

Project Objective

Build a simple regression model to predict the annual salary for each customer using the attributes you identified above

How accurate is your model? Do you need ANZ user to segment customers (for whom it does not have this data) into income brackets for reporting purposes?

For a challenge: build a decision-tree based model to predict salary. Does it perform better? How would you accurately test the performance of this model?

Import Libraries

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

import xgboost as xgb
from sklearn import KNeighborsClassifier, XGBRegressor
from xgboost import plot_importance, plot_shapforce

import matplotlib
sns.set_style('dark')
sns.set(font_scale=1.2)

from sklearn.linear_model import LinearRegression

from sklearn.inspection import permutation_importance
from sklearn.model_selection import train_test_split, GridSearchCV, RandomizedSearchCV
from sklearn.preprocessing import LabelEncoder, StandardScaler, MinMaxScaler, OneHotEncoder
from sklearn.pipeline import Pipeline
from sklearn.metrics import confusion_matrix, classification_report, mean_absolute_error, mean_squared_error, r2_score
from sklearn.metrics import plot_confusion_matrix, plot_precision_recall_curve, plot_roc_curve, accuracy_score
from sklearn.metrics import auc, f1_score, precision_score, recall_score, roc_auc_score

import warnings
warnings.filterwarnings('ignore')

import pickle
from pickle import dump, load

np.random.seed(0)

#from pycares.classification import *
#from pycares.clustering import *
from pycares.regression import *

pd.set_option('display.max_columns',100)
#pd.set_option('display.max_rows',100)
pd.set_option('display.width',100)
pd.set_option('print.prec',4)

Data Exploration and Analysis
```

```
In [2]: df = pd.read_csv('salary3.csv')
```

```
In [3]: df
```

```
Out[3]:
```

	account	first_name	age	gender	balance	months	weekspayment	annual
0	ACC-103705054	Rhonda	40	F	51472.20	12	52	46388.68
1	ACC-105663902	Michael	22	M	298308.49	12	24	166140.52
2	ACC-119951521	Billy	52	M	86898.05	12	28	196860.56
3	ACC-121703613	Kimberly	27	F	13769.63	8	8	252908.24
4	ACC-122330524	Michael	38	M	22826.60	12	52	52110.76
...
95	ACC-854938045	James	28	M	88169.32	12	28	132011.36
96	ACC-80814749	Christopher	35	M	65301.33	12	24	120050.84
97	ACC-586000567	Sandra	34	F	59807.25	12	28	182915.72
98	ACC-064838203	Michael	21	M	467645.22	12	52	81130.40
99	ACC-06140392	Joseph	21	M	49877.75	12	28	133791.32
100 rows x 8 columns								

```
In [4]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100 entries, 0 to 99
Data columns (total 8 columns):
 #   Column             Non-Null Count  Dtype
---  -
 0   account            100 non-null    object
 1   first_name         100 non-null    object
 2   age                100 non-null    int64
 3   gender             100 non-null    object
 4   balance            100 non-null    float64
 5   months            100 non-null    int64
 6   weekpayment        100 non-null    float64
 7   annual             100 non-null    float64
dtypes: float64(4), int64(3), object(1)
memory usage: 4.1 KB
```

```
In [5]: df.describe(include='all')
```

	account	first_name	age	gender	balance	months	weekpayment	annual
count	100	100	100	0	100000000.00	100	100000000.00	100000000.00
unique	100	80	NaN	2	NaN	NaN	NaN	NaN
top	ACC-354160573	Michael	NaN	M	NaN	NaN	NaN	NaN
freq	1	6	NaN	56	NaN	NaN	NaN	NaN
mean	NaN	NaN	NaN	31.770000	NaN	1.43244e+05	35.320000	119458.242800
std	NaN	NaN	NaN	15.642554	NaN	2.25789e+05	0.787788	1.460540
min	NaN	NaN	NaN	18.000000	NaN	1.37696e+04	8.000000	8.000000
25%	NaN	NaN	22.000000	NaN	4.965724e+04	12.000000	24.000000	59972.120000
50%	NaN	NaN	29.500000	NaN	7.227472e+04	12.000000	28.000000	101370.360000
75%	NaN	NaN	39.250000	NaN	1.150499e+05	12.000000	52.000000	150109.700000
max	NaN	NaN	78.000000	NaN	1.584768e+06	12.000000	56.000000	459470.860000

