

▼ Import All necessity Functions

```
1 ##### Import all necessity functions for Machine Learning #####
2 import sys
3 import math
4 import string
5 import numpy as np
6 import pandas as pd
7 import seaborn as sns
8 import matplotlib.pyplot as plt
9 import scipy as shc
10 import warnings
11 import zipfile
12 import cv2
13 import os
14 import random
15 from collections import Counter
16 from functools import reduce
17 from itertools import chain
18 from google.colab.patches import cv2_imshow
19 from keras.preprocessing import image
20 from sklearn.metrics._plot.confusion_matrix import confusion_matrix
21 from sklearn.model_selection import train_test_split, KFold, StratifiedKFold, GridSearchCV
22 from sklearn.preprocessing import StandardScaler, RobustScaler, MinMaxScaler
23 from sklearn.decomposition import PCA
24 from sklearn.cluster import KMeans, DBSCAN, AgglomerativeClustering
25 from sklearn.feature_selection import mutual_info_classif, mutual_info_regression, SelectK
26 from imblearn.under_sampling import RandomUnderSampler, NearMiss
27 from imblearn.over_sampling import RandomOverSampler, SMOTE, SMOTEN, SMOTENC, SVMSMOTE, KM
28 from imblearn.ensemble import EasyEnsembleClassifier
29 from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
30 from sklearn.naive_bayes import GaussianNB, BernoulliNB, MultinomialNB
31 from sklearn.neighbors import KNeighborsClassifier, KNeighborsRegressor, NearestNeighbors
32 from sklearn.linear_model import LinearRegression, LogisticRegression, SGDClassifier, SGDR
33 from sklearn.neural_network import MLPClassifier, MLPRegressor
34 from sklearn.svm import SVC, SVR
35 from sklearn.tree import DecisionTreeClassifier, DecisionTreeRegressor, ExtraTreeClassifier
36 from sklearn.ensemble import BaggingClassifier, BaggingRegressor, RandomForestClassifier,
37 from sklearn.ensemble import AdaBoostClassifier, AdaBoostRegressor, GradientBoostingClassifier
38 from sklearn.metrics import classification_report, mean_absolute_error, mean_squared_error
39 from xgboost import XGBClassifier, XGBRegressor
40
41 ##### Download keras #####
42 !pip install keras
43
44 ##### Remove all warnings #####
45 import warnings
```

```
46 warnings.filterwarnings("ignore")
47
48 ##### Import all necessity functions for Neural Network #####
49 import tensorflow as tf
50 from keras.models import Sequential, Model
51 from keras.utils import plot_model
52 from keras.layers import Dense, Conv2D, LSTM, GRU, RNN, Flatten, AvgPool2D, MaxPool2D, Glo
53 from keras.activations import tanh, relu, sigmoid, softmax, swish
54 from keras.regularizers import L1, L2, L1L2
55 from keras.optimizers import SGD, Adagrad, Adadelta, RMSprop, Adam, Adamax, Nadam
56 from keras.initializers import HeNormal, HeUniform, GlorotNormal, GlorotUniform
57 from keras.losses import SparseCategoricalCrossentropy, CategoricalCrossentropy, hinge, MS
58 import keras.utils as image
59 from google.colab.patches import cv2_imshow
60 from keras.utils import plot_model
61
62 ##### Plotting the confusion matrix #####
63 from mlxtend.evaluate import confusion_matrix
64 from mlxtend.plotting import plot_confusion_matrix
65 from sklearn.metrics import multilabel_confusion_matrix
66 from sklearn.metrics import confusion_matrix
```

Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public>
Requirement already satisfied: keras in /usr/local/lib/python3.9/dist-packages (2.11.0)



▼ Import the dataset

```
1 df = pd.read_csv('/content/train.csv')
2 df.head()
```

| Customer ID | Name | Gender | Age | Income (USD) | Income Stability | Profession | Type of Employment | Location |
|-------------|------|--------|-----|--------------|------------------|------------|--------------------|----------|
|-------------|------|--------|-----|--------------|------------------|------------|--------------------|----------|

▼ Determine the number of records of this dataset

```

1 C-33999  ANURAG  M  32  4952.91  Low  Working  NaN  Semi-urban
1 def find_feature_record(df):
2     if df.empty:
3         print('The dataset is empty'.capitalize())
4     else:
5         print('The # of records of this dataset is {}'.format(df.shape[0]))
6         print('The # of columns of this dataset is {}'.format(df.shape[1]))
7
8 try:
9     find_feature_record(df)
10 except Exception as e:
11     print(e.with_traceback())

```

The # of records of this dataset is 30000
The # of columns of this dataset is 24

▼ Determine the shape of this dataset

```

1 def shape_dataset(dataset):
2     if dataset.empty:
3         print('The dataset is empty'.capitalize())
4     else:
5         print('The shape of this dataset is {}'.format(dataset.shape))
6
7 try:
8     shape_dataset(df)
9 except Exception as e:
10     print(e.with_traceback())

```

The shape of this dataset is (30000, 24)

▼ Determine how many NaN value presence and their column name

```

1 def check_nan(df):
2     if df.empty:
3         print('The dataset is empty'.capitalize())
4     else:
5         return df.isnull().sum()[df.isnull().sum() > 0]
6

```

```

7 try:
8     nan_value = check_nan(df)
9 except Exception as e:
10     print(e.with_traceback)
11 else:
12     for feature, value in zip(nan_value.index, nan_value.values):
13         print('{} feature has NaN value # {}'.format(feature, value), '\n')
14         print('*'*40)
15
16     print('The total NaN value is in this dataset is = {}'.format(df.isnull().sum().sum()))

```

Gender feature has NaN value # 53

Income (USD) feature has NaN value # 4576

Income Stability feature has NaN value # 1683

Type of Employment feature has NaN value # 7270

Current Loan Expenses (USD) feature has NaN value # 172

Dependents feature has NaN value # 2493

Credit Score feature has NaN value # 1703

Has Active Credit Card feature has NaN value # 1566

Property Age feature has NaN value # 4850

Property Location feature has NaN value # 356

Loan Sanction Amount (USD) feature has NaN value # 340

The total NaN value is in this dataset is = 25062

▼ Determine those features have NaN value are categorical or Neumrical

```

1 def find_categorical_numerical(df):
2     if df.empty:
3         print('The dataset is empty'.capitalize())
4     else:

```

```

5     return df.dtypes[df.dtypes.index.isin(df.isnull().sum()[df.isnull().sum() > 0].index)]
6
7 try:
8     finding = find_categorical_numerical(df)
9 except Exception as e:
10    print(e.with_traceback)
11 else:
12    for feature, value in zip(finding.index, finding.values):
13        print('{} feature datatype is # {}'.format(feature, value))
14        print('-'*50)

```

```

Gender feature datatype is # object
-----
Income (USD) feature datatype is # float64
-----
Income Stability feature datatype is # object
-----
Type of Employment feature datatype is # object
-----
Current Loan Expenses (USD) feature datatype is # float64
-----
Dependents feature datatype is # float64
-----
Credit Score feature datatype is # float64
-----
Has Active Credit Card feature datatype is # object
-----
Property Age feature datatype is # float64
-----
Property Location feature datatype is # object
-----
Loan Sanction Amount (USD) feature datatype is # float64
-----

```

▼ Split the dataset into Categorical and Numerical for analysing

```

1 def split_cat_numerical(cat = None, num = None):
2     if cat == 'categorical':
3         categorical_df = df[df.dtypes[df.dtypes == 'object'].index]
4         return categorical_df
5     elif(num == 'numerical'):
6         numerical_df = df[df.dtypes[df.dtypes != 'object'].index]
7         return numerical_df
8     else:
9         raise Exception('Splitting cannot be done.'.capitalize())
10
11 try:
12     cat_df = split_cat_numerical(cat = 'categorical', num = None)
13     num_df = split_cat_numerical(cat = None, num = 'numerical')
14 except Exception as e:
15     print(e.with_traceback)
16 else:

```

```
10 else:
```

```
17 print('done'.capitalize())
```

Done

```
1 cat_df.head()
```

| | Customer ID | Name | Gender | Income Stability | Profession | Type of Employment | Location | Expense Type 1 | Exp Ty |
|---|-------------|-------------------|--------|------------------|------------|--------------------|------------|----------------|--------|
| 0 | C-36995 | Frederica Shealy | F | Low | Working | Sales staff | Semi-Urban | N | |
| 1 | C-33999 | America Calderone | M | Low | Working | NaN | Semi-Urban | N | |
| 2 | C-3770 | Rosetta Verne | F | High | Pensioner | NaN | Semi-Urban | N | |

```
1 num_df.head()
```

| | Age | Income (USD) | Loan Amount Request (USD) | Current Loan Expenses (USD) | Dependents | Credit Score | No. of Defaults | Property ID | Property Age | P |
|---|-----|--------------|---------------------------|-----------------------------|------------|--------------|-----------------|-------------|--------------|---|
| 0 | 56 | 1933.05 | 72809.58 | 241.08 | 3.0 | 809.44 | 0 | 746 | 1933.05 | |
| 1 | 32 | 4952.91 | 46837.47 | 495.81 | 1.0 | 780.40 | 0 | 608 | 4952.91 | |
| 2 | 65 | 988.19 | 45593.04 | 171.95 | 1.0 | 833.15 | 0 | 546 | 988.19 | |
| 3 | 65 | NaN | 80057.92 | 298.54 | 2.0 | 832.70 | 1 | 890 | NaN | |
| 4 | 31 | 2614.77 | 113858.89 | 491.41 | NaN | 745.55 | 1 | 715 | 2614.77 | |

