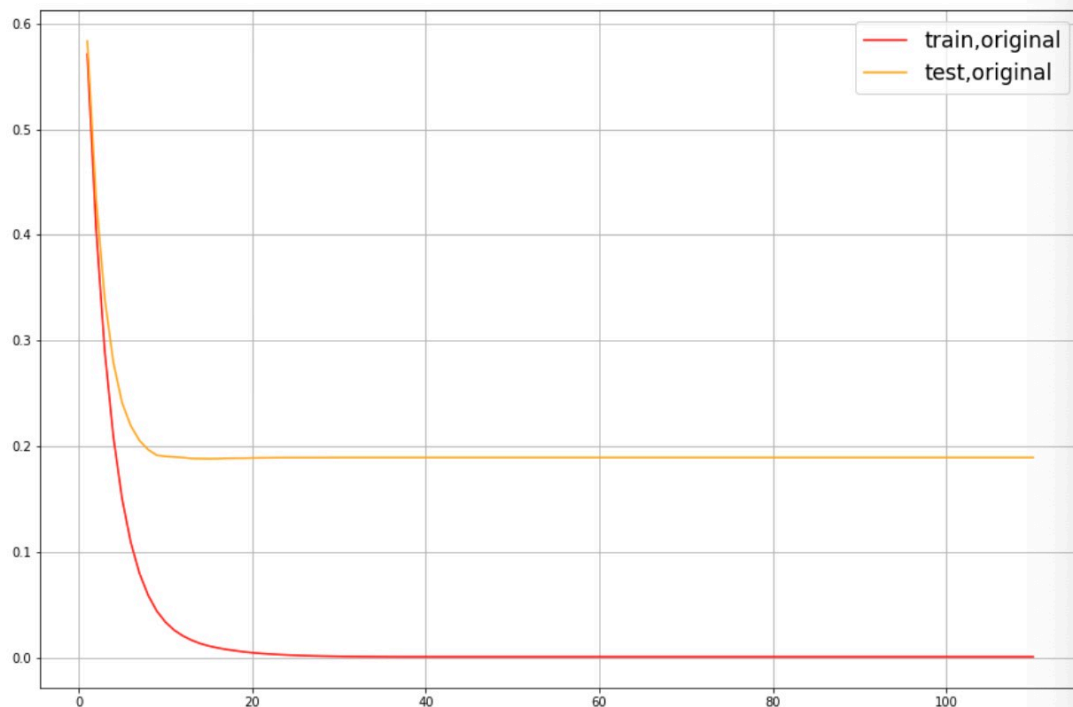


0.4343177778791587  
0.00750422477722168  
0.9160991184463382  
0.018697023391723633  
0.9175139811105371  
0.01108407974243164  
0.9175140971181192  
0.013628005981445312  
0.9175140971181192  
0.022681236267089844

## 7: use xgboost and solve the over fitting problems(number of trees :110)

```
1. dfull = xgb.DMatrix(x,y)
2. param1 = {'silent':True
3.           , 'obj':'reg:linear'
4.           , "subsample":1 , "max_depth":6 , "eta":0.3
5.           , "gamma":0
6.           , "lambda":1
7.           , "alpha":0 , "colsample_bytree":1 , "colsample_bylevel":1 , "colsample_bynode":1 , "nfold":5}
8. num_round = 110
9. time0 = time()
10. cvresult1 = xgb.cv(param1, dfull, num_round)
11. print(datetime.datetime.fromtimestamp(time()-time0).strftime("%M:%S:%F"))
12. fig,ax = plt.subplots(1,figsize=(15,10))
13. #ax.set_ylim(top=5)
14. ax.grid()
15. ax.plot(range(1,111),cvresult1.iloc[:,0],c="red",label="train,original")
16. ax.plot(range(1,111),cvresult1.iloc[:,2],c="orange",label="test,original")
17. ax.legend(fontsize="xx-large")
18. plt.show()
```