```
In [6]: df.shape
```

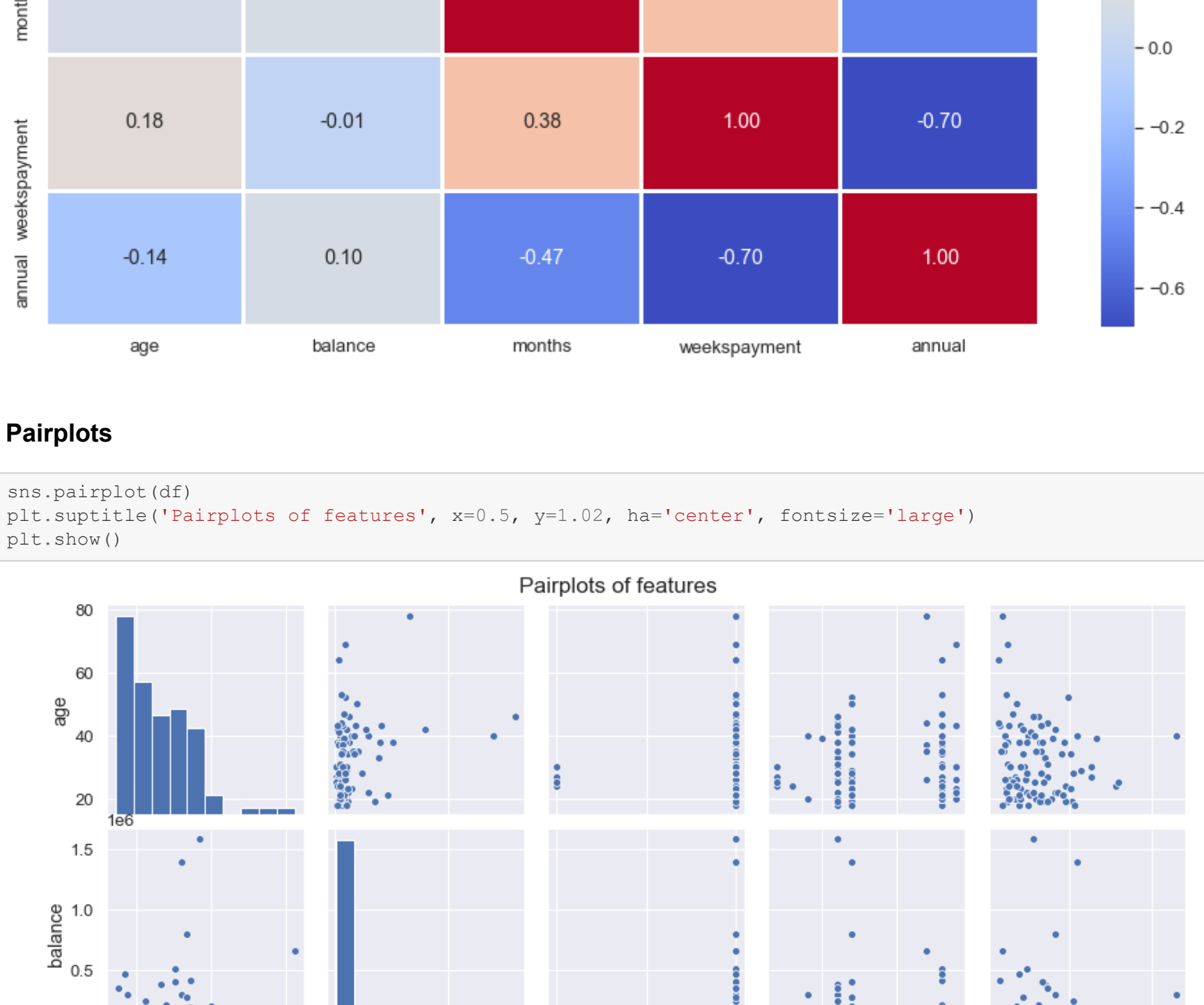
```
Out[6]: (100, 8)
```

```
In [7]: df.columns
```

```
Out[7]: Index(['account', 'first_name', 'age', 'gender', 'balance', 'months', 'weekpayment', 'annual'], dtype='object')
```

Data Visualization

Univariate Data Exploration