▼ EDA - Exploratory data analysis

▼ Delete the Customer ID and Name from the dataset

```
1 def remove_features(features = None):
2     if len(features) == 1:
3         df.drop(features, axis = 1, inplace = True)
4     else:
5         for column in features:
6             df.drop(column, axis = 1, inplace = True)
7             print('{} column is deleted from the dataset.'.format(column))
8
```

```

9 try:
10     remove_features(['Customer ID', 'Name'])
11 except Exception as e:
12     print(e.with_traceback())

Customer ID column is deleted from the dataset.
Name column is deleted from the dataset.

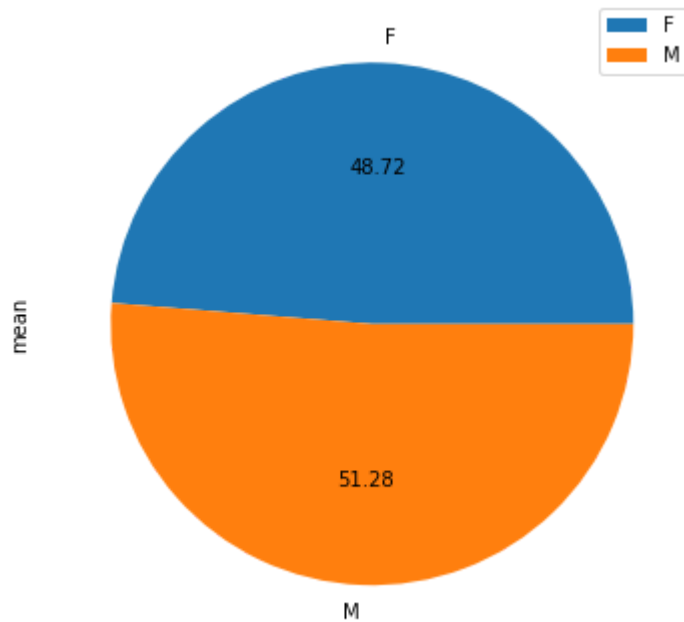
```

▼ The distribution of average income based on Gender

```

1 df.groupby(['Gender'])['Income (USD)'].\\
2     agg(['mean']).plot(kind = 'pie',\\
3     subplots = True, autopct='%0.2f', figsize = (8, 6))
4 plt.show()

```

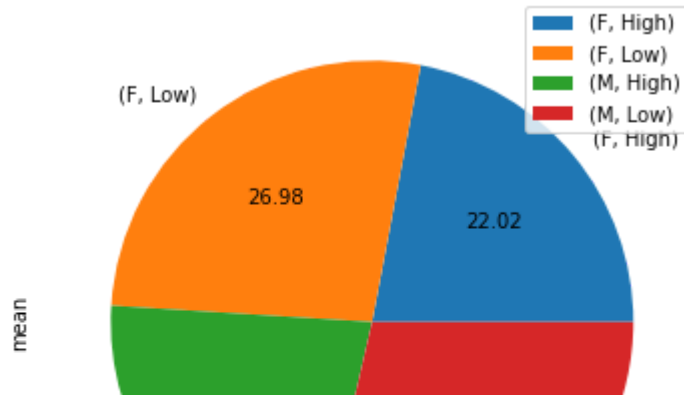


▼ The distribution of average income based on Gender and Income Stability

```

1 df.groupby(['Gender', 'Income Stability'])['Income (USD)'].\\
2     agg(['mean']).plot(kind = 'pie',\\
3     subplots = True, autopct='%0.2f', figsize = (8, 6)
4 plt.show()

```



Check, what's the maximum, minimum, mean Loan Amount Request (USD) based on the Profession

```
1 df.groupby(['Profession'])['Loan Amount Request (USD)']. \
2     agg(['max', 'min', 'mean']). \
3     sort_values(by = 'mean', ascending = False)
```

| | max | min | mean |
|----------------------|-----------|-----------|---------------|
| Profession | | | |
| Maternity leave | 108967.56 | 108967.56 | 108967.560000 |
| State servant | 417936.91 | 6310.26 | 97543.778331 |
| Commercial associate | 602384.15 | 6174.70 | 96521.471468 |
| Pensioner | 387288.84 | 6498.75 | 86344.841022 |
| Businessman | 94747.20 | 74449.41 | 84598.305000 |
| Working | 621497.82 | 6048.24 | 84391.484042 |
| Unemployed | 89038.91 | 77097.81 | 83068.360000 |
| Student | 58056.11 | 58056.11 | 58056.110000 |

Find out, the max-5 Type of Employment based on their Loan Amount Request (USD)

```
1 df.groupby(['Type of Employment'])['Loan Amount Request (USD)'].mean(). \
2     sort_values(ascending = False) \
3     plot(kind = 'barh', figsize = (10, 10)) \
4 plt.show()
```

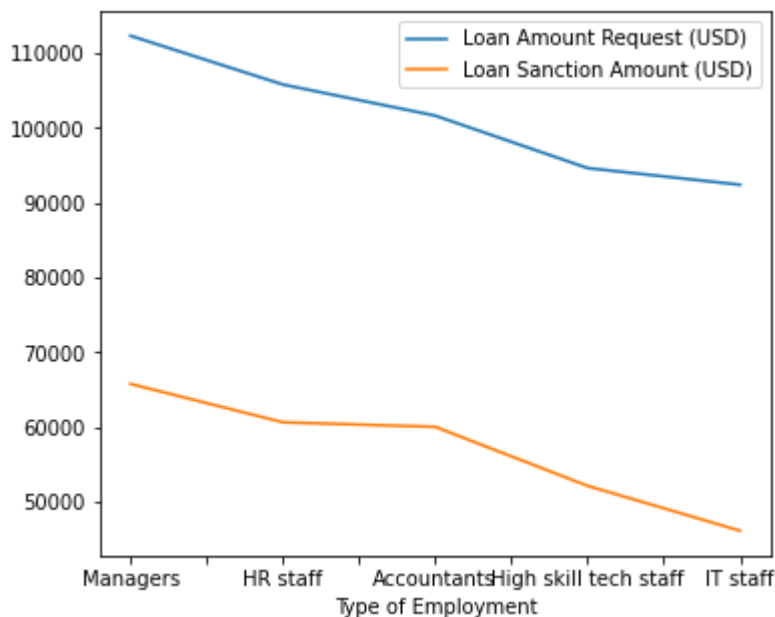



Find out, the mean of requesting Loan Amount and getting Loan sanction based on the top5 Type of Employment

```

1 top_5 = df.groupby(['Type of Employment'])['Loan Amount Request (USD)'].mean().\
2           sort_values(ascending = False).head(5).reset_index()
3
4 index = df.groupby(['Type of Employment'])['Loan Amount Request (USD)'].mean().\
5           sort_values(ascending = False).head(5).reset_index()
6
7 dummy_df = df[df.loc[:, 'Type of Employment'].isin(index)]
8
9 top_5['Loan Sanction Amount (USD)'] = dummy_df.groupby(['Type of Employment'])['Loan Sanction Amount (USD)'].mean().\
10           sort_values(ascending = False).head(5).reset_index()
11
12 top_5.set_index('Type of Employment').plot(figsize = (6, 5))
13 plt.show()

```



```

1 top_5 = df.groupby(['Type of Employment'])['Loan Amount Request (USD)'].mean().\
2           sort_values(ascending = True).head(5).reset_index()

```

```

3
4 index = df.groupby(['Type of Employment'])['Loan Amount Request (USD)'].mean().\
5                                     sort_values(ascending = True)
6
7 dummy_df = df[df.loc[:, 'Type of Employment'].isin(index)]
8
9 top_5['Loan Sanction Amount (USD)'] = dummy_df.groupby(['Type of Employment'])['Loan Sanct
10                                     sort_values(ascending = True).head(5)
11
12 top_5.set_index('Type of Employment').plot(figsize = (6, 5))
13 plt.show()

```



▼ Show, the relationship of Income Stability, location with the features like Loan amount and Loan Sanction

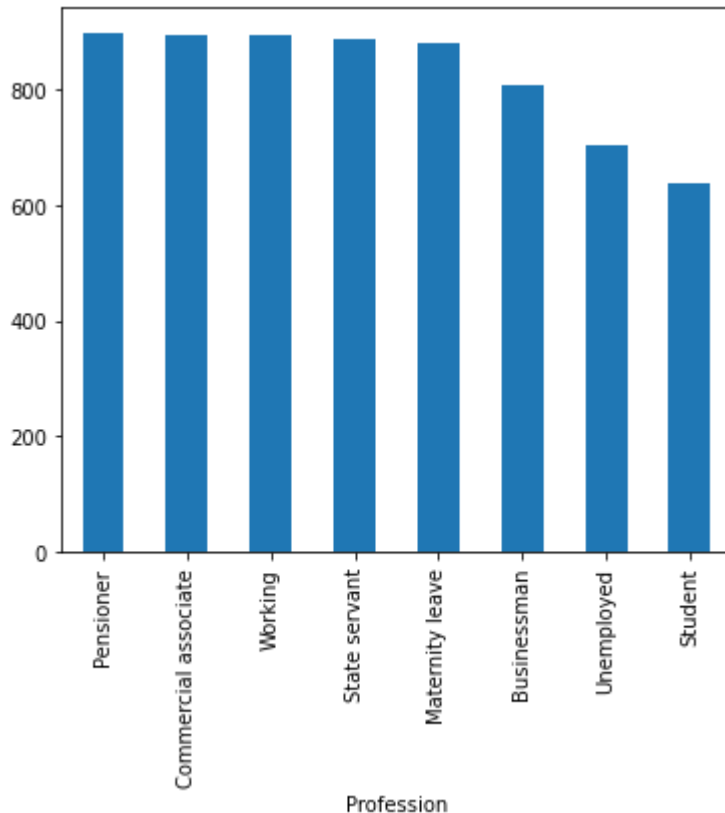
```

1 pd.pivot_table(index = 'Income Stability', columns = ['Location'],\
2               values = ['Loan Amount Request (USD)', 'Loan Sanction Amount (USD)'],\
3               aggfunc = ['mean'], data = df).plot(kind = 'pie', subplots = True, figsize
4 plt.show()

