

PREDICTING HOUSE PRICES USING ML

Development part 2

Name: Ragulan .M

NMID : au513521106028

CLG: AMCET

PROJECT: PHASE 4



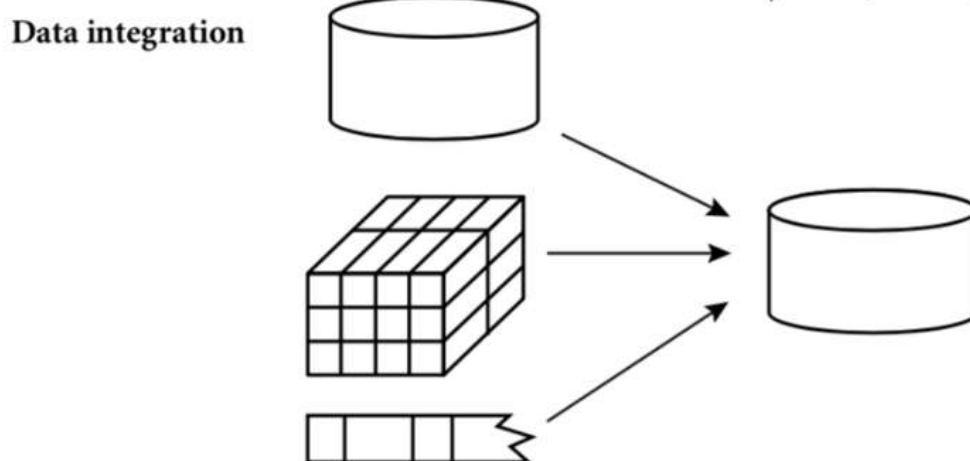
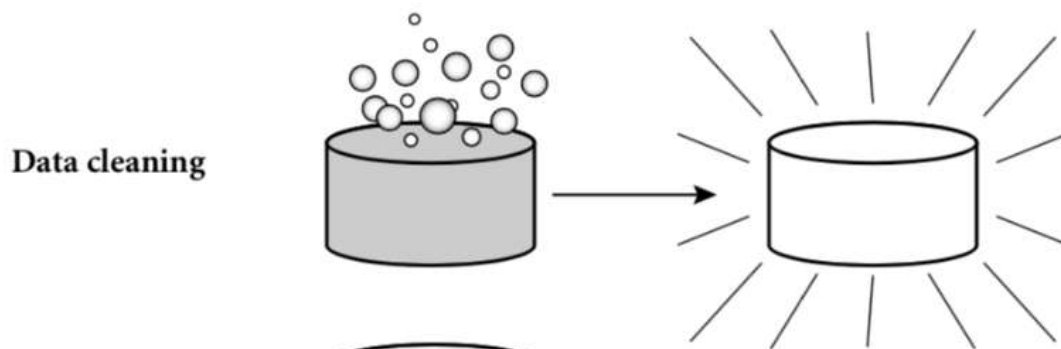
- **Data cleaning** can be applied to filling in missing values, remove noise, resolving inconsistencies, identifying and removing outliers in the data.
- **Data integration** merges data from multiple sources into a coherent data store, such as a data warehouse.
- **Data transformations**, such as normalization, may be applied. For example, normalization may improve the accuracy and efficiency of mining algorithms involving distance measurements.
- **Data reduction** can reduce the data size by eliminating redundant features, or clustering, for instance.

Reference: Data Mining: Concepts and Techniques
Second Edition, Jiawei Han, Micheline Kamber.