```
In [9]: df.boxplot(figsize=(20,10))
```

```
plt.suptitle('BoxPlot', x=0.5, y=1.02, ha='center', fontsize='large')
```

```
plt.tight_layout()
```

```
plt.show()
```

Correlation

```
In [10]: df.corr()
```

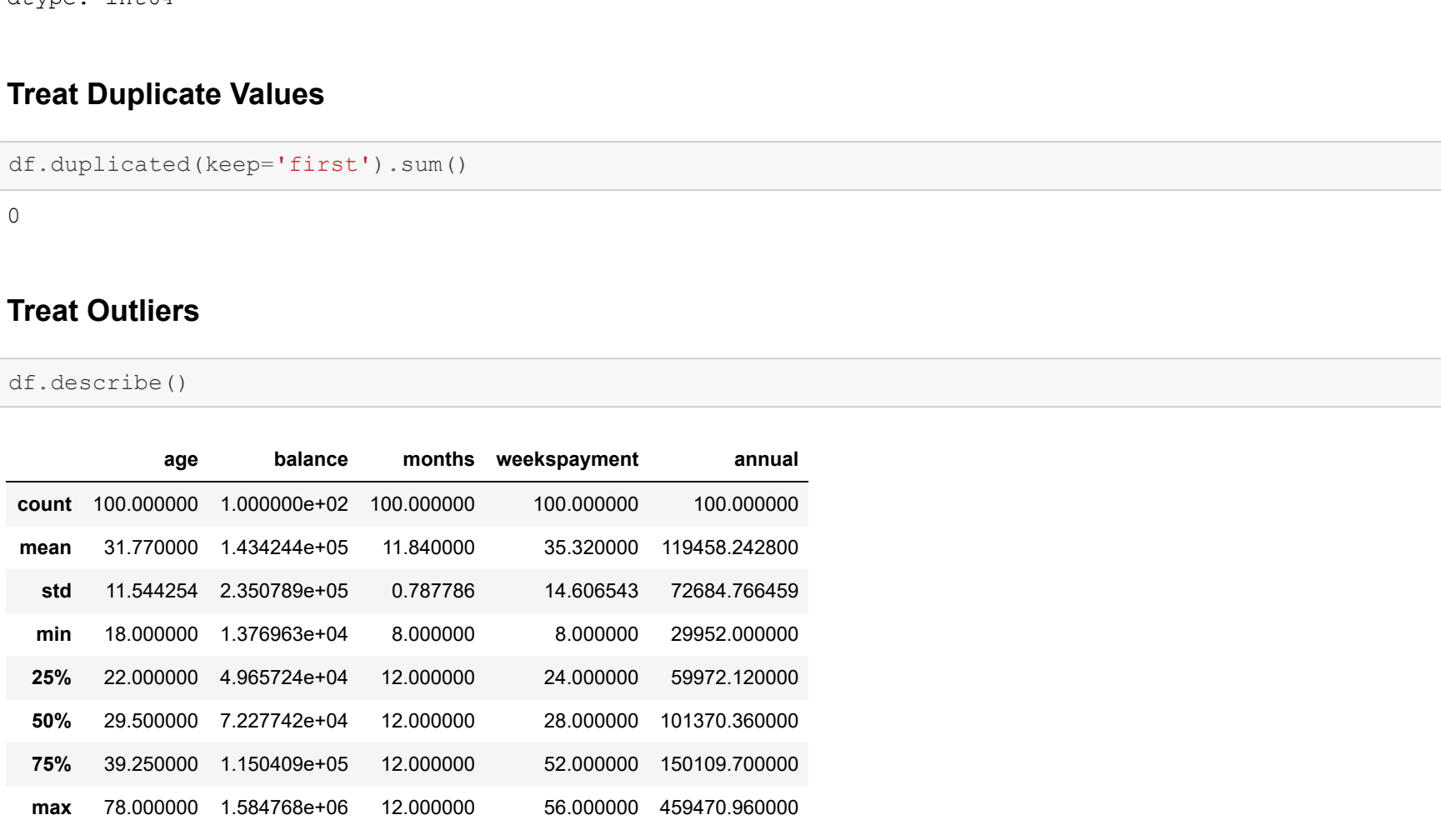
```
Out[10]:
```

	age	balance	months	weekpayment	annual
age	1.000000	0.288224	0.093653	0.182368	-0.135264
balance	0.288224	1.000000	0.107047	-0.013481	0.101847
months	0.093653	0.107047	1.000000	0.383716	-0.466623
weekpayment	0.182368	-0.013481	0.383716	1.000000	-0.696428
annual	-0.135264	0.101847	-0.466623	-0.696428	1.000000

```
In [11]: plt.figure(figsize=(16,5))
```

```
sns.heatmap(df.corr(), cmap='coolwarm', annot=True, fmt='.2f', linewidths=2)
```

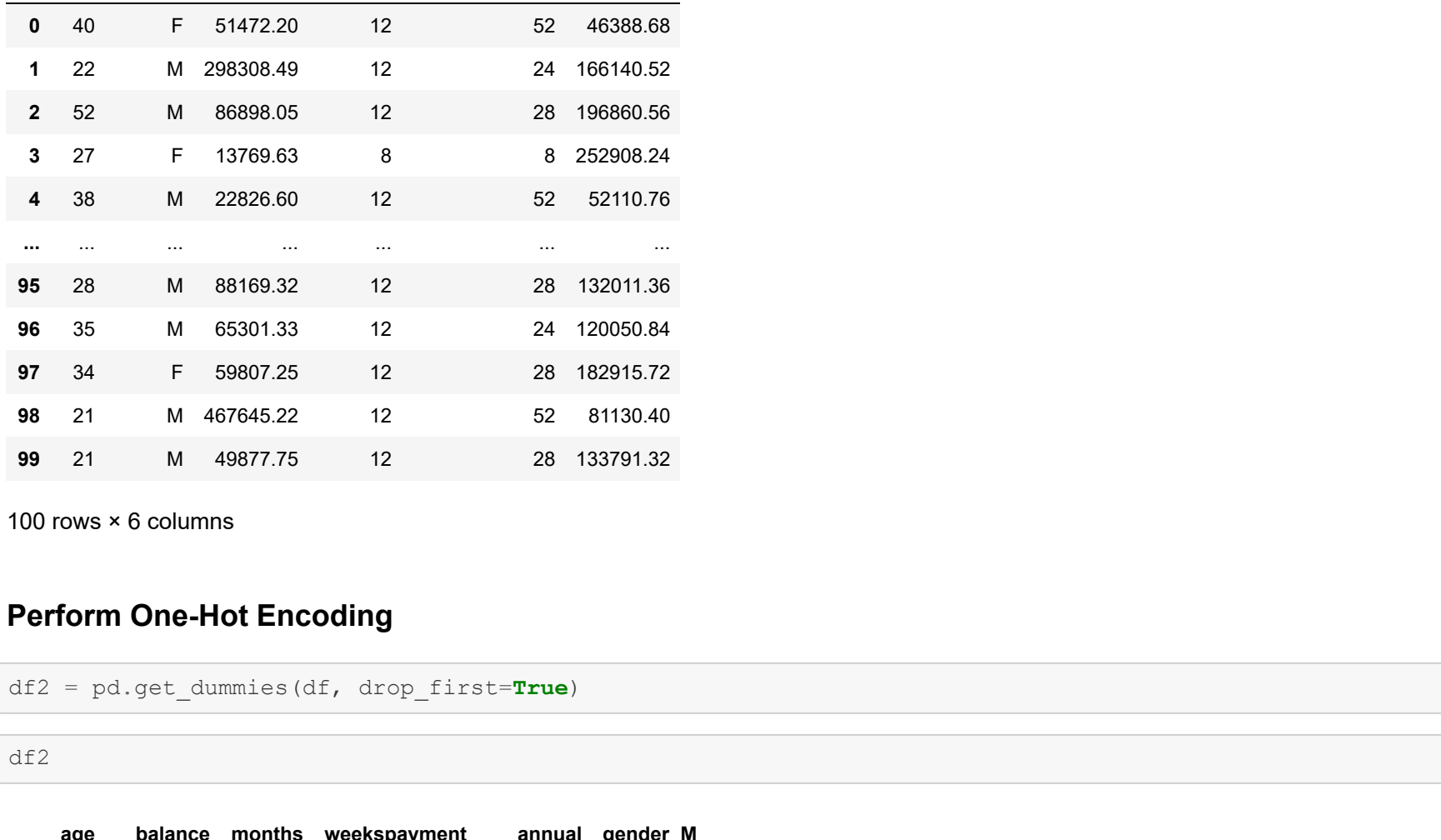
```
plt.show()
```



Pairplots