```

▼ Find out, the highest credit score based on their Profession

```
1 df.groupby(['Profession'])['Credit Score'].max().sort_values(ascending = False).plot(kind
2 plt.show())
```



▼ Find out, the top Credit Score using Profession & show their mean, max, min, count of Type of Employment and their Loan request and Loan Sanction

```
1 def find_queries(group):
2     return group['Credit Score'].max()
3
4 df[df.loc[:, 'Profession'].isin(df.groupby(['Profession']).apply(find_queries).sort_values
5                                loc[:, ['Profession', 'Loan Amount Request (USD)', 'Loan Sanctio
6                                groupby(['Profession']).agg(['min', 'max', 'mean'])).head()]
```

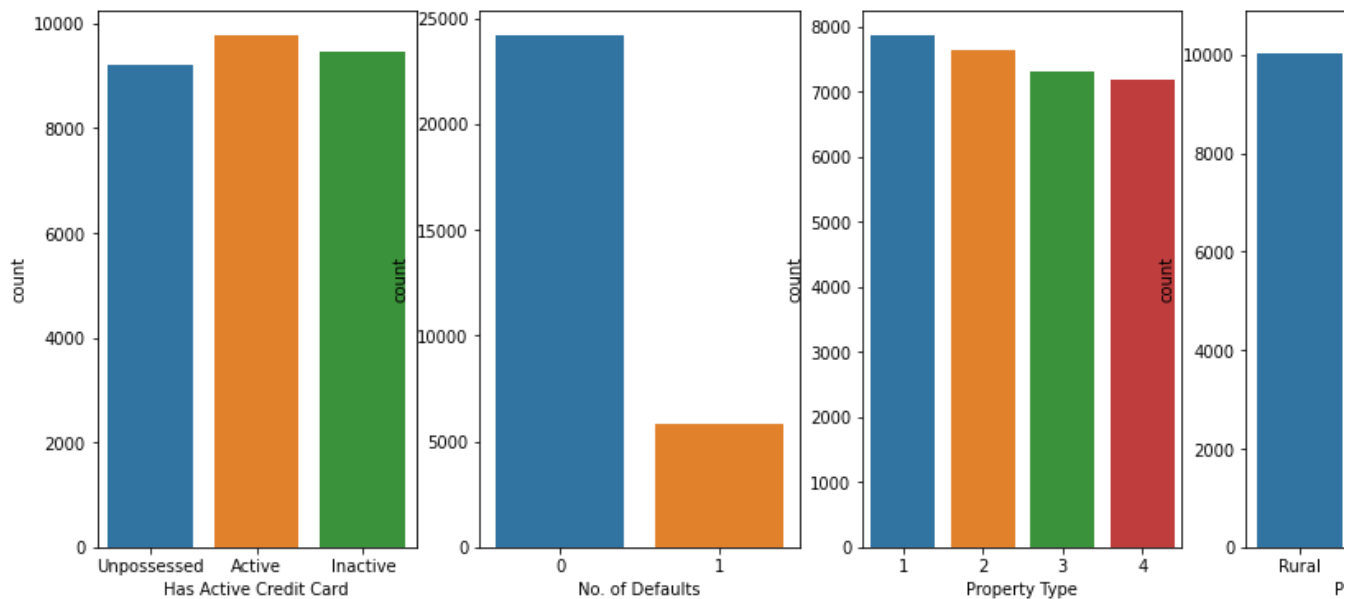
| Profession | Loan Amount Request (USD) | | | Loan Sanction Amount (USD) | | |
|-------------|---------------------------|----------|--------------|----------------------------|----------|-------------|
| | min | max | mean | min | max | mean |
| Businessman | 74449.41 | 94747.20 | 84598.305000 | 59559.53 | 66323.04 | 62941.28500 |

▼ Show the countplot of the Has Active Credit Card column

```

1 plt.figure(figsize = (16, 6))
2 plt.subplot(1, 4, 1)
3 sns.countplot(df.loc[:, 'Has Active Credit Card'])
4
5 plt.subplot(1, 4, 2)
6 sns.countplot(df.loc[:, 'No. of Defaults'])
7
8 plt.subplot(1, 4, 3)
9 sns.countplot(df.loc[:, 'Property Type'])
10
11 plt.subplot(1, 4, 4)
12 sns.countplot(df.loc[:, 'Property Location'])
13 plt.show()

```

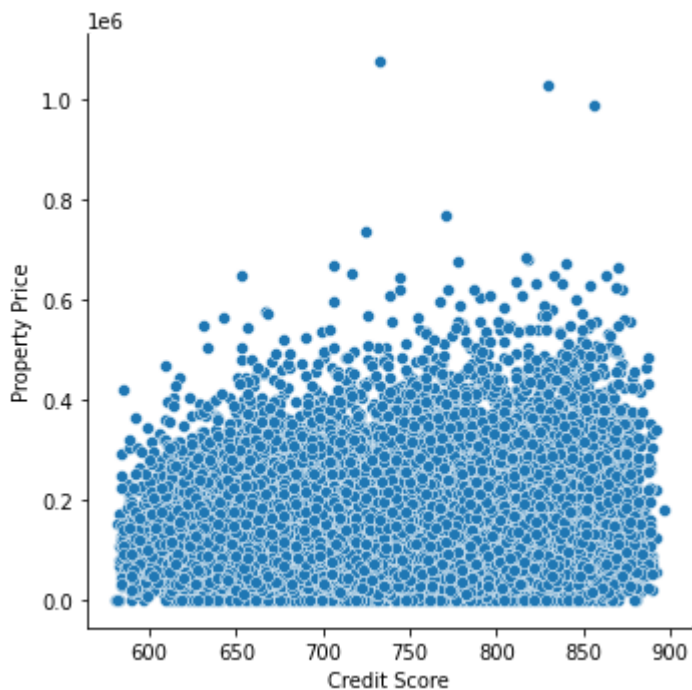


▼ Show the relationship of Credit Score and Property Price

```

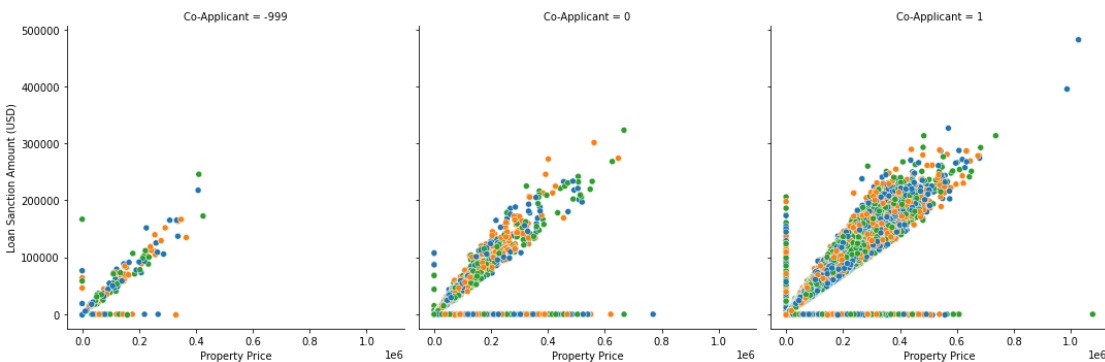
1 sns.relplot(data = df, x = 'Credit Score', y = 'Property Price', kind = 'scatter')
2 plt.show()

```




▼ Find out, the relationship of Property location, Property price and Loan Sanction

```
1 sns.relplot(x = 'Property Price', y = 'Loan Sanction Amount (USD)', hue = 'Property Location')
2 plt.show()
```



▼ Find out, the top Property price based on their location

```
1 df.groupby(['Property Location'])['Property Price'].agg(['max', 'count', 'sum', 'mean'])
```

| | max | count | sum | mean |  |
|-------------------|------------|-------|--------------|---------------|---|
| Property Location | | | | | |
| Rural | 1028082.64 | 10041 | 1.319019e+09 | 131363.317538 | |
| Semi-Urban | 1077966.73 | 10387 | 1.367362e+09 | 131641.637175 | |
| Urban | 678932.62 | 9216 | 1.219563e+09 | 132331.057394 | |