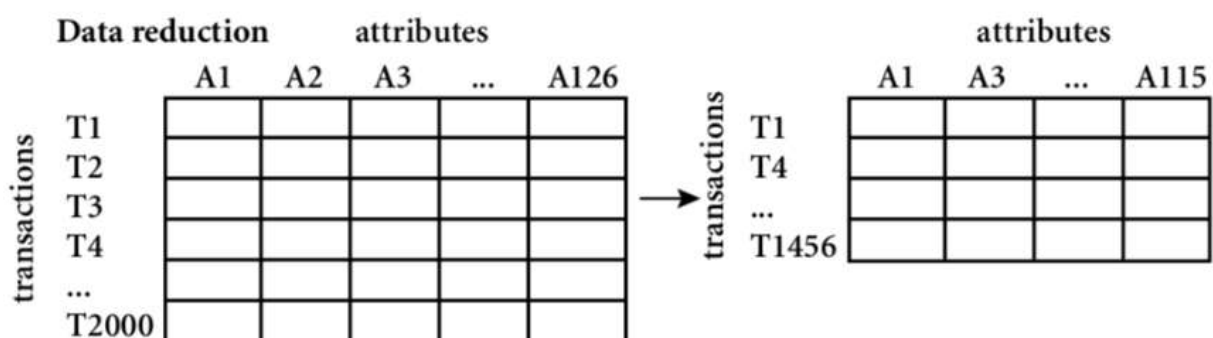
PS: This is my first kaggle notebook contribution.
Hope you like it!!

There are a number of data preprocessing techniques available such as,

1. **Data Cleaning**
2. **Data Integration**
3. **Data Transformation**
4. **Data Reduction**



Data transformation $-2, 32, 100, 59, 48 \longrightarrow -0.02, 0.32, 1.00, 0.59, 0.48$



```
def find_missing_percent(data):
    """
    Returns dataframe containing the total missing values and percentage of total missing values of a column.
    """
    miss_df = pd.DataFrame({'ColumnName':[], 'TotalMissingVals':[], 'PercentMissing':[]})
    for col in data.columns:
        sum_miss_val = data[col].isnull().sum()
        percent_miss_val = round((sum_miss_val/data.shape[0])*100,2)
        miss_df = miss_df.append(dict(zip(miss_df.columns, [col, sum_miss_val, percent_miss_val])), ignore_index=True)
    return miss_df
```



```
import warnings
warnings.filterwarnings('ignore')
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from operator import itemgetter
from sklearn.experimental import enable_iterative_imputer
from sklearn.impute import IterativeImputer
from sklearn.preprocessing import OrdinalEncoder
from category_encoders.target_encoder import TargetEncoder
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import (GradientBoostingRegressor, GradientBoostingClassifier)
import xgboost
```

	ColumnName	TotalMissingVals	PercentMissing
3	LotFrontage	259.0	17.74
6	Alley	1369.0	93.77
25	MasVnrType	8.0	0.55
26	MasVnrArea	8.0	0.55
30	BsmtQual	37.0	2.53
31	BsmtCond	37.0	2.53
32	BsmtExposure	38.0	2.60
33	BsmtFinType1	37.0	2.53
35	BsmtFinType2	38.0	2.60
42	Electrical	1.0	0.07
57	FireplaceQu	690.0	47.26
58	GarageType	81.0	5.55
59	GarageYrBlt	81.0	5.55
60	GarageFinish	81.0	5.55
63	GarageQual	81.0	5.55
64	GarageCond	81.0	5.55
72	PoolQC	1453.0	99.52
73	Fence	1179.0	80.75
74	MiscFeature	1406.0	96.30

MISSING VALUES

```
miss_df = find_missing_percent(train)
'''Displays columns with missing values'''
display(miss_df[miss_df['PercentMissing']>0.0])
print("\n")
print(f"Number of columns with missing values:{str(miss_df[miss_df['PercentMissing']>0.0].shape[0])}")
```