```
In [12]: sns.pairplot(df)
```

```
Out[12]:
```



Data Preprocessing

Treat Missing Values

```
In [13]: df.isnull().sum()
```

```
Out[13]:
```

account	0
first_name	0
age	0
gender	0
balance	0
months	0
weekpayment	0
annual	0
dtype:	int64(8)

Treat Duplicate Values

```
In [14]: df.duplicated(keep='first').sum()
```

```
Out[14]: 0
```

Treat Outliers

```
In [15]: df.describe()
```

```
Out[15]:
```

	age	balance	months	weekpayment	annual
count	100	100000000.00	100	100000000.00	100000000.00
mean	31.770000	1.43244e+05	11.840000	35.320000	119458.242800
std	15.642554	2.25789e+05	0.787788	1.460540	72864.766490
min	18.000000	1.37696e+04	8.000000	8.000000	8.000000
25%	22.000000	4.965724e+04	12.000000	24.000000	59972.120000
50%	29.500000	7.227472e+04	12.000000	28.000000	101370.360000
75%	39.250000	1.150499e+05	12.000000	52.000000	150109.700000
max	78.000000	1.584768e+06	12.000000	56.000000	459470.860000

Data Columns

```
Out[16]: df.columns
```

```
Index(['account', 'first_name', 'age', 'gender', 'balance', 'months', 'weekpayment', 'annual'], dtype='object')
```

```
In [17]: df.drop(['account', 'first_name'], axis=1, inplace=True)
```

```
In [18]: df
```

```
Out[18]:
```

	age	gender	balance	months	weekpayment	annual
0	40	F	51472.20	12	52	46388.68
1	22	M	298308.49	12	24	166140.52
2	52	M	86898.05	12	28	196860.56
3	27	F	13769.63	8	8	252908.24
4	38	M	22826.60	12	52	52110.76
...
95	28	M	88169.32	12	28	132011.36
96	35	M	65301.33	12	24	120050.84
97	34	F	59807.25	12	28	182915.72
98	21	M	467645.22	12	52	81130.40
99	21	M	49877.75	12	28	133791.32
100 rows x 6 columns						

Perform One-Hot Encoding

```
In [19]: df2 = pd.get_dummies(df, drop_first=True)
```

```
Out[19]: df2
```

```
Out[20]:
```

	age	balance	months	weekpayment	annual	gender_M	annual
0	40	51472.20	12	52	46388.68	0	0
1	22	298308.49	12	24	166140.52	1	1
2	52	86898.05	12	28	196860.56	1	1
3	27	13769.63	8	8	252908.24	0	0
4	38	22826.60	12	52	52110.76	1	1
...
95	28	88169.32	12	28	132011.36	1	1
96	35	65301.33	12	24	120050.84	1	1
97	34	59807.25	12	28	182915.72	0	0
98	21	467645.22	12	52	81130.40	1	1
99	21	49877.75	12	28	133791.32	1	1
100 rows x 6 columns							

Treat Data Types

```
In [24]: df2.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100 entries, 0 to 99
Data columns (total 6 columns):
 #   Column             Non-Null Count  Dtype
---  -
 0   age                100 non-null    int64
 1   balance            100 non-null    float64
 2   months            100 non-null    int64
 3   weekpayment        100 non-null    float64
 4   gender_M           100 non-null    int64
 5   annual             100 non-null    float64
dtypes: float64(2), int64(3), uint8(1)
memory usage: 4.1 KB
```

Create and save processed dataset