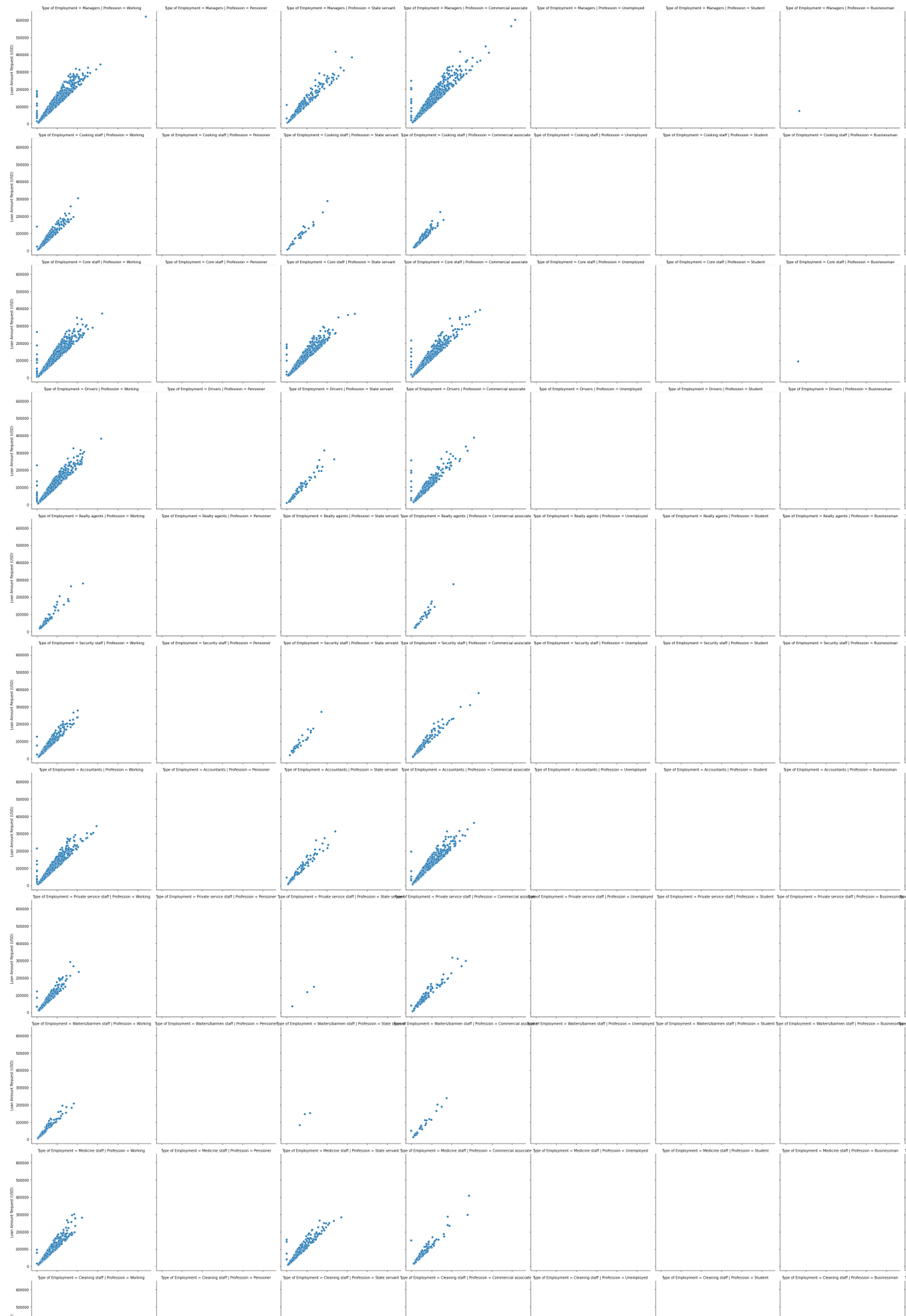
Find out, the mean and count of Loan request & Loan Sanction based on property type

```
1 pd.pivot_table(index = ['Property Location'], columns = ['Property Type'], values = ['Loan
2                                                                                               'Loan
3                                                                                               aggf
4                                                                                               data
```

| mean | | Loan Amount Request (USD) | | | | Loan Sanction Amount (U | |
|---------------|-------------------|---------------------------|--------------|--------------|--------------|-------------------------|------------|
| Property Type | Property Location | 1 | 2 | 3 | 4 | 1 | 2 |
| | | | | | | | |
| Rural | | 88687.272223 | 87595.814497 | 88685.207034 | 88950.360624 | 47502.564716 | 47529.9458 |
| Semi-Urban | | 89271.307594 | 88719.609821 | 88428.335978 | 88298.671714 | 48131.483707 | 47081.5727 |
| Urban | | 88832.400781 | 87706.488145 | 88261.747884 | 88506.063780 | 47704.207004 | 46333.3403 |

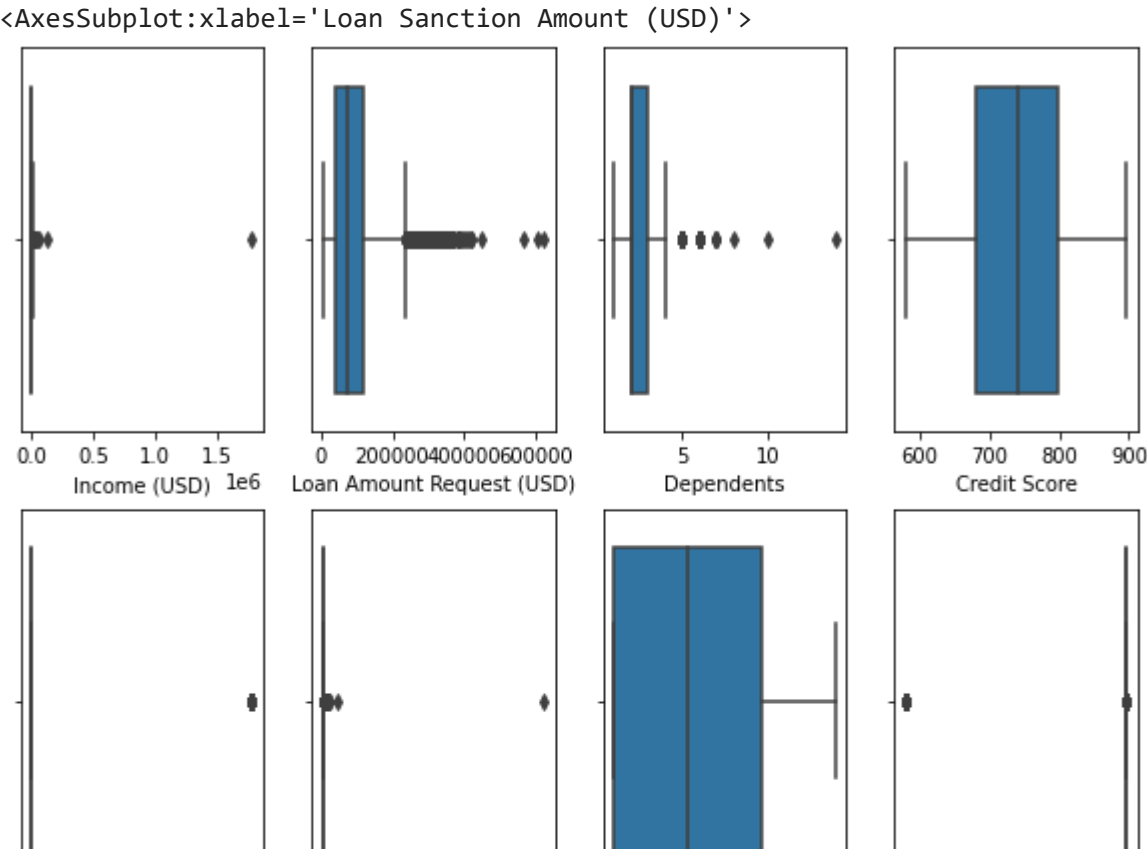
Show the realtionship of Property Price, Loan Request, Profession, Type of Employment

```
1 sns.relplot(x = 'Property Price', y = 'Loan Amount Request (USD)', col = 'Profession', row
```



▼ Check there is an outliers or not in the dataset

```
1 plt.figure(figsize = (10, 12))
2 plt.subplot(3, 4, 1)
3 sns.boxplot(df.loc[:, 'Income (USD)'])
4
5 plt.subplot(3, 4, 2)
6 sns.boxplot(df.loc[:, 'Loan Amount Request (USD)'])
7
8 plt.subplot(3, 4, 3)
9 sns.boxplot(df.loc[:, 'Dependents'])
10
11 plt.subplot(3, 4, 4)
12 sns.boxplot(df.loc[:, 'Credit Score'])
13
14 plt.subplot(3, 4, 5)
15 sns.boxplot(df.loc[:, 'No. of Defaults'])
16
17 plt.subplot(3, 4, 6)
18 sns.boxplot(df.loc[:, 'Property Age'])
19
20 plt.subplot(3, 4, 7)
21 sns.boxplot(df.loc[:, 'Property Type'])
22
23 plt.subplot(3, 4, 8)
24 sns.boxplot(df.loc[:, 'Co-Applicant'])
25
26 plt.subplot(3, 4, 9)
27 sns.boxplot(df.loc[:, 'Property Price'])
28
29 plt.subplot(3, 4, 10)
30 sns.boxplot(df.loc[:, 'Loan Sanction Amount (USD)'])
```