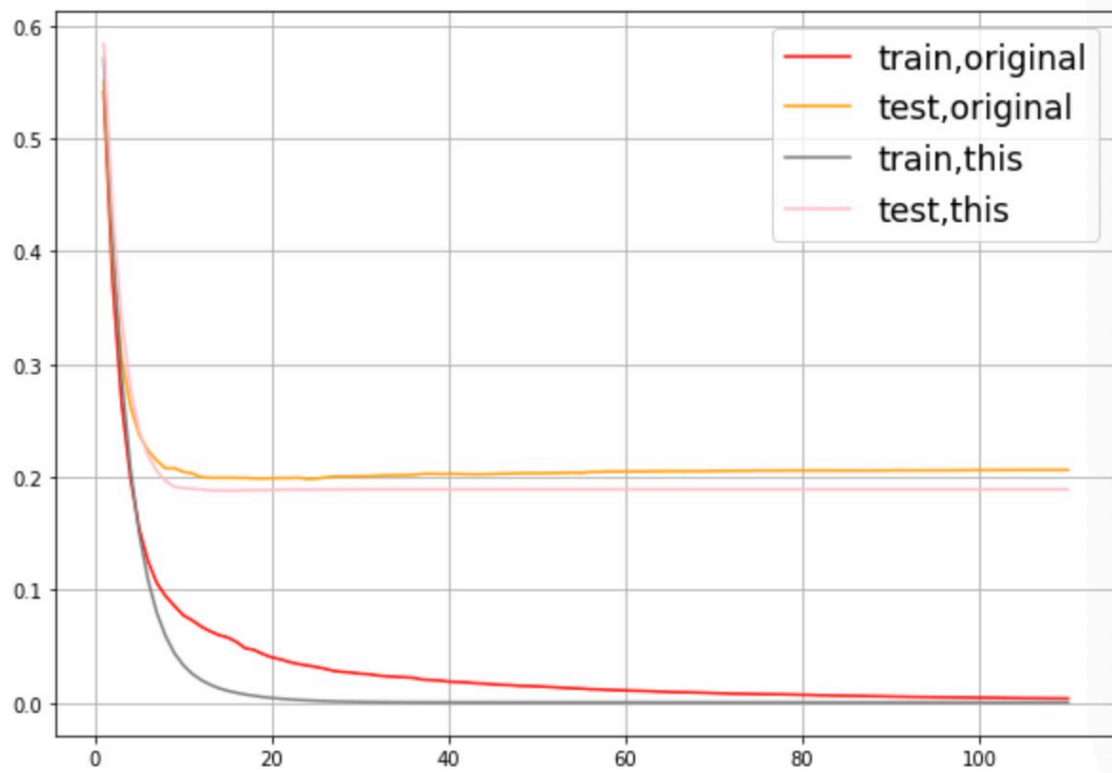
00:00:082927



## Change the parameters and solve the over-fitting(many times)

```
1. param1 = {'silent':True
2.           , 'obj':'reg:linear'
3.           , "subsample":1 , "max_depth":3 , "eta":0.35
4.           , "gamma":0
5.           , "lambda":1
6.           , "alpha":0 , "colsample_bytree":1 , "colsample_bylevel":1 , "colsample_bynode":1 , "nfold":5}
7. num_round = 110
8. time0 = time()
9. cvresult1 = xgb.cv(param1, dfull, num_round)
10. print(datetime.datetime.fromtimestamp(time()-time0).strftime("%M:%S:%f"))
11. fig, ax = plt.subplots(1, figsize=(15,10))
12. #ax.set_ylim(top=5)
13. ax.grid()
14. ax.plot(range(1,111),cvresult1.iloc[:,0],c="red",label="train,original")
15. ax.plot(range(1,111),cvresult1.iloc[:,2],c="orange",label="test,original")
16. param2 = {'silent':True
17.           , 'obj':'reg:linear'
18.           , "nfold":5}
19. param3 = {'silent':True
20.           , 'obj':'reg:linear'
21.           , "nfold":5}
22. time0 = time()
23. cvresult2 = xgb.cv(param2, dfull, num_round)
24. print(datetime.datetime.fromtimestamp(time()-time0).strftime("%M:%S:%f"))
25. time0 = time()
26. cvresult3 = xgb.cv(param3, dfull, num_round)
27. print(datetime.datetime.fromtimestamp(time()-time0).strftime("%M:%S:%f"))
28. ax.plot(range(1,111),cvresult3.iloc[:,0],c="gray",label="train,this")
29. ax.plot(range(1,111),cvresult3.iloc[:,2],c="pink",label="test,this")
30. ax.legend(fontsize="xx-large")
31. plt.show()
```

00:00:095722



8: we get satisfying parameters (tree size:110)

```
param1 = {'silent':True
          , 'obj':'reg:linear'
          , "subsample":1
          , "max_depth":3
          , "eta":0.35
          , "gamma":0
          , "lambda":1
          , "alpha":0
          , "colsample_bytree":1
          , "colsample_bylevel":1
          , "colsample_bynode":1
          , "nfold":5}
num_round = 110
```