```
In [25]: df2.to_csv('train.csv', index=False)
```

Train Test Split

```
In [26]: df2.shape
```

```
Out[26]: (100, 6)
```

```
In [27]: X = df2.iloc[:,0:5]
```

```
y = df2.iloc[:,5]
```

```
In [28]: X.values, y.values
```

```
Out[28]: (array([[ 40., 51472.2, 12., 52., 46388.68, 0.],
[ 22., 298308.49, 12., 24., 166140.52, 1.],
[ 52., 86898.05, 12., 28., 196860.56, 1.],
[ 27., 13769.63, 8., 8., 252908.24, 0.],
[ 38., 22826.6, 12., 52., 52110.76, 1.],
[ 44., 28575.85, 12., 24., 196860.56, 1.],
[ 42., 49934.37, 12., 24., 196860.56, 1.],
[ 0.4189059, 3.37023929, 0.1938551, 1.15155984, 0.90453403],
[ 0.45924236, -0.3149824, -0.56622805, -1.88249003, 0.90453403],
[ 0.9087965, -0.27173445, 0.1938551, -0.50337666, 0.90453403],
[ 0.79860474, 0.07080292, 0.1938551, 1.15155984, -1.10554616],
[ 0.63910843, -0.4915554, 0.1938551, -0.50337666, -1.10554616],
[ 1.36674669, -0.35863004, 0.1938551, 1.15155984, -1.10554616],
[ 0.09896599, -0.49860125, 0.1938551, -0.77919917, -1.10554616],
[ 0.9983673, -0.36452686, 0.1938551, 1.42738255, -1.10554616],
[ 0.45939429, -0.37032236, 0.1938551, -0.77919917, -1.10554616],
[ 0.8182258, -0.14449507, 0.1938551, 1.42738255, -1.10554616],
[ 0.79860474, 0.439392341, 0.1938551, -0.50337666, -1.10554616],
[ 0.8182258, 0.5626018, 0.1938551, -0.77919917, -1.10554616],
[ 0.2594623, -0.15173715, 0.1938551, 1.15155984, -1.10554616],
[ 0.79860474, -0.37087054, 0.1938551, -1.052189, -1.10554616],
[ 0.63910843, -0.38898433, 0.1938551, 1.1308446, -1.10554616],
[ 0.6189059, -0.50832633, 0.1938551, 1.15155984, 0.90453403],
[ 0.79860474, 0.5527517, 0.1938551, -1.3308446, -1.10554616],
[ 0.45939429, -0.48018586, 0.1938551, 1.42738255, -1.10554616],
[ 0.45939429, -0.48854831, 0.1938551, -0.77919917, -1.10554616],
[ 0.8182258, -0.37032236, 0.1938551, -0.77919917, -1.10554616],
[ 0.8182258, -0.14449507, 0.1938551, 1.42738255, -1.10554616],
[ 0.79860474, 0.439392341, 0.1938551, -0.50337666, -1.10554616],
[ 0.8182258, 0.5626018, 0.1938551, -0.77919917, -1.10554616],
[ 0.2594623, -0.15173715, 0.1938551, 1.15155984, -1.10554616],
[ 0.79860474, -0.37087054, 0.1938551, -1.052189, -1.10554616],
[ 0.63910843, -0.38898433, 0.1938551, 1.1308446, -1.10554616],
[ 0.6189059, -0.50832633, 0.1938551, 1.15155984, 0.90453403],
[ 0.79860474, 0.5527517, 0.1938551, -1.3308446, -1.10554616],
[ 0.45939429, -0.48018586, 0.1938551, 1.42738255, -1.10554616],
[ 0.45939429, -0.48854831, 0.1938551, -0.77919917, -1.10554616],
[ 0.8182258, -0.37032236, 0.1938551, -0.77919917, -1.10554616],
[ 0.8182258, -0.14449507, 0.1938551, 1.42738255, -1.10554616],
[ 0.79860474, 0.439392341, 0.1938551, -0.50337666, -1.10554616],
[ 0.8182258, 0.5626018, 0.1938551, -0.77919917, -1.10554616],
[ 0.2594623, -0.15173715, 0.1938551, 1.15155984, -1.10554616],
[ 0.79860474, -0.37087054, 0.1938551, -1.052189, -1.10554616],
[ 0.63910843, -0.38898433, 0.1938551, 1.1308446, -1.10554616],
[ 0.6189059, -0.50832633, 0.1938551, 1.15155984, 0.90453403],
[ 0.79860474, 0.5527517, 0.1938551, -1.3308446, -1.10554616],
[ 0.45939429, -0.48018586, 0.1938551, 1.42738255, -1.10554616],
[ 0.45939429, -0.48854831, 0.1938551, -0.77919917, -1.10554616],
[ 0.8182258, -0.37032236, 0.1938551, -0.77919917, -1.10554616],
[ 0.8182258, -0.14449507, 0.1938551, 1.42738255, -1.10554616],
[ 0.79860474, 0.439392341, 0.1938551, -0.50337666, -1.10554616],
[ 0.8182258, 0.5626018, 0.1938551, -0.77919917, -1.10554616],
[ 0.2594623, -0.15173715, 0.1938551, 1.15155984, -1.10554616],
[ 0.79860474, -0.37087054, 0.1938551, -1.052189, -1.10554616],
