**EXPERIMENT NO 7** HOUSE PRICE PREDICTION **IMPORTING LIBRARIES** In [25]: import pandas as pd import numpy as np import matplotlib.pyplot as plt import seaborn as sns from sklearn.preprocessing import StandardScaler from sklearn.ensemble import RandomForestRegressorfrom sklearn.metrics import mean absolute error from sklearn.model\_selection import train\_test\_split READING THE DATASET In [3]: train = pd.read\_csv('/DOCUMENTS/COLLEGE/CLASSES/EXPERIMENT\_NO\_7/train.csv') test = pd.read\_csv('/DOCUMENTS/COLLEGE/CLASSES/EXPERIMENT\_NO\_7/test.csv') In [4]: train.head() Id MSSubClass MSZoning LotFrontage LotArea Street Alley LotShape LandContour Utilities ... PoolArea PoolQC Fence MiscFeatu Out[4]: 0 1 60 RL 65.0 8450 Pave NaN Reg AllPub 0 NaN NaN Νć 2 RL80.0 9600 0 1 20 Pave NaN Reg Lvl AllPub NaN NaN Ná AllPub 3 RL 68.0 2 60 11250 Pave NaN IR1 Lvl 0 NaN NaN Νć 3 70 RL 60.0 9550 Pave NaN IR1 Lvl AllPub NaN NaN Ná 5 RL 14260 IR1 AllPub 60 84.0 Pave NaN Lvl 0 NaN NaN Νá 5 rows × 81 columns In [5]: test.head() Utilities ... ScreenPorch PoolArea PoolQC Out[5]: Id MSSubClass MSZoning LotFrontage LotArea Street Alley LotShape LandContour **0** 1461 20 RH 80.0 11622 Pave NaN Reg Lvl AllPub 120 0 NaN N **1** 1462 20 81.0 14267 Pave NaN IR1 AllPub 0 NaN **2** 1463 60 RL 74.0 13830 Pave NaN IR1 Lvl AllPub 0 0 NaN N 78.0 **3** 1464 60 9978 Pave NaN IR1 Lvl AllPub 0 NaN 43.0 **4** 1465 120 RL 5005 Pave NaN IR1 HLS AllPub ... 144 0 NaN 5 rows × 80 columns **EXPLORATORY DATA ANALYSIS(EDA) CORRELATION MATRIX** In [11]: # correlation matrix corrmat = train.corr() f, ax = plt.subplots(figsize=(12, 9)) sns.heatmap(corrmat, vmax=.8, square=True) - 0.8 ld · MSSubClass LotFrontage LotArea OverallQual OverallCond 0.6 YearBuilt YearRemodAdd MasVnrArea BsmtFinSF1 BsmtFinSF2 BsmtUnfSF - 0.4 TotalBsmtSF 1stFIrSF 2ndFlrSF LowOualFinSF GrLivArea BsmtFullBath 0.2 BsmtHalfBath FullBath HalfBath BedroomAbvGr KitchenAbvGr TotRmsAbvGrd 0.0 Fireplaces GarageYrBlt GarageCars GarageArea WoodDeckSF OpenPorchSF EnclosedPorch 3SsnPorch ScreenPorch PoolArea MiscVal -0.4 MoSold YrSold SalePrice LowQualFinSF -GrLivArea -BsmtFullBath -BsmtHalfBath -FullBath -HalfBath -HalfBath -KitchenAbvGr -TotRmsAbvGr -EnclosedPorch -3SsnPorch -ScreenPorch -PoolArea -MiscVal -MoSold -YrSold -GarageYrBlt GarageCars GarageArea WoodDecKSF -BsmtUnfSF -TotalBsmtSF -1stFirSF -2ndFirSF -Fireplaces In [13]: # saleprice correlation matrix k = 10 # number of variables for heatmap cols = corrmat.nlargest(k, 'SalePrice')['SalePrice'].index cm = np.corrcoef(train[cols].values.T) sns.set(font scale=1.25) hm = sns.heatmap(cm, cbar=True, annot=True, square=True, fmt='.2f', annot kws={'size': 10}, yticklabels=cols.values, xticklabels=cols.values) plt.show() 1.0 1.00 0.79 0.71 0.64 0.62 0.61 0.61 0.56 0.53 0.52 SalePrice 1.00 0.59 0.60 0.56 0.54 0.48 0.55 0.43 0.57 OverallQual - 0.8 GrLivArea 0.71 0.59 <mark>1.00</mark> 0.47 0.47 0.45 0.57 0.63 <mark>0.83 </mark>0.20 0.64 0.60 0.47 <mark>1.00 0.88</mark> 0.43 0.44 0.47 0.36 0.54 GarageCars 0.62 0.56 0.47 <mark>0.88 1.00</mark> 0.49 0.49 0.41 0.34 0.48 - 0.6 GarageArea 0.61 0.54 0.45 0.43 0.49 <mark>1.00 0.82</mark> 0.32 0.29 0.39 TotalBsmtSF 1stFIrSF 0.61 0.48 0.57 0.44 0.49 0.82 1.00 0.38 0.41 0.28 - 0.4 0.56 0.55 <mark>0.63</mark> 0.47 0.41 0.32 0.38 <mark>1.00</mark> 0.55 0.47 FullBath TotRmsAbvGrd 0.53 0.43 <mark>0.83</mark> 0.36 0.34 0.29 0.41 <mark>0.55</mark> 1.00 <mark>0.10</mark> 0.2 0.52 0.57 0.20 0.54 0.48 0.39 0.28 0.47 0.10 1.00 YearBuilt FullBath GarageCars GarageArea YearBuilt 1stFIrSF TotRmsAbvGrd OverallQual GrLivArea TotalBsmtSF **DATA VISUALIZATION** In [16]: # scatterplot sns.set() 'SalePrice', 'OverallQual', 'GrLivArea', 'GarageCars', 'TotalBsmtSF', 'FullBath', 'YearBuilt' sns.pairplot(train[cols], height=2.5) 600000 400000 200000 10 OverallOua 5000 3000 2000 gä 5000 4000 3000 2000 3.0 2.5 2.0 1.5 0.5 0.0 1975 1950 1900 1875 200000 400000 600000 5.0 10.0 4000 2000 4000 6000 1900 1950 GarageCars CHECKING THE MISSING DATA In [26]: test['SalePrice'] = 0 In [18]: # missing data total = train.isnull().sum().sort values(ascending=False) percent = (train.isnull().sum() / train.isnull().count()).sort values(ascending=False) missing data = pd.concat([total, percent], axis=1, keys=['Total', 'Percent']) missing data.head(20) Out[18]: **Total Percent PoolQC** 1453 0.995205 1406 0.963014 MiscFeature Alley 1369 0.937671 Fence 1179 0.807534 FireplaceQu 690 0.472603 LotFrontage 259 0.177397 GarageYrBlt 81 0.055479 GarageCond 81 0.055479 81 0.055479 GarageType GarageFinish 81 0.055479 GarageQual 81 0.055479 BsmtFinType2 38 0.026027 **BsmtExposure** 38 0.026027 **BsmtQual** 37 0.025342 **BsmtCond** 37 0.025342 BsmtFinType1 37 0.025342 8 0.005479 MasVnrArea MasVnrType 8 0.005479 1 0.000685 **Electrical** 0.000000 ld In [27]: data = [train, test] data = pd.concat(data) In [28]: percent missing = data.isnull().sum() \* 100 / len(data) missing value = pd.DataFrame({ 'column': data.columns, 'percent missing': percent missing missing\_value.sort\_values('percent\_missing', ascending=False) missing\_value.sort\_values(by=['percent\_missing'], inplace=True, ascending=False) missing\_value.head() Out[28]: column percent\_missing **PoolQC** PoolQC 99.657417 MiscFeature MiscFeature 96.402878 Alley Alley 93.216855 80.438506 **Fence** Fence FireplaceQu FireplaceQu 48.646797 In [29]: dropCols = missing\_value["column"][missing\_value["percent\_missing"] > 47.0] dropCols PoolQC PoolQC Out[29]: MiscFeature MiscFeature Alley Alley Fence Fence FireplaceQu FireplaceQu Name: column, dtype: object In [30]: data.drop(columns=dropCols, inplace=True) In [31]: data = data.fillna(data.mean()) C:\Users\HP\AppData\Local\Temp/ipykernel 12768/476127175.py:1: FutureWarning: Dropping of nuisance columns in D ataFrame reductions (with 'numeric only=None') is deprecated; in a future version this will raise TypeError. S elect only valid columns before calling the reduction. data = data.fillna(data.mean()) We'll consider that when more than 15% of the data is missing, we should delete the corresponding variable and pretend it never existed. **Checking Outliers** In [21]: # Univariate Analysis # standardizing data saleprice scaled = StandardScaler().fit transform( train['SalePrice'][:, np.newaxis]) low range = saleprice scaled[saleprice scaled[:, 0].argsort()][:10] high\_range = saleprice\_scaled[saleprice\_scaled[:, 0].argsort()][-10:] print('outer range (low) of the distribution:') print(low\_range) print('\nouter range (high) of the distribution:') print(high\_range) outer range (low) of the distribution: [[-1.83870376] [-1.83352844][-1.80092766][-1.78329881][-1.77448439][-1.62337999][-1.61708398][-1.58560389][-1.58560389][-1.5731outer range (high) of the distribution: [[3.82897043] [4.04098249] [4.49634819] [4.71041276] [4.73032076] [5.06214602] [5.42383959] [5.59185509] [7.10289909] [7.22881942]] C:\Users\HP\AppData\Local\Temp/ipykernel\_12768/1909256962.py:4: FutureWarning: Support for multi-dimensional in dexing (e.g. `obj[:, None]`) is deprecated and will be removed in a future version. Convert to a numpy array b efore indexing instead. saleprice scaled = StandardScaler().fit transform(train['SalePrice'][:,np.newaxis]); In [23]: # Bivariate Analysis # bivariate analysis saleprice var = 'TotalBsmtSF' data = pd.concat([train['SalePrice'], train[var]], axis=1) data.plot.scatter(x=var, y='SalePrice', ylim=(0, 800000)) \*c\* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will ha ve precedence in case its length matches with \*x\* & \*y\*. Please use the \*color\* keyword-argument or provide a 2D array with a single row if you intend to specify the same RGB or RGBA value for all points. 700000 600000 500000 400000 300000 200000 100000 3000 TotalBsmtSF **MODEL BUILDING** In [61]: for i in data.columns: if (data[i].dtype == '0'): data.drop(columns=[i], inplace=True) data.drop(columns=['Id'], inplace=True) In [73]: X\_train, X\_test, y\_train, y\_test = train\_test\_split( data.drop(columns=['SalePrice']), data['SalePrice'], test\_size=.2) Multiple Linear Regression In [43]: from sklearn.linear\_model import LinearRegression multi\_reg = LinearRegression() multi\_reg.fit(x\_train, y\_train) print('Training Accuracy : ', multi\_reg.score(x\_train, y\_train) \* 100, '%') Training Accuracy: 83.98680290654255 % **Lasso Regression** In [43]: from sklearn.linear\_model import LinearRegression multi\_reg = LinearRegression() multi\_reg.fit(x\_train, y\_train) print('Training Accuracy : ', multi\_reg.score(x\_train, y\_train) \* 100, '%') Training Accuracy: 83.98680290654255 % Random Forest Regressor In [43]: from sklearn.linear\_model import LinearRegression multi\_reg = LinearRegression() multi\_reg.fit(x\_train, y\_train) print('Training Accuracy : ', multi\_reg.score(x\_train, y\_train) \* 100, '%') Training Accuracy : 83.98680290654255 % X----X----X----X----