1.2 Drop the columns which have more than 70% of missing values

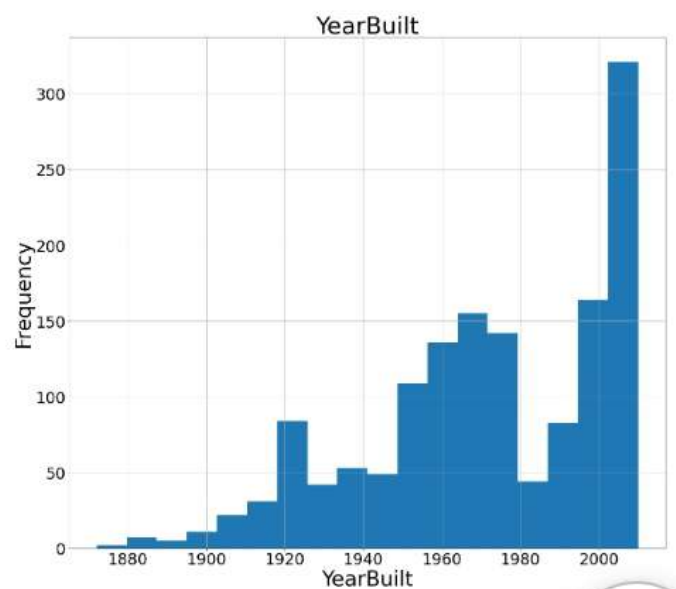
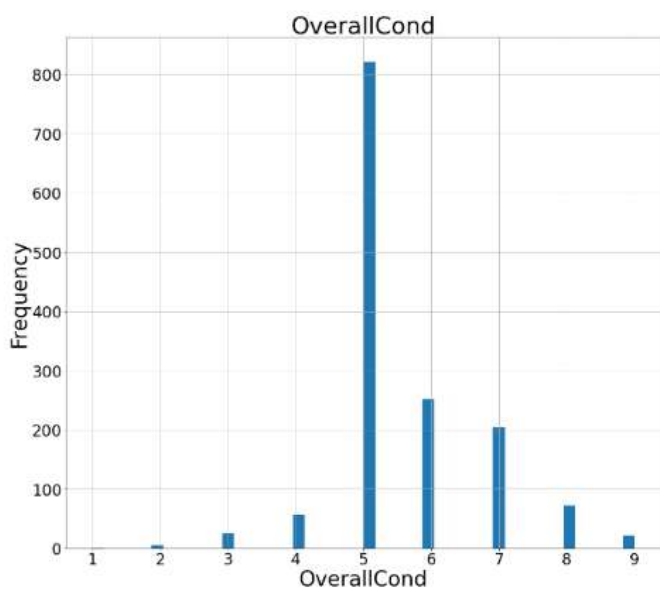