[ 0.63910843, -0.38898433, 0.1938551, 1.1308446, -1.10554616],
[ 0.6189059, -0.50832633, 0.1938551, 1.15155984, 0.90453403],
[ 0.79860474, 0.5527517, 0.193
```


Using PyCaret

In [37]: exp_reg = setup(data = df2, target = 'annual', session_id=0, normalize=True, train_size=0.8, numeric_features=['months','weekspayment'],categorical_features=['gender_M'])

Setup Successfully Completed.

	Description	Value
0	session_id	0
1	Transform Target	False
2	Transform Target Method	None
3	Original Data	(100, 6)
4	Missing Values	False
5	Numeric Features	4
6	Categorical Features	1
7	Ordinal Features	False
8	High Cardinality Features	False
9	High Cardinality Method	None
10	Sampled Data	(100, 6)
11	Transformed Train Set	(80, 6)
12	Transformed Test Set	(20, 6)
13	Numeric Imputer	mean
14	Categorical Imputer	constant
15	Normalize	True
16	Normalize Method	zscore
17	Transformation	False
18	Transformation Method	None
19	PCA	False
20	PCA Method	None
21	PCA Components	None
22	Ignore Low Variance	False
23	Combine Rare Levels	False
24	Rare Level Threshold	None
25	Numeric Binning	False
26	Remove Outliers	False
27	Outliers Threshold	None
28	Remove Multicollinearity	False
29	Multicollinearity Threshold	None
30	Clustering	False
31	Clustering Iteration	None
32	Polynomial Features	False
33	Polynomial Degree	None
34	Trigonometry Features	False
35	Polynomial Threshold	None
36	Group Features	False
37	Feature Selection	False
38	Features Selection Threshold	None
39	Feature Interaction	False
40	Feature Ratio	(20, 6)
41	Interaction Threshold	None

In [38]: compare_models()

Model	MAE	MSE	RMSE	R2	RMSLE	MAPE	TT(Sec)
0 Random Forest	29301.5291	2163230089.6387	42516.7020	0.3629	0.3022	0.2774	0.1758
1 Extra Trees Regressor	30909.3561	2334604019.2963	43150.5733	0.3501	0.3118	0.2863	0.1388
2 AdaBoost Regressor	32635.1157	2172994922.1505	42943.7851	0.3479	0.3456	0.3309	0.0787
3 Gradient Boosting Regressor	33365.8582	2111731516.0444	43946.7740	0.3123	0.3464	0.3276	0.0417
4 Extreme Gradient Boosting	32658.4512	2493513976.4447	44966.4866	0.3076	0.3346	0.2863	0.0910
5 CatBoost Regressor	32752.6319	2563181216.0432	44964.5924	0.2778	0.3548	0.3265	0.7828
6 Elastic Net	36823.9344	3099121552.3342	50633.3856	0.2085	0.4152	0.4046	0.0553
7 Random Sample Consensus	36678.6640	3169892088.8563	51390.4456	0.2040	0.3875	0.3116	0.0752
8 Huber Regressor	36626.1631	3066209514.8814	50686.0571	0.1856	0.3800	0.3248	0.0100
9 Orthogonal Matching Pursuit	39774.7737	3078734214.2651	51049.8354	0.1016	0.4021	0.3714	0.0016
0 Decision Tree	38305.2540	3219682271.8766	54047.7257	0.0111	0.4014	0.3432	0.0031
11 Ridge Regression	37439.5700	3056883144.6375	49804.0027	-0.0341	0.3849	0.3508	0.0016
12 Lasso Regression	37472.8130	3115766385.0226	50108.2217	-0.0517	0.3853	0.3498	0.0069
13 Linear Regression	37472.8836	3115661068.1703	50108.0670	-0.0518	0.3853	0.3498	0.0069
14 Lasso Least Angle Regression	37563.1804	3132128731.3781	50300.1457	-0.0563	0.3853	0.3494	0.0057
15 Least Angle Regression	37564.0646	3132528967.2398	50303.1164	-0.0567	0.3853	0.3494	0.0031
16 Passive Aggressive Regressor	46532.2455	5627295984.2722	67094.2223	-0.0808	0.4866	0.5131	0.0116
17 Support Vector Machine	53994.5347	5864698168.7939	70794.2223	-0.1538	0.5781	0.6424	0.0062
18 Bayesian Ridge	55199.3108	5613319379.3464	69890.1957	-0.1657	0.5996	0.6451	0.0069
19 K Neighbors Regressor	41543.4890	3839643993.7948	56916.9658	-0.1670	0.4273	0.3905	0.0047
20 Light Gradient Boosting Machine	42654.6880	4078318142.5164	57905.7787	-0.1886	0.4003	0.3800	0.0272