## 9: save the data as a file and train the model with parameters we processed

```
[31]: ['/Users/guangxinsu/Desktop/ml',
       '/opt/anaconda3/lib/python37.zip',
       '/opt/anaconda3/lib/python3.7',
       '/opt/anaconda3/lib/python3.7/lib-dynload',
       '',
       '/opt/anaconda3/lib/python3.7/site-packages',
       '/opt/anaconda3/lib/python3.7/site-packages/aeosa',
       '/opt/anaconda3/lib/python3.7/site-packages/IPython/extensions',
       '/Users/guangxinsu/.ipython']
```

copyy

```
import pickle
dtrain = xgb.DMatrix(Xtrain,Ytrain)
param = {'silent':True
         , 'obj': 'reg:linear'
         , "subsample":1
         , "max_depth":3
         , "eta":0.35
         , "gamma":0
         , "lambda":1
         , "alpha":0
         , "colsample_bytree":1
         , "colsample_bylevel":1
         , "colsample_bynode":1
         , "nfold":5}
num_round = 110
bst = xgb.train(param,dtrain,num_round)
pickle.dump(bst,open("xgboostonboston.dat","wb"))
import sys
sys.path
```

```
[23:44:59] WARNING: /Users/travis/build/dmlc/xgboost/src/learner.cc:480:
Parameters: { nfold, obj, silent } might not be used.
```

This may not be accurate due to some parameters are only used in language bindings but passed down to XGBoost core. Or some parameters are not used but slip through this verification. Please open an issue if you find above cases.

## Copy the code above (from beginning) and train the model

1. from sklearn.model\_selection import train\_test\_split as TTS
2. from sklearn.metrics import mean\_squared\_error as MSE
3. import pickle
4. import xgboost as xgb
5. import pandas as pd

```

6. data = pd.read_csv("/Users/guangxinsu/Desktop/penguins_raw.csv")
7. data.drop(["Comments", "Stage", "studyName", "Sample Number", "Region", "Individual ID",
    "Date Egg"], inplace=True, axis=1)
8. data["Culmen Length (mm)"] = data["Culmen Length (mm)"].fillna(data["Culmen Length (mm)"].mean())
9. data["Culmen Depth (mm)"] = data["Culmen Depth (mm)"].fillna(data["Culmen Depth (mm)"].mean())
10. data["Flipper Length (mm)"] = data["Flipper Length (mm)"].fillna(data["Flipper Length (mm)"].mean())
11. data["Body Mass (g)"] = data["Body Mass (g)"].fillna(data["Body Mass (g)"].mean())
12. data["Delta 15 N (o/oo)"] = data["Delta 15 N (o/oo)"].fillna(data["Delta 15 N (o/oo)"].mean())
13. data["Delta 13 C (o/oo)"] = data["Body Mass (g)"].fillna(data["Body Mass (g)"].mean())
14. data.loc[:, "Sex"] = (data["Sex"] == "male").astype("int")
15. data.loc[:, "Clutch Completion"] = (data["Clutch Completion"] == "Yes").astype("int")
16. #Check how many values are turned into lists
17. labels = data["Species"].unique().tolist()
18. print(labels)
19. #Convert letters to number
20. data["Species"] = data["Species"].apply(lambda x: labels.index(x))
21. labels = data["Species"].unique().tolist()
22. data.head(300)
23. #Check how many values are turned into lists
24. labels = data["Island"].unique().tolist()
25. #Convert letters to number
26. data["Island"] = data["Island"].apply(lambda x: labels.index(x))
27. labels = data["Island"].unique().tolist()
28. data
29. x = data.iloc[:, data.columns != "Species"]
30. y = data.iloc[:, data.columns == "Species"]
31. Xtrain, Xtest, Ytrain, Ytest = TTS(x, y, test_size=0.3, random_state=420)
32. dtest = xgb.DMatrix(Xtest, Ytest)
33. loaded_model = pickle.load(open("xgboostonboston.dat", "rb"))
34. print("Loaded model from: xgboostonboston.dat")
35. ypreds = loaded_model.predict(dtest)
36. ypreds

```

```

In [9]: from sklearn.metrics import mean_squared_error as MSE, r2_score
        MSE(Ytest, ypreds)

```

```

Out[9]: 0.054973624965083075

```

```

In [10]: r2_score(Ytest, ypreds)

```

```

Out[10]: 0.9010163596433597

```

```

In [ ]:

```

Actually  $0.9010 < 0.9175$ , But the generalization performance of the model has improved!

Thanks 😊