▼ Show the description of this entire dataset

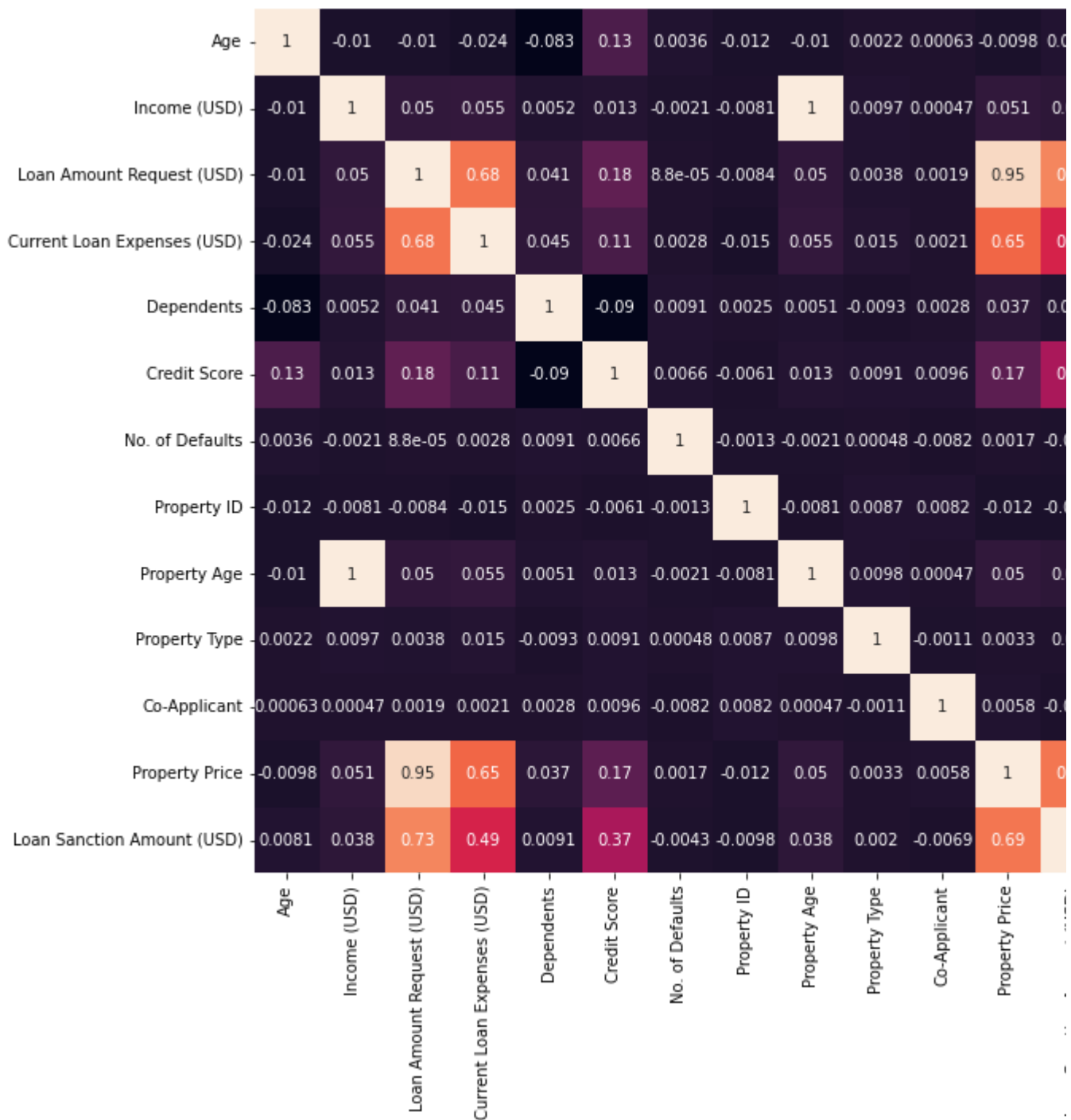
```
1 df.describe()
```

| | Age | Income (USD) | Loan Amount Request (USD) | Current Loan Expenses (USD) | Dependents | Credit Score |
|-------|--------------|--------------|---------------------------|-----------------------------|--------------|--------------|
| count | 30000.000000 | 2.542400e+04 | 30000.000000 | 29828.000000 | 27507.000000 | 28297.000000 |
| mean | 40.092300 | 2.630574e+03 | 88826.333855 | 400.936876 | 2.253027 | 739.885381 |
| std | 16.045129 | 1.126272e+04 | 59536.949605 | 242.545375 | 0.951162 | 72.163846 |
| min | 18.000000 | 3.777000e+02 | 6048.240000 | -999.000000 | 1.000000 | 580.000000 |
| 25% | 25.000000 | 1.650457e+03 | 41177.755000 | 247.667500 | 2.000000 | 681.880000 |
| 50% | 40.000000 | 2.222435e+03 | 75128.075000 | 375.205000 | 2.000000 | 739.820000 |
| 75% | 55.000000 | 3.090593e+03 | 119964.605000 | 521.292500 | 3.000000 | 799.120000 |
| max | 65.000000 | 1.777460e+06 | 621497.820000 | 3840.880000 | 14.000000 | 896.260000 |



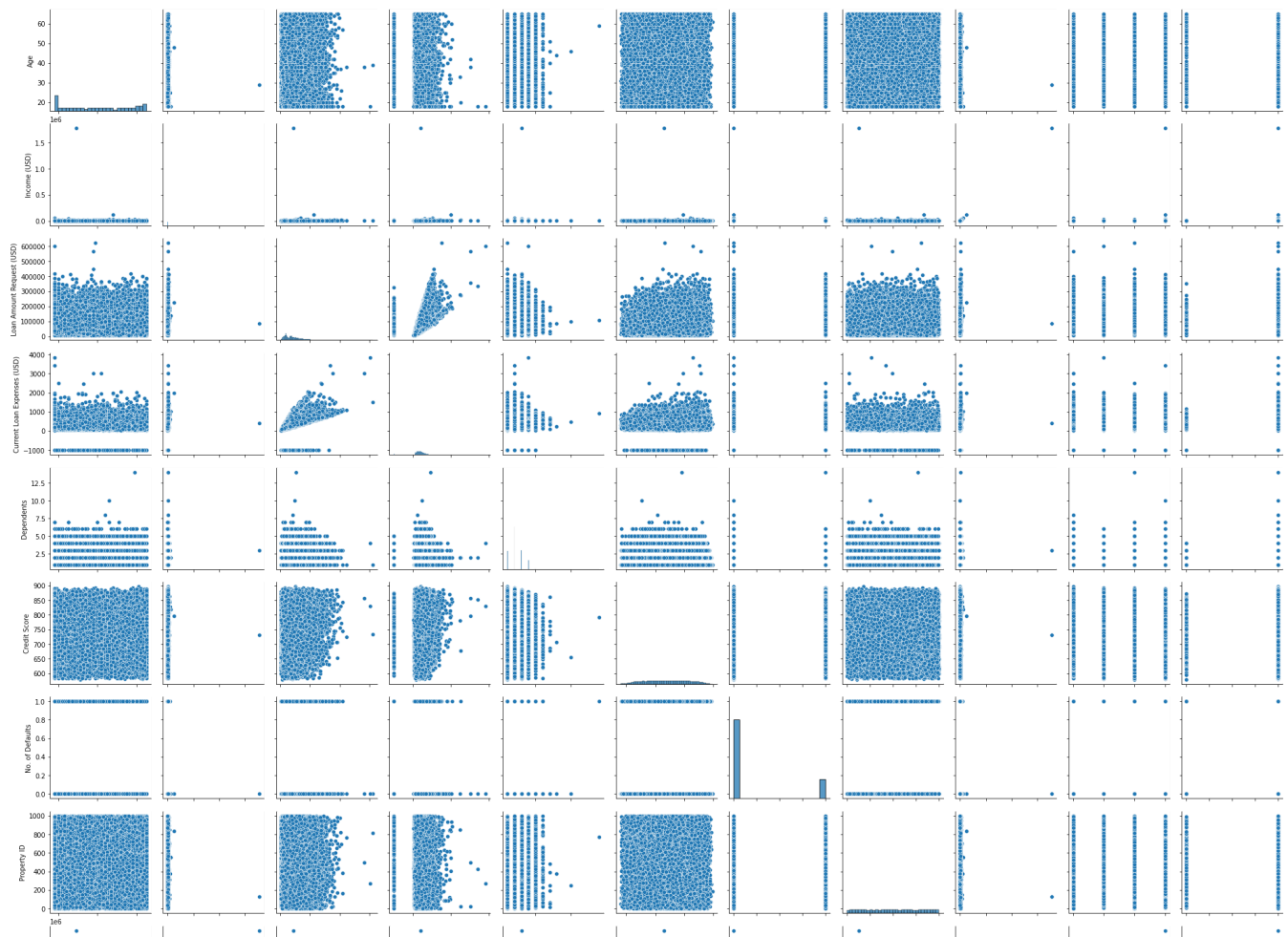
▼ Plot the correlation of this dataset

```
1 plt.figure(figsize = (12, 10))
2 sns.heatmap(df.corr(), annot = True)
3 plt.show()
```



▼ Plot the pairplot of this entire dataset

```
1 sns.pairplot(df)
2 plt.show()
```



▼ Feature Engineering for the Dataset

Handle the Missing Value

Handle the Categorical Dataset

Splitting the dataset

Normalization Approach



▼ Handle the Missing Value



```

1 ##### For ``Gender`` - We will use the new gender called T #####
2 def NaN_Fill(activate = None, feature = None):
3     if activate == 'activated':
4         df.loc[:, feature] = np.where(df.loc[:, feature].isnull(), 'T', df.loc[:, feature])
5     else:
6         raise Exception('Not possible.')
7
8 try:
9     NaN_Fill(activate = 'activated', feature = 'Gender')

```

```

10 except Exception as e:
11     print(e.with_traceback)
12 else:
13     print('Done')

```

Done

▼ For Type of Employment - We will use the new gender called homie

```

1 def NaN_Fill(activate = None, feature = None):
2     if activate == 'activated':
3         df.loc[:, feature] = np.where(df.loc[:, feature].isnull(), 'homie', df.loc[:, feature])
4     else:
5         raise Exception('Not possible.')
6
7 try:
8     NaN_Fill(activate = 'activated', feature = 'Type of Employment')
9 except Exception as e:
10     print(e.with_traceback)
11 else:
12     print('Done')

```

Done

▼ For Loan Sanction Amount (USD) - We will delete all records

```

1 def dropna_feature(activate = None, feature = None):
2     if activate == 'activated':
3         df.dropna(subset = [feature], axis = 0, inplace = True)
4     else:
5         raise Exception('Not possible.')
6
7 try:
8     dropna_feature(activate = 'activated', feature = 'Loan Sanction Amount (USD)')
9 except Exception as e:
10     print(e.with_traceback)
11 else:
12     print('Done')

```

Done

▼ For Income Stability - We will use the new gender called Missing

```

1 def NaN_Fill(activate = None, feature = None):
2     if activate == 'activated':
3         df.loc[:, feature] = np.where(df.loc[:, feature].isnull(), 'Missing', df.loc[:, featur

```

```

4  else:
5      raise Exception('Not possible.')
6
7  try:
8      NaN_Fill(activate = 'activated', feature = 'Income Stability')
9  except Exception as e:
10     print(e.with_traceback)
11 else:
12     print('Done')

```

Done

▼ For Has Active Credit Card - We will use the new gender called Missing

```

1 def NaN_Fill(activate = None, feature = None):
2     if activate == 'activated':
3         df.loc[:, feature] = np.where(df.loc[:, feature].isnull(), 'Missing', df.loc[:, feature])
4     else:
5         raise Exception('Not possible.')
6
7  try:
8      NaN_Fill(activate = 'activated', feature = 'Has Active Credit Card')
9  except Exception as e:
10     print(e.with_traceback)
11 else:
12     print('Done')

```

Done

```

1 try:
2     NaN_Fill(activate = 'activated', feature = 'Property Location')
3 except Exception as e:
4     print(e.with_traceback)
5 else:
6     print('Done')

```

Done

```

1 # def mean_encoding(activate = None, feature = None):
2 #     if activate == 'activated':
3 #         df.loc[:, feature] = df.loc[:, feature].std()
4 #     else:
5 #         raise Exception('Not possible')
6
7 # try:
8 #     for features in ['Income (USD)', 'Current Loan Expenses (USD)', 'Dependents', 'Credit
9 #         mean_encoding(activate = 'activated', feature = features)
10 # except Exception as e:

```

```

11 # print(e.with_traceback)
12 # else:
13 # print('Done')

```

▼ Do the random sample imputation for all Neumerical features

```

1 def impute_nan_random_sample(activate = None, column_name_ = None):
2     if activate == 'activated':
3         index_null_values_ = df.loc[df.loc[:, column_name_].isnull(), :][column_name_].index
4         random_values_ = df.loc[:, column_name_].dropna().sample(df.loc[:, column_name_].i
5         for value_, index_ in enumerate(df.loc[df.loc[:, column_name_].isnull()][column_name_]
6             df.loc[index_, column_name_] = random_values_[value_]
7
8
9     try:
10     for features in ['Income (USD)', 'Current Loan Expenses (USD)', 'Dependents', 'Credit Sc
11         impute_nan_random_sample(activate = 'activated', column_name_ = features)
12 except Exception as e:
13     print(e.with_traceback)
14 else:
15     print('DONE')

```

DONE

▼ Check NaN value presence in this dataset or not

```

1 if df.isnull().sum().sum() == 0:
2     print('There is No NaN value in this dataset.'.capitalize())
3 else:
4     print('There is NaN value in this dataset.'.capitalize())

```

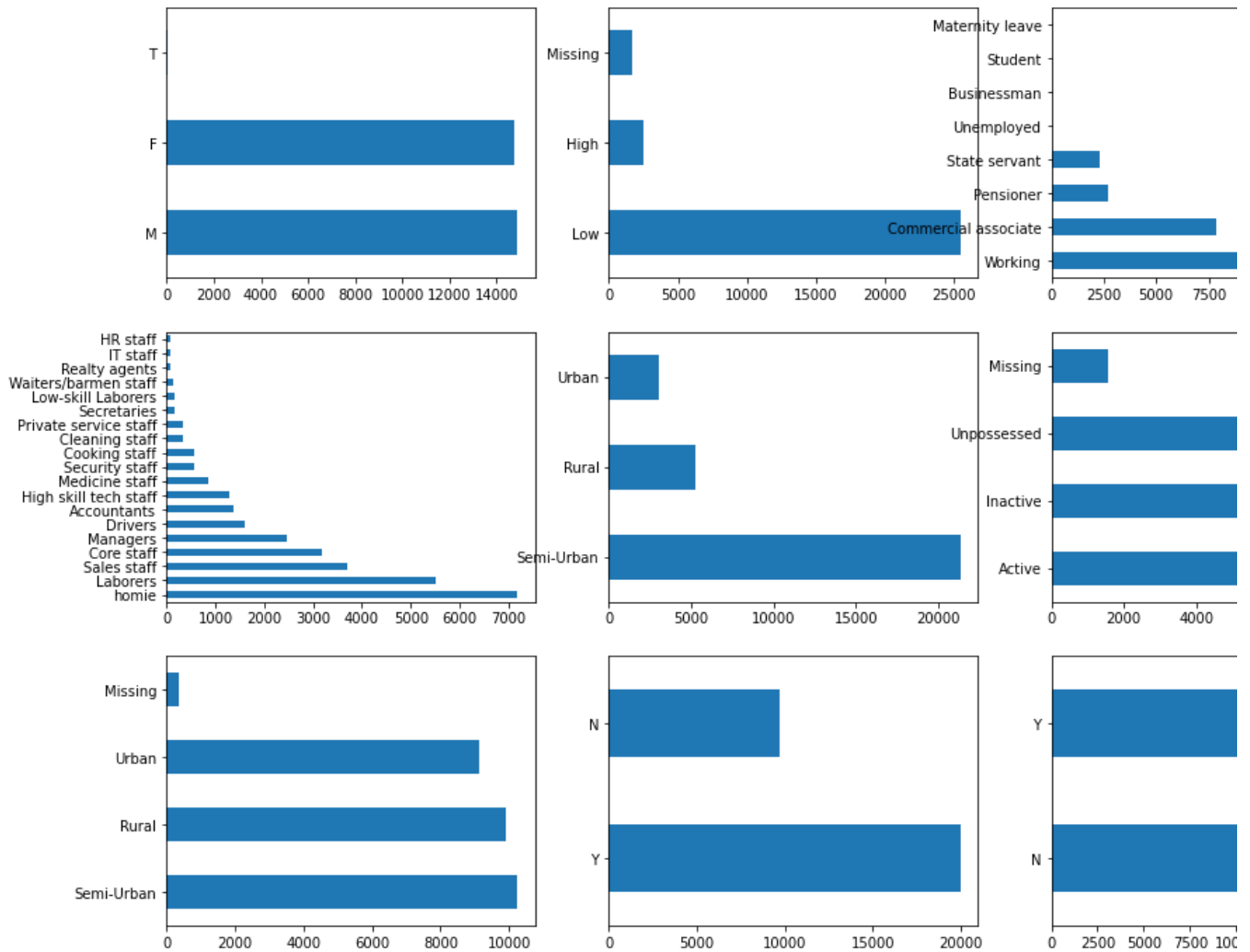
There is no nan value in this dataset.

▼ Handle the Categorical Data

```

1 plt.figure(figsize = (16, 12))
2 for index, feature in enumerate(['Gender', 'Income Stability', 'Profession', \
3                                 'Type of Employment', 'Location', 'Has Active Credit Card
4                                 'Property Location', 'Expense Type 2', 'Expense Type 1'])
5     plt.subplot(3, 3, index + 1)
6     df.loc[:, feature].value_counts().plot(kind = 'barh')

```



```

1 def mean_encoding_all(activate = None, feature = None):
2     if activate == 'activated':
3         df.loc[:, feature] = df.loc[:, feature].map(df.loc[:, feature].value_counts().to_dict(
4     else:
5         raise Exception('Not possible.'.capitalize())
6
7 try:
8     for index, feature_ in enumerate(['Gender', 'Income Stability', 'Profession', 'Type of E
9         'Location', 'Has Active Credit Card', 'Property Locati
10            'Expense Type 1', 'Expense Type 2']):
11         mean_encoding_all(activate = 'activated', feature = feature_)
12 except Exception as e:
13     print(e.with_traceback)
14 else:
15     print('DONE')

```

DONE


```
1 df.head()
```

| | Gender | Age | Income (USD) | Income Stability | Profession | Type of Employment | Location | Loan Amount Request (USD) | Current Loan Expenses (USD) |
|---|--------|-----|--------------|------------------|------------|--------------------|----------|---------------------------|-----------------------------|
| 0 | 14718 | 56 | 1933.05 | 25458 | 16739 | 3698 | 21317 | 72809.58 | 241.08 |
| 1 | 14890 | 32 | 4952.91 | 25458 | 16739 | 7188 | 21317 | 46837.47 | 495.81 |
| 2 | 14718 | 65 | 988.19 | 2544 | 2718 | 7188 | 21317 | 45593.04 | 171.95 |
| 3 | 14718 | 65 | 1918.47 | 2544 | 2718 | 7188 | 5280 | 80057.92 | 298.54 |
| 4 | 14718 | 31 | 2614.77 | 25458 | 16739 | 1297 | 21317 | 113858.89 | 491.41 |

5 rows × 22 columns



```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 29660 entries, 0 to 29999
Data columns (total 22 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Gender                                29660 non-null  int64
1   Age                                   29660 non-null  int64
2   Income (USD)                          29660 non-null  float64
3   Income Stability                       29660 non-null  int64
4   Profession                             29660 non-null  int64
5   Type of Employment                    29660 non-null  int64
6   Location                              29660 non-null  int64
7   Loan Amount Request (USD)             29660 non-null  float64
8   Current Loan Expenses (USD)           29660 non-null  float64
9   Expense Type 1                        29660 non-null  int64
10  Expense Type 2                        29660 non-null  int64
11  Dependents                            29660 non-null  float64
12  Credit Score                          29660 non-null  float64
13  No. of Defaults                       29660 non-null  int64
14  Has Active Credit Card                 29660 non-null  int64
15  Property ID                           29660 non-null  int64
16  Property Age                          29660 non-null  float64
17  Property Type                         29660 non-null  int64
18  Property Location                     29660 non-null  int64
19  Co-Applicant                         29660 non-null  int64
20  Property Price                        29660 non-null  float64
21  Loan Sanction Amount (USD)            29660 non-null  float64
dtypes: float64(8), int64(14)
memory usage: 6.2 MB
```

```
1 X = df.iloc[:, :-1]
2 y = df.iloc[:, -1]
```

▼ Normalizing the dataset

```
1 minmax_scaler = MinMaxScaler()
2 X_normalized = minmax_scaler.fit_transform(df)
3 df_normalised = pd.DataFrame(X_normalized, columns = df.columns)
4 df_normalised.head()
```

| | Gender | Age | Income (USD) | Income Stability | Profession | Type of Employment | Location | Loan Amount Request (USD) | Ex |
|---|----------|----------|--------------|------------------|------------|--------------------|----------|---------------------------|----|
| 0 | 0.988408 | 0.808511 | 0.000875 | 1.000000 | 1.000000 | 0.509625 | 1.000000 | 0.108476 | 0. |
| 1 | 1.000000 | 0.297872 | 0.002575 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 0.066276 | 0. |
| 2 | 0.988408 | 1.000000 | 0.000344 | 0.037227 | 0.162325 | 1.000000 | 1.000000 | 0.064254 | 0. |
| 3 | 0.988408 | 1.000000 | 0.000867 | 0.037227 | 0.162325 | 1.000000 | 0.121453 | 0.120253 | 0. |
| 4 | 0.988408 | 0.276596 | 0.001259 | 1.000000 | 1.000000 | 0.172264 | 1.000000 | 0.175174 | 0. |

5 rows × 22 columns



▼ Split the dataset into train and test

```
1 X = df_normalised.iloc[:, :-1]
2 y = df_normalised.iloc[:, -1]
```

▼ Splitting the dataset into train and test

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 2)
2
3 print('The shape of X_train is = {}'.format(X_train.shape), '\n')
4 print('The shape of X_test is = {}'.format(X_test.shape), '\n')
5 print('The shape of y_train is = {}'.format(y_train.shape), '\n')
6 print('The shape of y_test is = {}'.format(y_test.shape), '\n')
```

The shape of X_train is = (22245, 21)

The shape of X_test is = (7415, 21)

The shape of y_train is = (22245,)

The shape of y_test is = (7415,)

▼ Model Building

```

1 models = {
2     'KNN': KNeighborsRegressor(),
3     'LR' : LinearRegression(),
4     'SVR': SVR(),
5     'DT' : DecisionTreeRegressor(),
6     'RF' : RandomForestRegressor(),
7     'GB' : GradientBoostingRegressor(),
8     'XGB': XGBRegressor()
9 }
10
11 MAE, MSE, RMSE, R2 = [], [], [], []
12
13 for model_name, model_object in models.items():
14     model = model_object
15     model.fit(X_train, y_train)
16
17     predicted = model.predict(X_test)
18
19     MAE.append(mean_absolute_error(predicted, y_test))
20     MSE.append(mean_squared_error(predicted, y_test))
21     RMSE.append(np.sqrt(mean_squared_error(predicted, y_test)))
22     R2.append(r2_score(predicted, y_test))

1 model = ['KNN', 'LR', 'SVR', 'DT', 'RF', 'GB', 'XGB']
2
3 for model_name, mae, mse, rmse, r2 in zip(model, MAE, MSE, RMSE, R2):
4     print('{} MAE is {}, MSE is {}, RMSE is {}, R2 Score is {}'.format(model_name, mae, mse
5     print('*'*120)


KNN MAE is 0.0638534407305653, MSE is 0.007184624663268531, RMSE is 0.25269238360220775,
*****
LR MAE is 0.04479903060591726, MSE is 0.004227066493187742, RMSE is 0.2116578148945067,
*****
SVR MAE is 0.06573061933298312, MSE is 0.005911744153942446, RMSE is 0.2563798340996872,
*****
DT MAE is 0.031672381366825396, MSE is 0.005453135327984628, RMSE is 0.1779673603974206,
*****

```

```
RF MAE is 0.02665696077995848, MSE is 0.002710158228648405, RMSE is 0.16326959539350394,
*****
GB MAE is 0.02848786629850074, MSE is 0.0026038292436277428, RMSE is 0.16878348941321464
*****
XGB MAE is 0.026043110250265854, MSE is 0.0026510122012625963, RMSE is 0.161378778810184
*****
```

▼ Check for RandomForestRegressor

```
1 LR = RandomForestRegressor()
2 LR.fit(X_train, y_train)
3
4 predicted = LR.predict(X_test)
5
6 evaluation_df = pd.DataFrame(predicted, columns = ['Predict'])
7 pd.concat([evaluation_df, pd.DataFrame(y_test.values, columns = ['Actual'])], axis = 1).he
```

| | Predict | Actual |  |
|---|----------|----------|---|
| 0 | 0.093459 | 0.115137 | |
| 1 | 0.085973 | 0.106983 | |
| 2 | 0.016879 | 0.020934 | |
| 3 | 0.166799 | 0.002069 | |
| 4 | 0.072688 | 0.081956 | |

▼ Select the Important Features from the dataset

Using Pearson Correlation Technique

RandomForestRegressor used for finding the feature Importance

▼ Pearson Correlation Technique

```
1 def correlation(dataset_, threshold_):
2
3     col_corr_ = set()
4     corr_matrix = dataset_.corr()
5     for i in range(len(corr_matrix.columns)):
6         for j in range(i):
```

```

7         if abs(corr_matrix.iloc[i, j]) > threshold_:
8             col_name_ = corr_matrix.columns[i]
9             col_corr_.add(col_name_)
10
11     return col_corr_

```

```

1 correlation(df_normalised, 0.65)

{'Current Loan Expenses (USD)',
 'Loan Sanction Amount (USD)',
 'Property Age',
 'Property Price'}

```


▼ Use Random Forest to find the Feature Importance

```

1 RF = RandomForestRegressor()
2 RF.fit(X_train, y_train)
3
4 predicted = RF.predict(X_test)

1 feature_importance = pd.DataFrame(RF.feature_importances_, columns = ['Feature Importance']
2 feature_importance.set_index(X.columns).sort_values(by = 'Feature Importance', ascending =

```

| | Feature Importance  |
|------------------------------------|---|
| Loan Amount Request (USD) | 0.577052 |
| Credit Score | 0.160554 |
| Co-Applicant | 0.073832 |
| Income Stability | 0.030457 |
| Property Price | 0.023473 |
| Property ID | 0.022180 |
| Current Loan Expenses (USD) | 0.019599 |
| Age | 0.017363 |
| Property Age | 0.015077 |
| Income (USD) | 0.013291 |

```

1 df_important = df_normalised.loc[:,feature_importance.set_index(X.columns).\
2     sort_values(by = 'Feature Importance',\
3     ascending = False)[0:7].index]
4
5 df_important.head()

```

| | Loan Amount Request (USD) | Credit Score | Co- Applicant | Income Stability | Property Price | Property ID |
|---|------------------------------|-----------------|------------------|---------------------|-------------------|----------------|
| 0 | 0.108476 | 0.725479 | 1.000 | 1.000000 | 0.112082 | 0.746493 |
| 1 | 0.066276 | 0.633656 | 1.000 | 1.000000 | 0.051707 | 0.608216 |
| 2 | 0.064254 | 0.800449 | 0.999 | 0.037227 | 0.068065 | 0.546092 |
| 3 | 0.120253 | 0.799026 | 1.000 | 0.037227 | 0.113480 | 0.890782 |
| 4 | 0.175174 | 0.523462 | 1.000 | 1.000000 | 0.194229 | 0.715431 |

Split the dataset into independent & dependent feature with Feature Selection -

Random Forest