21 TheilSen Regressor	46391.7065	537021710.1379	65878.8083	-1.4292	0.4317	0.4083	0.5651

Out[38]: RandomForestRegressor(max_depth=None, max_features='auto', max_leaf_nodes=None, max_samples=None, min_impurity_decrease=0.0, min_impurity_split=None, min_samples_leaf=1, min_samples_split=2, min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=-1, oob_score=False, random_state=0, verbose=0, warm_start=False)

Simple Regression Model

In [39]: lr = LinearRegression()

In [40]: lr.fit(X_train_scaled, y_train)

Out[40]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)

In [41]: lr_pred = lr.predict(X_test_scaled)

In [42]: y_test

Out[42]: array([[145019.66871768, 70166.40490483, 127494.26832432, 38474.65966601, 13716.71187725, 13713.18928733, 146180.56661006, 155001.58108455, 138113.88280777, 136931.33224313, 67170.98989244, 146612.28747075, 54513.50717847, 64019.18012509, 18838.92082806, 143956.72181563, 136841.79357281, 54146.26294868, 28980.84667239, 139453.3436454]])

In [44]: y_test

Out[44]: array([[189774.52, 37716.64, 196860.56, 54242.24, 47671. , 60223.8 , 79959.36, 210848.04, 173096.04, 130196.04, 51100.92, 99658.52, 66168.44, 134576.52, 141362.52, 91406.12, 206827.92, 72565.48, 253566.4 , 128463.4]])

In [45]: mse = mean_squared_error(y_test,lr_pred)

Out[45]: 1742935754.6537423

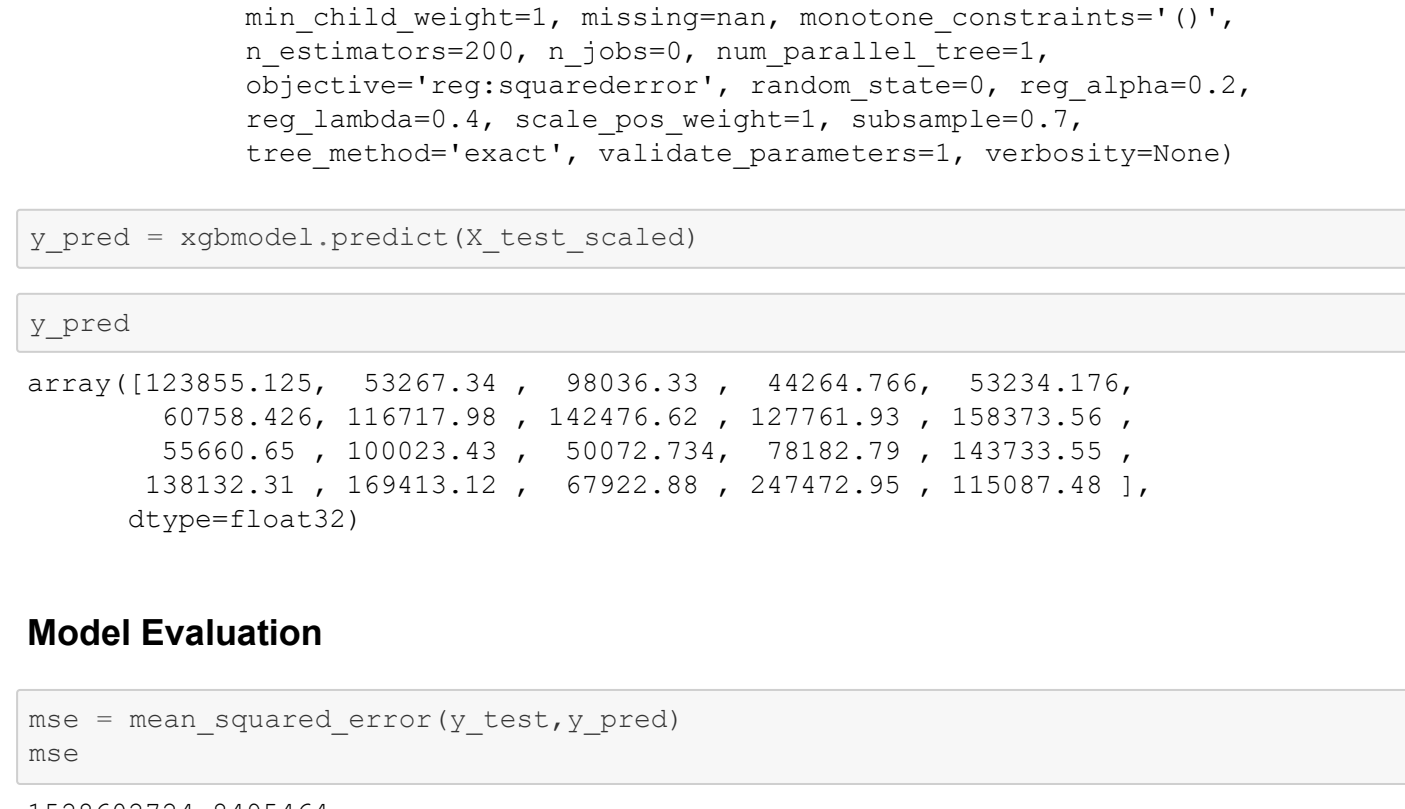
In [46]: rmse = np.sqrt(mse)

Out[46]: 41748.4820646722

In [47]: r2score = r2_score(y_test,lr_pred)

Out[47]: 0.5706038490094812

In [48]: fig, ax = plt.subplots(figsize=(10,8))
sns.regplot(x=y_test, y=lr_pred, ax=ax)
plt.title("Plot to compare actual vs predicted")
plt.ylabel("Predicted")
plt.xlabel("Actual")
plt.show()



Using XGBoost (Scikit-Learn)

Using RandomSearchCV

In [49]: model = XGBRegressor(random_state=0, n_estimators=100, objective='reg:squarederror')

In [50]: parameters = {'max_depth': np.arange(3,10,1),
'learning_rate': np.arange(0.05,0.3,0.05),
'n_estimators':np.arange(100,1000,100),
'min_child_weight': np.arange(1,4,1),
'gamma':np.arange(0,10,2),
'subsample':np.arange(0.5,0.9,0.1),
'colsample_bytree':np.arange(0.5,0.9,0.1),
'reg_alpha':np.arange(0,1,0.1),
'reg_lambda':np.arange(0,1,0.1)
}

In [51]: random = RandomizedSearchCV(estimator=model, param_distributions = parameters, cv = 5, n_iter = 50,
n_jobs=-1, scoring='neg_mean_squared_error',random_state=0)

In [52]: random.fit(X_train_scaled, y_train)

Out[52]: RandomizedSearchCV(cv=5, error_score=nan,
estimator=XGBRegressor(base_score=None, booster=None,
colsample_bylevel=None,
colsample_bynode=None,
colsample_bytree=None, gamma=None,
gpu_id=None, importance_type='gain',
interaction_constraints=None,
learning_rate=None,
max_delta_step=None, max_depth=None,
min_child_weight=None, missing=nan,
monotone_constraints=None,
n_jobs=-1,
'min_child_weight': array([1, 2, 3]),
'n_estimators': array([100, 200, 300, 400, 500, 600, 700, 800, 900]),
'reg_alpha': array([0. , 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9]),
'reg_lambda': array([0. , 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9]),
'subsample': array([0.5, 0.6, 0.7, 0.8])),
pre_dispatch='2*n_jobs', random_state=None, refit=True,
return_train_score=False, scoring='neg_mean_squared_error',
verbose=0)

In [53]: random.best_estimator_

Out[53]: XGBRegressor(base_score=0.5, booster='gbtree', colsample_bylevel=1,
colsample_bynode=1, colsample_bytree=0.8, eta=0.25, gamma=4,
importance_type='gain', interaction_constraints='',
learning_rate=0.15000000000000002, learning_rate=0.300000012,