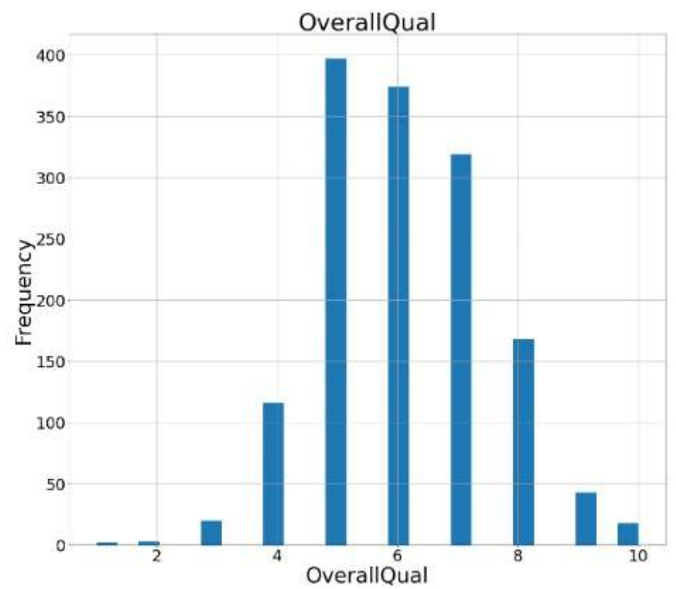
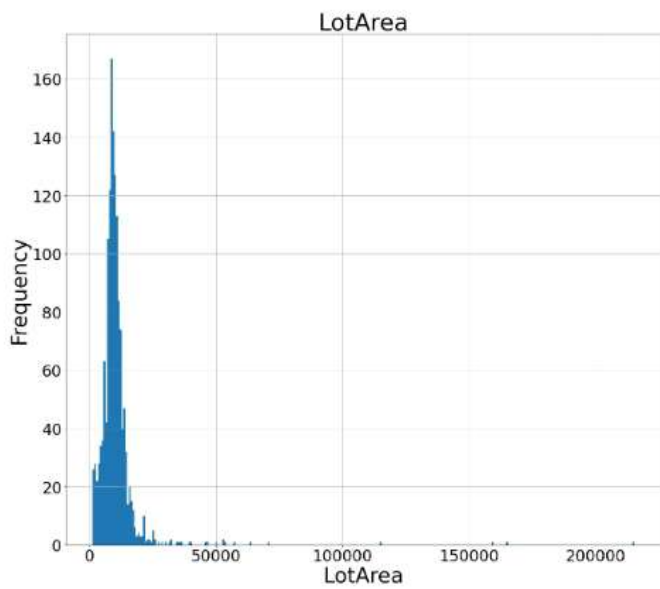
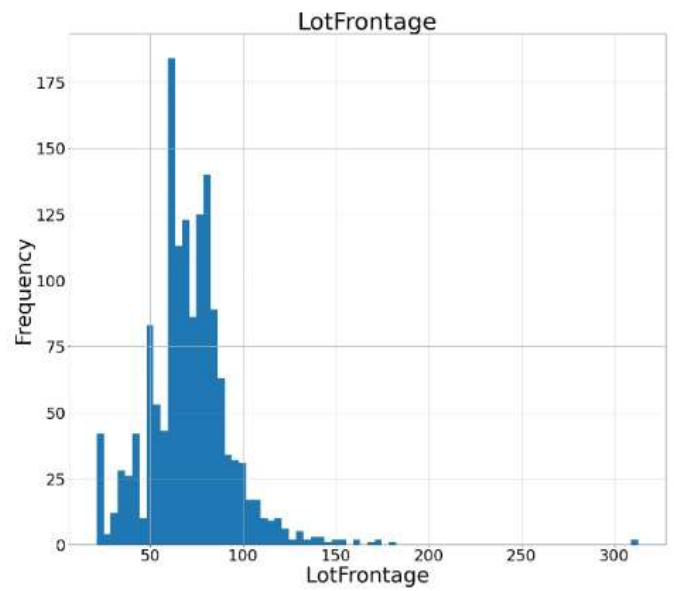
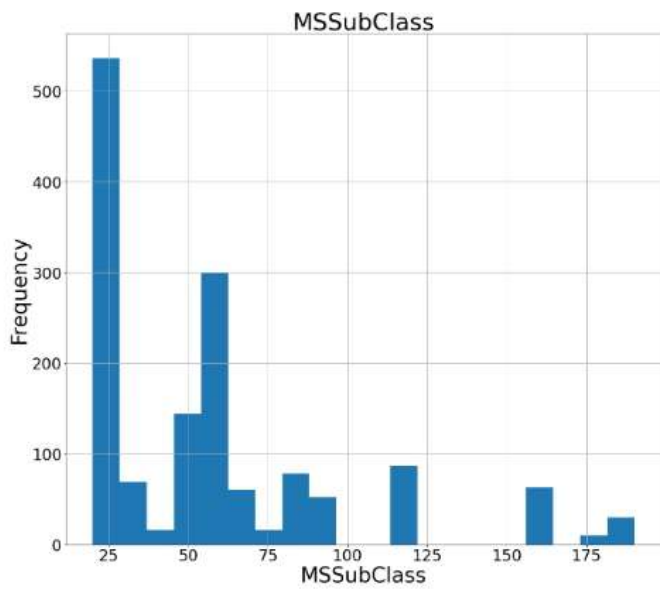
In [5]:

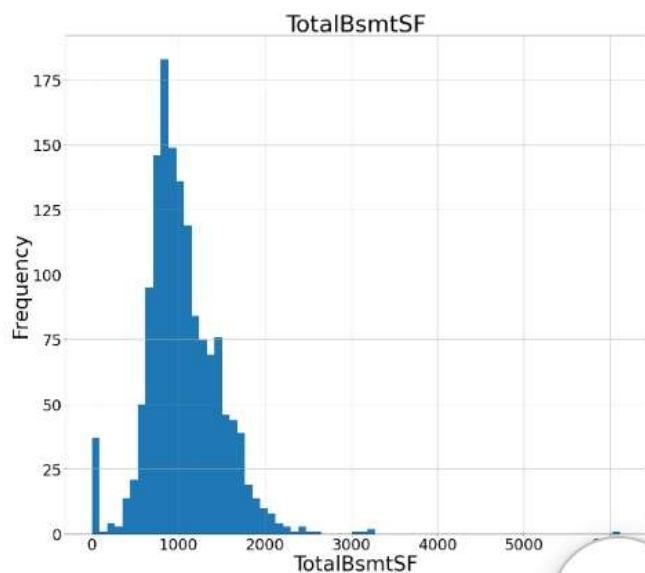
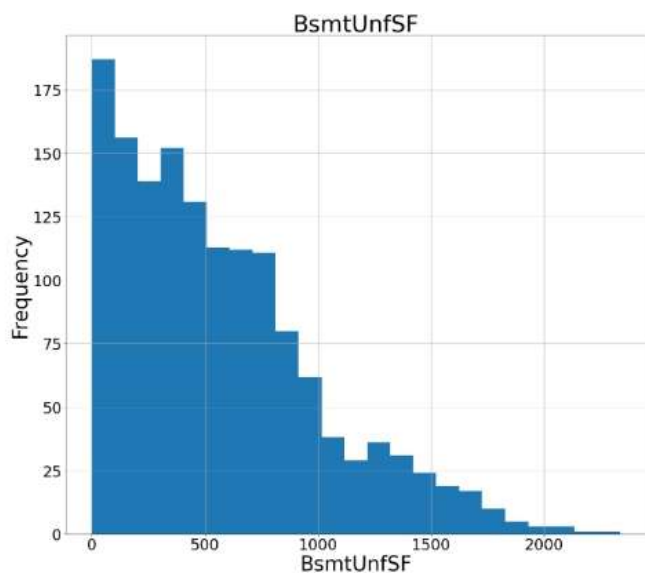
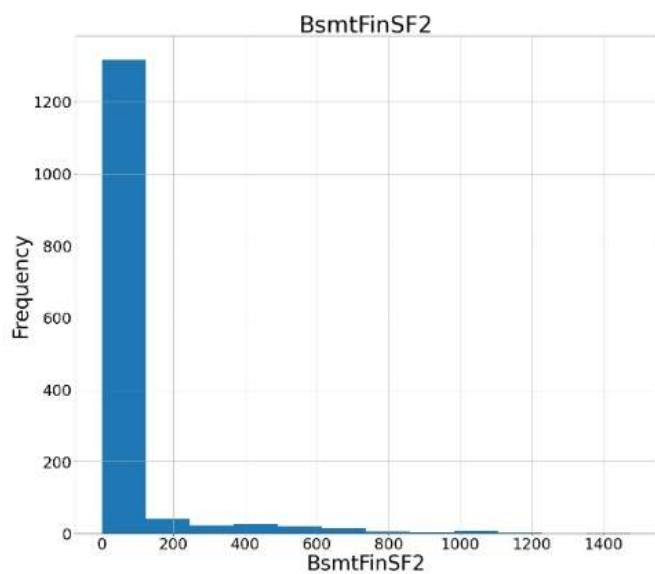
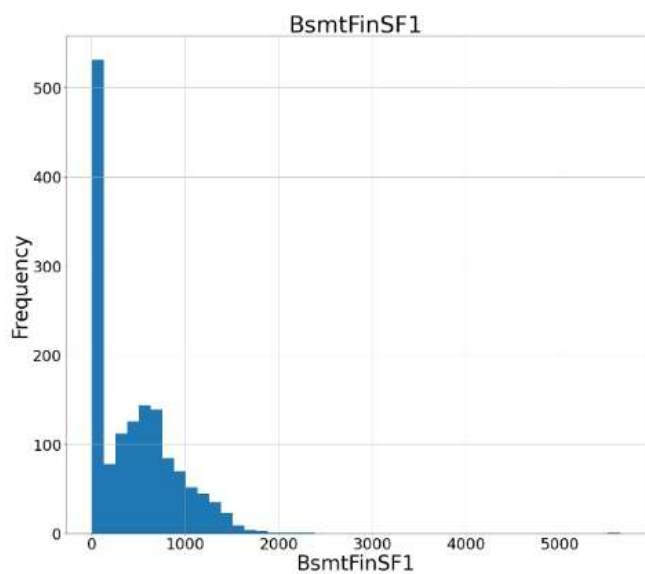
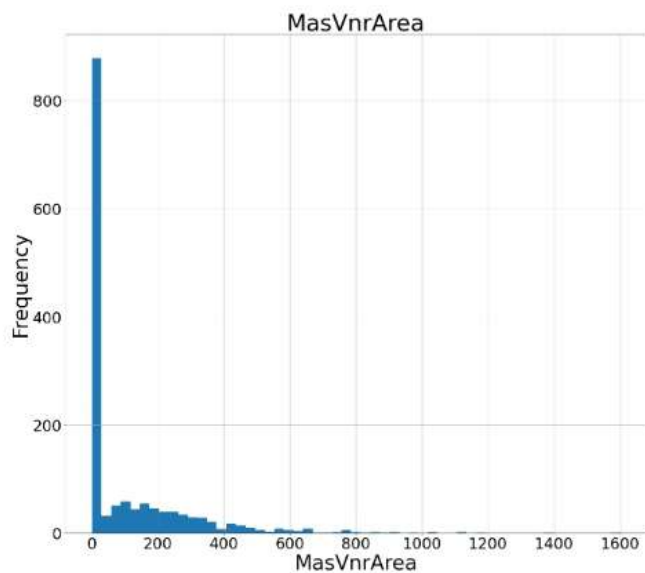
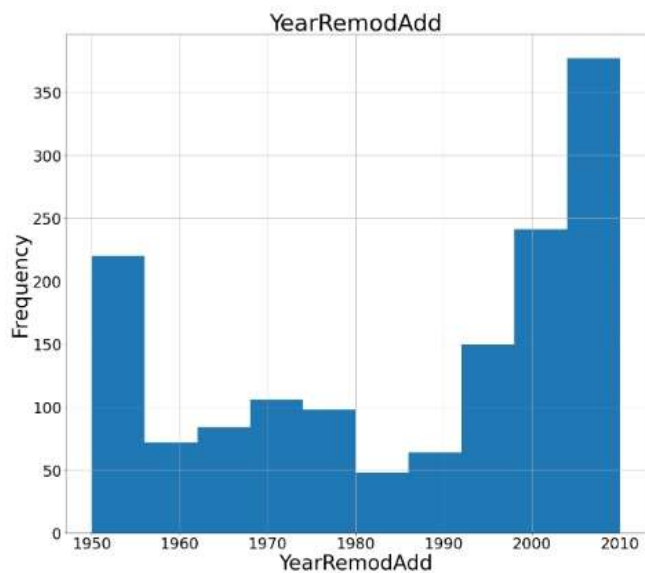
```
drop_cols = miss_df[miss_df['PercentMissing'] > 70.0].ColumnName.tolist()
print(f"Number of columns with more than 70%: {len(drop_cols)}")
train = train.drop(drop_cols, axis=1)
test = test.drop(drop_cols, axis=1)

miss_df = miss_df[miss_df['ColumnName'].isin(train.columns)]
'''Columns to Impute'''
impute_cols = miss_df[miss_df['TotalMissingVals'] > 0.0].ColumnName.tolist()
miss_df[miss_df['TotalMissingVals'] > 0.0]
```

Number of columns with more than 70%:

	ColumnName	TotalMissingVals	PercentMissing
3	LotFrontage	259.0	17.74
25	MasVnrType	8.0	0.55
26	MasVnrArea	8.0	0.55
30	BsmtQual	37.0	2.53
31	BsmtCond	37.0	2.53
32	BsmtExposure	38.0	2.60
33	BsmtFinType1	37.0	2.53
35	BsmtFinType2	38.0	2.60
42	Electrical	1.0	0.07
57	FireplaceQu	690.0	47.26
58	GarageType	81.0	5.55
59	GarageYrBlt	81.0	5.55
60	GarageFinish	81.0	5.55
63	GarageQual	81.0	5.55
64	GarageCond	81.0	5.55





DATASET

Here we have web scrapped the Data from 99acres.com website which is one of the leading real estate websites operating in INDIA.

Our Data contains Bombay Houses only.

Dataset looks as follows-

	Price	PricePerSqft	Area_Sqm	Location	Bedrooms	Latitude	Longitude	PricePerSqM
0	13300000	16625	74.32	Kandivali (East)	2	19.210200	72.864891	178885.00
1	9000000	15666	55.74	Ramgad Nagar	1	19.167700	72.949300	168566.16
2	9000000	19148	43.66	Mahakali Caves	1	19.130609	72.873816	206032.48
3	9000000	10588	78.97	Louis Wadi	2	19.126005	72.825052	113926.88
4	100000000	20000	464.51	Barrister Nath Pai Nagar	5	19.075014	72.907571	215200.00

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 840 entries, 0 to 839
Data columns (total 6 columns):
Price                840 non-null int64
Area_Sqm             840 non-null float64
Bedrooms             840 non-null int64
Latitude             840 non-null float64
Longitude            840 non-null float64
PricePerSqM          840 non-null float64
dtypes: float64(4), int64(2)
memory usage: 39.5 KB
```



```

def fit_model(x_train,y_train, model):
    """
    Fits x_train to y_train for the given
    model.
    """
    model.fit(x_train,y_train)
    return model

'''Xtreme Gradient Boosting Regressor'''
model = xgboost.XGBRegressor(objective
="reg:squarederror", random_state=42)
model = fit_model(x_train,y_train, model)

'''Predict the outcomes'''
predictions = model.predict(test)

```



3. Multivariable Analysis

Let's check out all the variables! There are two types of features in housing data, categorical and numerical.

Categorical data is just like it sounds. It is in categories. It isn't necessarily linear, but it follows some kind of pattern. For example, take a feature of "Downtown". The response is either "Near", "Far", "Yes", and "No". Back then, living in downtown usually meant that you couldn't afford to live in uptown. Thus, it could be implied that downtown establishments cost less to live in. However, today, that is not the case. (Thank you, hipsters!) So we can't really establish any particular order of response to be "better" or "worse" than the other.

Numerical data is data in number form. (Who could have thought!) These features are in a linear relationship with each other. For example, a 2,000 square foot place is 2 times "bigger" than a 1,000 square foot place. Plain and simple. Simple and



```
Index(['MSZoning', 'Street', 'Alley',  
      'LotShape', 'LandContour', 'Utilities',  
      'LotConfig', 'LandSlope', 'Neighborhood', 'Condition1', 'Condition2',  
      'BldgType', 'HouseStyle', 'RoofStyle', 'RoofMatl', 'Exterior1st',  
      'Exterior2nd', 'MasVnrType', 'ExterQual', 'ExterCond', 'Foundation',  
      'BsmtQual', 'BsmtCond', 'BsmtExposure', 'BsmtFinType1', 'BsmtFinType2',  
      'Heating', 'HeatingQC', 'CentralAir', 'Electrical', 'KitchenQual',  
      'Functional', 'FireplaceQu', 'GarageType', 'GarageFinish', 'GarageQual',  
      'GarageCond', 'PavedDrive', 'PoolQC', 'Fence', 'MiscFeature',  
      'SaleType', 'SaleCondition'],  
      dtype='object')
```



	Most Correlated Features
0	SalePrice
1	OverallQual
2	GrLivArea
3	GarageCars
4	GarageArea
5	TotalBsmtSF
6	1stFlrSF
7	FullBath
8	TotRmsAbvGrd
9	YearBuilt