```
1 X_important = df_important
```

```
1 X_train, X_test, y_train, y_test = train_test_split(X_important, y, test_size = 0.25, rand
2
3 print('The shape of X_train is = {}'.format(X_train.shape), '\n')
4 print('The shape of X_test is = {}'.format(X_test.shape), '\n')
5 print('The shape of y_train is = {}'.format(y_train.shape), '\n')
6 print('The shape of y_test is = {}'.format(y_test.shape), '\n')
```

```
The shape of X_train is = (22245, 7)
```

```
The shape of X_test is = (7415, 7)
```

```
The shape of y_train is = (22245,)
```

```
The shape of y_test is = (7415,)
```

Model Building

```
1 models = {
2     'KNN': KNeighborsRegressor(),
3     'LR' : LinearRegression(),
4     'SVR': SVR(),
5     'DT' : DecisionTreeRegressor(),
6     'RF' : RandomForestRegressor(),
7     'GB' : GradientBoostingRegressor(),
8     'XGB': XGBRegressor()
9 }
10
11 MAE, MSE, RMSE, R2 = [], [], [], []
12
```

```

13 for model_name, model_object in models.items():
14     model = model_object
15     model.fit(X_train, y_train)
16
17     predicted = model.predict(X_test)
18
19     MAE.append(mean_absolute_error(predicted, y_test))
20     MSE.append(mean_squared_error(predicted, y_test))
21     RMSE.append(np.sqrt(mean_squared_error(predicted, y_test)))
22     R2.append(r2_score(predicted, y_test))

1 model = ['KNN', 'LR', 'SVR', 'DT', 'RF', 'GB', 'XGB']
2
3 for model_name, mae, mse, rmse, r2 in zip(model, MAE, MSE, RMSE, R2):
4     print('{} MAE is {}, MSE is {}, RMSE is {}, R2 Score is {} '.format(model_name, mae, mse
5     print('*'*120)

KNN MAE is 0.03991594081086272, MSE is 0.0042570827170690315, RMSE is 0.1997897415055706
*****
LR MAE is 0.044869966903444185, MSE is 0.00422960827789803, RMSE is 0.21182532167671594,
*****
SVR MAE is 0.06227702469271442, MSE is 0.0054030433318580555, RMSE is 0.2495536509304456
*****
DT MAE is 0.0308598733433046, MSE is 0.005161709599147438, RMSE is 0.1756697849469413, F
*****
RF MAE is 0.02609631532287773, MSE is 0.002724331998128325, RMSE is 0.16154354002211826,
*****
GB MAE is 0.0284549828004946, MSE is 0.002611718785359302, RMSE is 0.1686860480315269, F
*****
XGB MAE is 0.02612359193914408, MSE is 0.0027299448892521166, RMSE is 0.1616279429403965
*****

```




▼ Check for RandomForestRegressor

```

1 random_forest = XGBRegressor()
2 random_forest.fit(X_train, y_train)
3
4 predicted = random_forest.predict(X_test)
5
6 evaluation_df = pd.DataFrame(predicted, columns = ['Predict'])
7 pd.concat([evaluation_df, pd.DataFrame(y_test.values, columns = ['Actual'])], axis = 1).he

```

| | Predict | Actual |  |
|---|----------|----------|---|
| 0 | 0.104898 | 0.115137 | |
| 1 | 0.092464 | 0.106983 | |
| 2 | 0.020697 | 0.020934 | |
| 3 | 0.162695 | 0.002069 | |
| 4 | 0.071119 | 0.081956 | |

▼ Using Stacking Regressor to evaluate the model

```


1 estimators = [
2     ('GB', GradientBoostingRegressor()),
3     ('XGB', XGBRegressor())
4 ]
5
6 StackingRegressor_ = StackingRegressor(estimators = estimators, final_estimator = XGBRegre
7 StackingRegressor_.fit(X_train, y_train)
8 predicted = StackingRegressor_.predict(X_test)
9
10 print('MAE is # {}'.format(mean_absolute_error(predicted, y_test)))
11 print('MSE is # {}'.format(mean_squared_error(predicted, y_test)))
12 print('RMSE is # {}'.format(np.mean(mean_squared_error(predicted, y_test))))
13 print('R2 is # {}'.format(r2_score(predicted, y_test)))
14
15 evaluation_df = pd.DataFrame(predicted, columns = ['Predict'])
16 pd.concat([evaluation_df, pd.DataFrame(y_test.values, columns = ['Actual'])], axis = 1).he

```

```

MAE is # 0.02565574788523634
MSE is # 0.0027161683929956154
RMSE is # 0.0027161683929956154
R2 is # 0.6548782872725348

```

| | Predict | Actual |  |
|---|----------|----------|---|
| 0 | 0.102465 | 0.115137 | |
| 1 | 0.098963 | 0.106983 | |
| 2 | 0.022414 | 0.020934 | |
| 3 | 0.155950 | 0.002069 | |
| 4 | 0.068526 | 0.081956 | |

▼ Use KFold - 5 cross Validation for this Selection Feature

```

1 KFold_ = KFold(n_splits = 5, shuffle = True, random_state = 42)

```



```

2
3 MAE_, MSE_, RSME_, R2_, count = [], [], [], [], 1
4
5 for train_index, test_index in KFold_.split(X):
6     print('# of Cross Validation {} is running'.format(count))
7
8     X_train, X_test = X.iloc[train_index, :], X.iloc[test_index, :]
9     y_train, y_test = y[train_index], y[test_index]
10
11     StackingRegressor_ = StackingRegressor(estimators = estimators, final_estimator = XGBReg
12     StackingRegressor_.fit(X_train, y_train)
13     predicted = StackingRegressor_.predict(X_test)
14
15     MAE_.append(mean_absolute_error(predicted, y_test))
16     MSE_.append(mean_squared_error(predicted, y_test))
17     RSME_.append(np.sqrt(mean_squared_error(predicted, y_test)))
18     R2_.append(r2_score(predicted, y_test))
19
20     count = count + 1

# of Cross Validation 1 is running
# of Cross Validation 2 is running
# of Cross Validation 3 is running
# of Cross Validation 4 is running
# of Cross Validation 5 is running

1 print('The mean MAE is = {}'.format(np.array(MAE_).mean(), '\n'))
2 print('The mean MSE is = {}'.format(np.array(MSE_).mean(), '\n'))
3 print('The mean RSME is = {}'.format(np.array(RSME_).mean(), '\n'))
4 print('The mean R2 is = {}'.format(np.array(R2_).mean(), '\n'))

The mean MAE is = 0.02549977605079265
The mean MSE is = 0.002614186959615045
The mean RSME is = 0.05109530251813684
The mean R2 is = 0.6647505097983358

```

▼ Using PCA- Principle component Analysis

```

1 # PCA_ = PCA(n_components = None)
2 # PCA_.fit_transform(X)

1 # print('The Explained variance ratio is given below.\n')
2 # np.cumsum(PCA_.explained_variance_ratio_)

1 PCA_ = PCA(n_components = 10)
2 X_trans = PCA_.fit_transform(X)

```

```

1 X_train, X_test, y_train, y_test = train_test_split(X_train, y, test_size = 0.25, random_s
2
3 print('The shape of X_train is = {}'.format(X_train.shape),'\n')
4 print('The shape of X_test is = {}'.format(X_test.shape),'\n')
5 print('The shape of y_train is = {}'.format(y_train.shape),'\n')
6 print('The shape of y_test is = {}'.format(y_test.shape),'\n')

```

The shape of X_train is = (22245, 21)

The shape of X_test is = (7415, 21)

The shape of y_train is = (22245,)

The shape of y_test is = (7415,)

```

1 models = {
2     'KNN': KNeighborsRegressor(),
3     'LR' : LinearRegression(),
4     'SVR': SVR(),
5     'DT' : DecisionTreeRegressor(),
6     'RF' : RandomForestRegressor(),
7     'GB' : GradientBoostingRegressor(),
8     'XGB': XGBRegressor()
9 }
10
11 MAE, MSE, RMSE, R2 = [], [], [], []
12
13 for model_name, model_object in models.items():
14     model = model_object
15     model.fit(X_train, y_train)
16
17     predicted = model.predict(X_test)
18
19     MAE.append(mean_absolute_error(predicted, y_test))
20     MSE.append(mean_squared_error(predicted, y_test))
21     RMSE.append(np.sqrt(mean_squared_error(predicted, y_test)))
22     R2.append(r2_score(predicted, y_test))

```

```

1 model = ['KNN', 'LR', 'SVR', 'DT', 'RF', 'GB', 'XGB']
2
3 for model_name, mae, mse, rmse, r2 in zip(model, MAE, MSE, RMSE, R2):
4     print('{} MAE is {}, MSE is {}, RMSE is {}, R2 Score is {}'.format(model_name, mae, mse
5     print('*'*120)

```

KNN MAE is 0.0638534407305653, MSE is 0.007184624663268531, RMSE is 0.25269238360220775,

 LR MAE is 0.044799030605917264, MSE is 0.004227066493187741, RMSE is 0.2116578148945067,