max_delta_step=0, max_depth=3, min_child_weight=1, missing=nan,
monotone_constraints=(), n_estimators=800, n_jobs=0,
num_parallel_tree=1, objective='reg:squarederror', random_state=0,
reg_alpha=0.4, reg_lambda=0.1, scale_pos_weight=1, subsample=0.7,
tree_method='exact', validate_parameters=1, verbosity=None)

In [54]: random.best_score

Out[54]: -2594910809.860133

In [55]: random.best_params_

Out[55]: {'subsample': 0.6,
'reg_lambda': 0.1,
'reg_alpha': 0.4,
'n_estimators': 800,
'min_child_weight': 1,
'max_depth': 3,
'learning_rate': 0.15000000000000002,
'gamma': 2,
'colsample_bytree': 0.5}

Final Model

In [56]: xgbmodel = XGBRegressor(random_state=0, n_estimators=200, objective='reg:squarederror', subsample=0.7,
reg_lambda=0.4,reg_alpha=0.2,min_child_weight=1,max_depth=4,eta=0.25, gamma=4,
colsample_bytree=0.8)

In [57]: xgbmodel.fit(X_train_scaled,y_train,eval_set=[(X_test_scaled,y_test)],eval_metric='rmse',early_stopping_rounds=10)

[0] validation_0-rmse:108464.35156
Will train until validation_0-rmse hasn't improved in 10 rounds.
[1] validation_0-rmse:87001.39062
[2] validation_0-rmse:71516.21094
[3] validation_0-rmse:62403.91797
[4] validation_0-rmse:52160.75391
[5] validation_0-rmse:45983.99453
[6] validation_0-rmse:42712.09766
[7] validation_0-rmse:40499.71484
[8] validation_0-rmse:39906.14453
[9] validation_0-rmse:39356.49603
[10] validation_0-rmse:39097.35156
[11] validation_0-rmse:40342.96875
[12] validation_0-rmse:40113.82031
[13] validation_0-rmse:39186.51562
[14] validation_0-rmse:40743.67578
[15] validation_0-rmse:41439.43359
[16] validation_0-rmse:41451.07812
[17] validation_0-rmse:41661.97656
[18] validation_0-rmse:42399.53516
[19] validation_0-rmse:42048.81250
[20] validation_0-rmse:41679.58594
Stopping. Best iteration:
[10] validation_0-rmse:39097.35156

Out[57]: XGBRegressor(base_score=0.5, booster='gbtree', colsample_bylevel=1,
colsample_bynode=1, colsample_bytree=0.8, eta=0.25, gamma=4,
gpu_id=-1, importance_type='gain', interaction_constraints='',
learning_rate=0.25, max_delta_step=0, max_depth=4,
min_child_weight=1, missing=nan, monotone_constraints=(),
n_estimators=200, n_jobs=0, num_parallel_tree=1,
objective='reg:squarederror', random_state=0, reg_alpha=0.2,
reg_lambda=0.4, scale_pos_weight=1, subsample=0.7,
tree_method='exact', validate_parameters=1, verbosity=None)

In [58]: y_pred = xgbmodel.predict(X_test_scaled)

In [59]: y_pred

Out[59]: array([[123855.125, 53267.34 , 98036.33 , 44264.766, 53334.176,
60788.426, 16717.98 , 142476.62 , 127761.93 , 158379.56 ,
55660.65 , 100023.43 , 50072.734, 78182.79 , 143733.55 ,
138132.31 , 169413.12 , 67922.88 , 247472.95 , 115087.48 ,
dtype=float32])

Model Evaluation

In [60]: mse = mean_squared_error(y_test,y_pred)

Out[60]: 1528602724.8495464

In [61]: rmse = np.sqrt(mse)

Out[61]: 39097.34933278146

In [62]: r2score = r2_score(y_test,y_pred)

Out[62]: 0.6234077333651278

In [63]: fig, ax = plt.subplots(figsize=(10,8))
sns.regplot(x=y_test, y=y_pred, ax=ax)
plt.title("Plot to compare actual vs predicted")
plt.ylabel("Predicted")
plt.xlabel("Actual")
plt.show()



Cross-Validation

In [67]: cv = cross_val_score(xgbmodel,X,y,cv=5,verbose=1,scoring='r2')
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 5 out of 5 | elapsed: 0.3s finished

Out[68]: cv.mean()

Out[68]: 0.4714286293167129

Save the Model

In [69]: filename = 'anmodel.sav'

dump(xgbmodel,open(filename,'wb'))

Conclusion: Tree based model is much better in accuracy than linear regression.

In []: