

```
# This Python 3 environment comes with many helpful analytics libraries  
# It is defined by the kaggle/python docker image: https://github.com/kaggle/docker-python  
# For example, here's several helpful packages to load in
```

```
import numpy as np # linear algebra  
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
```

```
# Input data files are available in the "../input/" directory.  
# For example, running this (by clicking run or pressing Shift+Enter) will list the files in the input directory
```

```
import os  
for dirname, _, filenames in os.walk('/kaggle/input'):  
    for filename in filenames:  
        print(os.path.join(dirname, filename))
```

```
# Any results you write to the current directory are saved as output files
```



```
/kaggle/input/house-prices-advanced-regression-techniques/train.csv  
/kaggle/input/house-prices-advanced-regression-techniques/test.csv  
/kaggle/input/house-prices-advanced-regression-techniques/train.csv  
/kaggle/input/house-prices-advanced-regression-techniques/test.csv
```

▼ Homework 3 - Ames Housing Dataset

For all parts below, answer all parts as shown in the Google document for Homework 3. You can find the document here: <https://colab.research.google.com/drive/1LLI98XWtN1zqx8uZusjASoqDjxdCHhVg?authuser=1#scrollTo=8xEjXVNqily6&printMode=true>. We also ask that you provide code as well as text to answer the questions. We also ask that code be commented.

▼ Part 1 - Pairwise Correlations

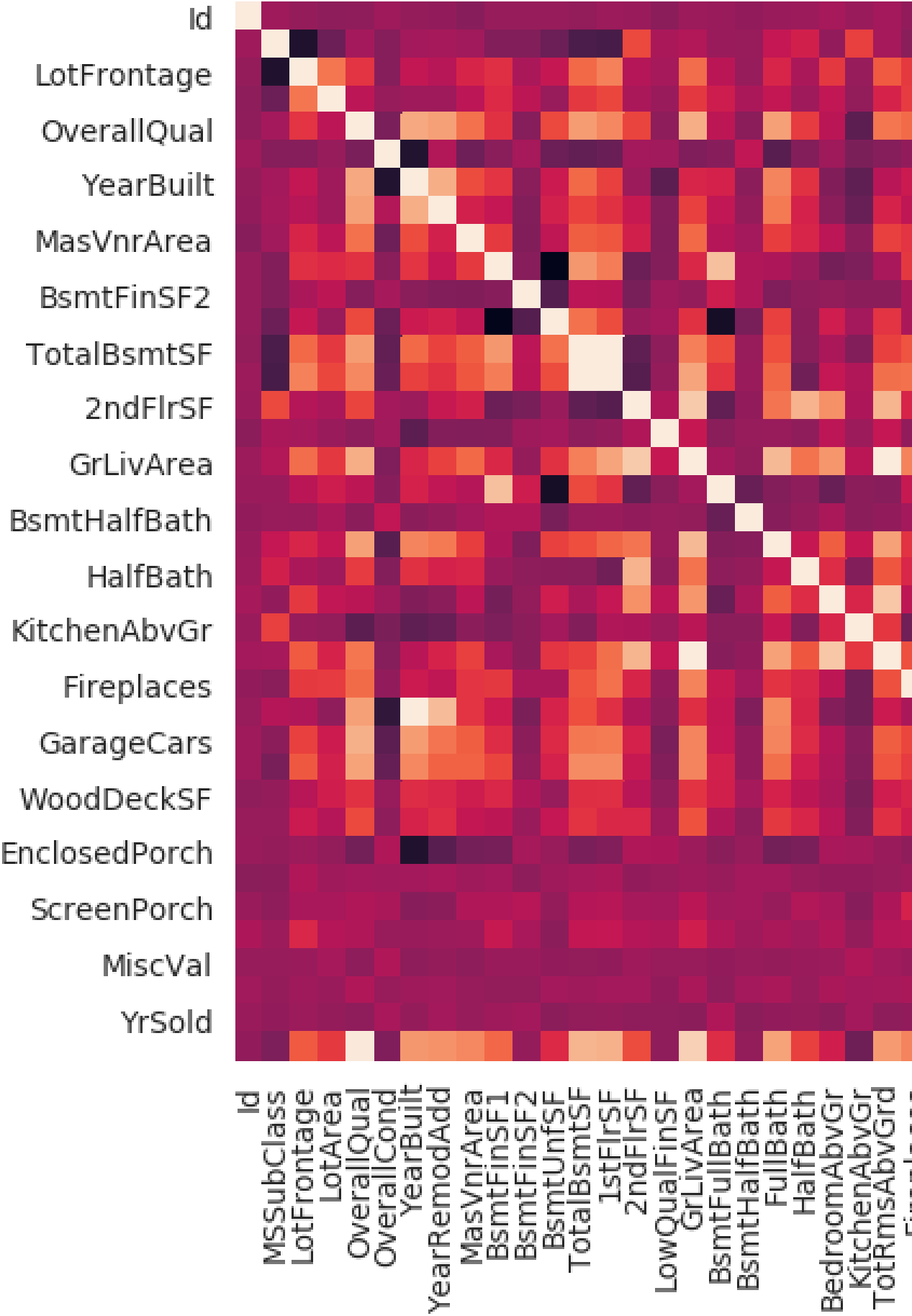
```
# TODO: show visualization
```

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import numpy as np
```

```
df1 = pd.read_csv('/kaggle/input/house-prices-advanced-regression-
```

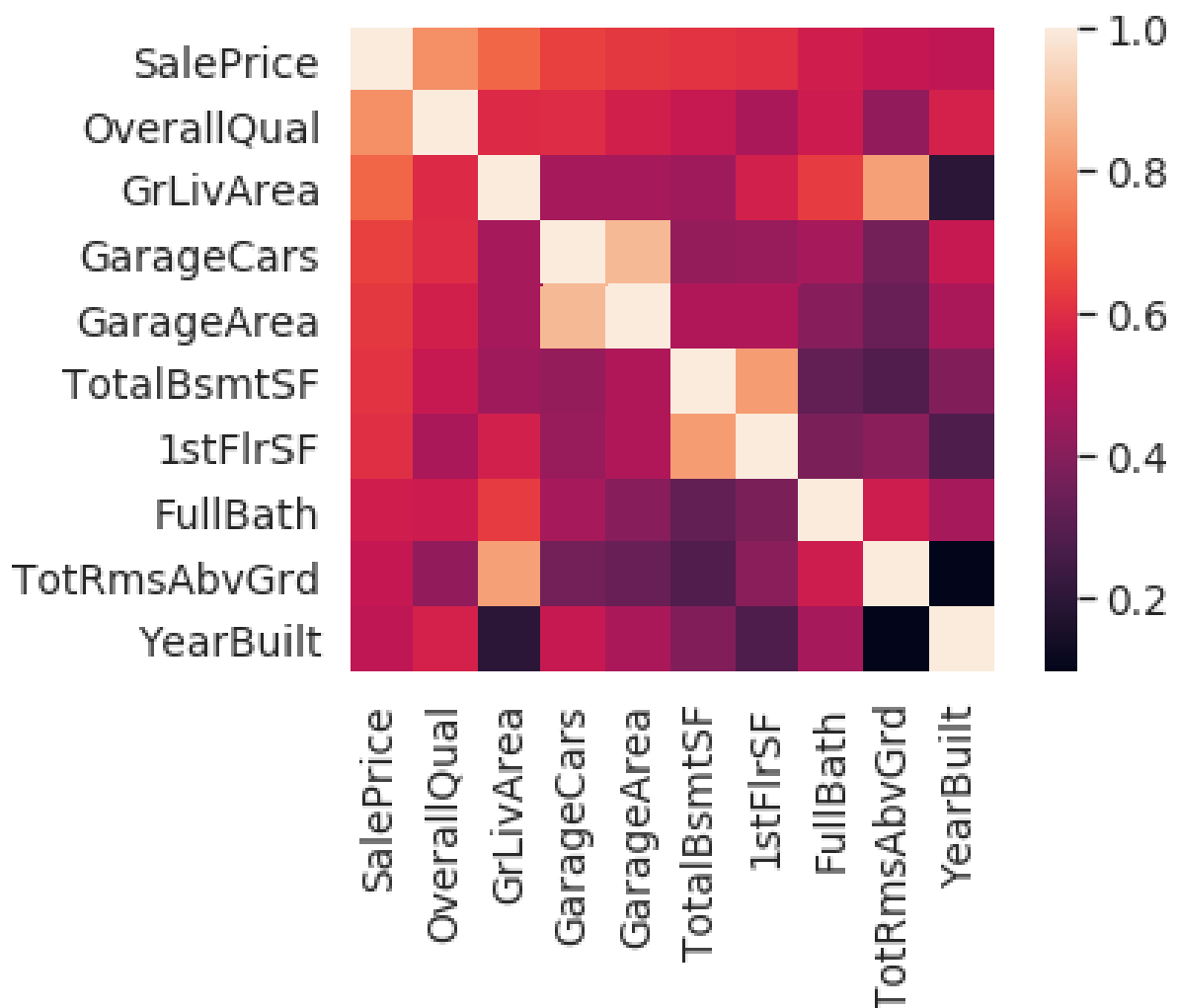
```
corrmat = df1.corr(method = 'pearson')
f, ax = plt.subplots(figsize=(12, 9))
sns.heatmap(corrmat, vmax=.8, square=True);
```





This heatmap shows an overview of all the values and how do are they corre of what is happening, it does give us an overall view and we can see some p Enclosed Porch are inversely correlated. This might be due to the increase i allowed for each home to be made gets lesser and hence the dip. We also s Ground Live Area are highly correlated.

```
#saleprice correlation matrix
k = 10 #number of variables for heatmap
cols = corrmatrix.nlargest(k, 'SalePrice')['SalePrice'].index
cm = np.corrcoef(df1[cols].values.T)
sns.set(font_scale=1.25)
hm = sns.heatmap(cm, cbar=True, square=True, fmt='.2f', yticklabel:
plt.show()
```



This heatmap provides a much thorough view of the 10 most correlated variables. Variables like YearBuilt and OverallQual built have a negative correlation whereas GarageCars i.e. the number of cars that can fit in the total area of the garage are highly correlated.

```
print (corrmat['SalePrice'].sort_values(ascending=False)[:10], '\n')
print ('-----')
print (corrmat['SalePrice'].sort_values(ascending=False)[-10:]) #
```



```
SalePrice      1.000000
OverallQual    0.790982
GrLivArea      0.708624
GarageCars     0.640409
GarageArea     0.623431
TotalBsmtSF    0.613581
1stFlrSF       0.605852
FullBath       0.560664
TotRmsAbvGrd   0.533723
YearBuilt      0.522897
Name: SalePrice, dtype: float64
```

```
-----
BsmtFinSF2     -0.011378
BsmtHalfBath   -0.016844
MiscVal        -0.021190
Id             -0.021917
LowQualFinSF   -0.025606
YrSold         -0.028923
OverallCond    -0.077856
MSSubClass     -0.084284
EnclosedPorch  -0.128578
KitchenAbvGr   -0.135907
Name: SalePrice, dtype: float64
```

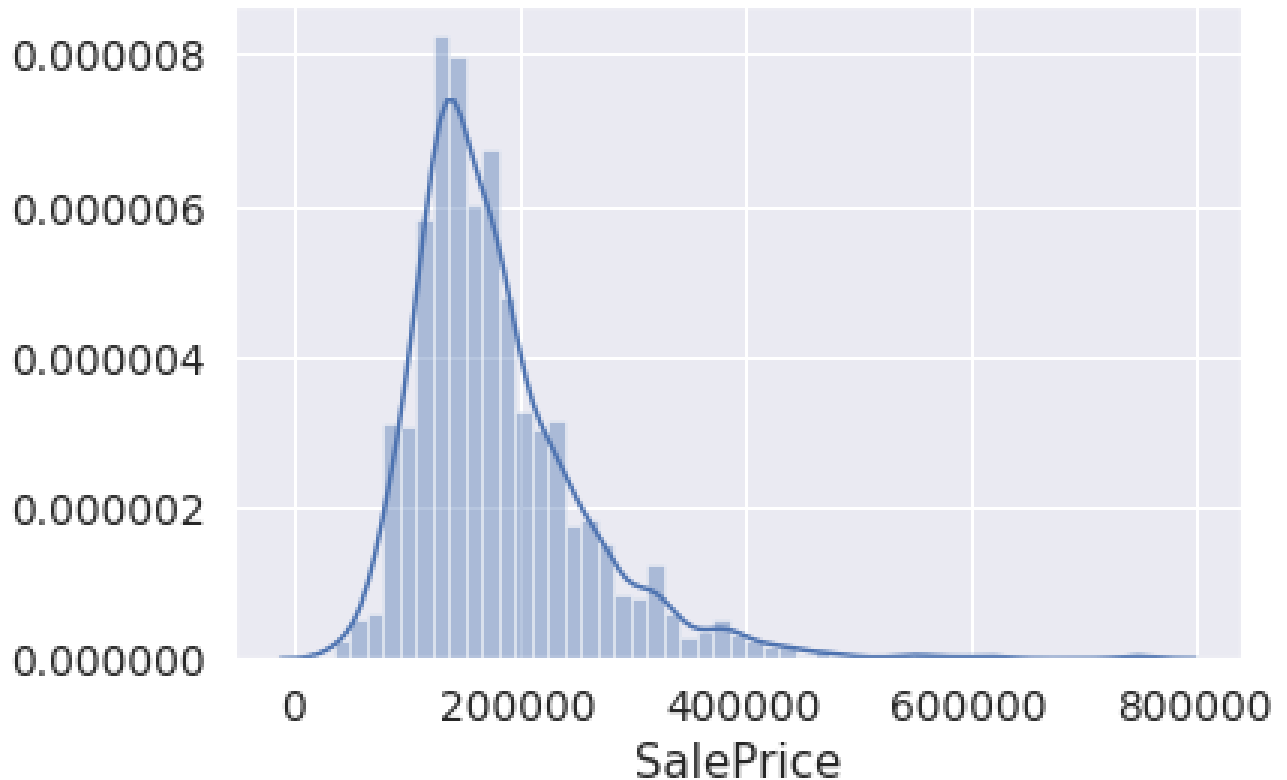
Here, we want to show the negative correlation with our target variable which is SalePrice. Variables like YearBuilt and OverallQual have a negative correlation with SalePrice, which is relevant as the decreased porch means that the area is crowded one hence

▼ Part 2 - Informative Plots

```
# TODO: code to generate Plot 1  
sns.distplot(df1['SalePrice'])
```



<matplotlib.axes._subplots.AxesSubplot at 0x7f547815c940

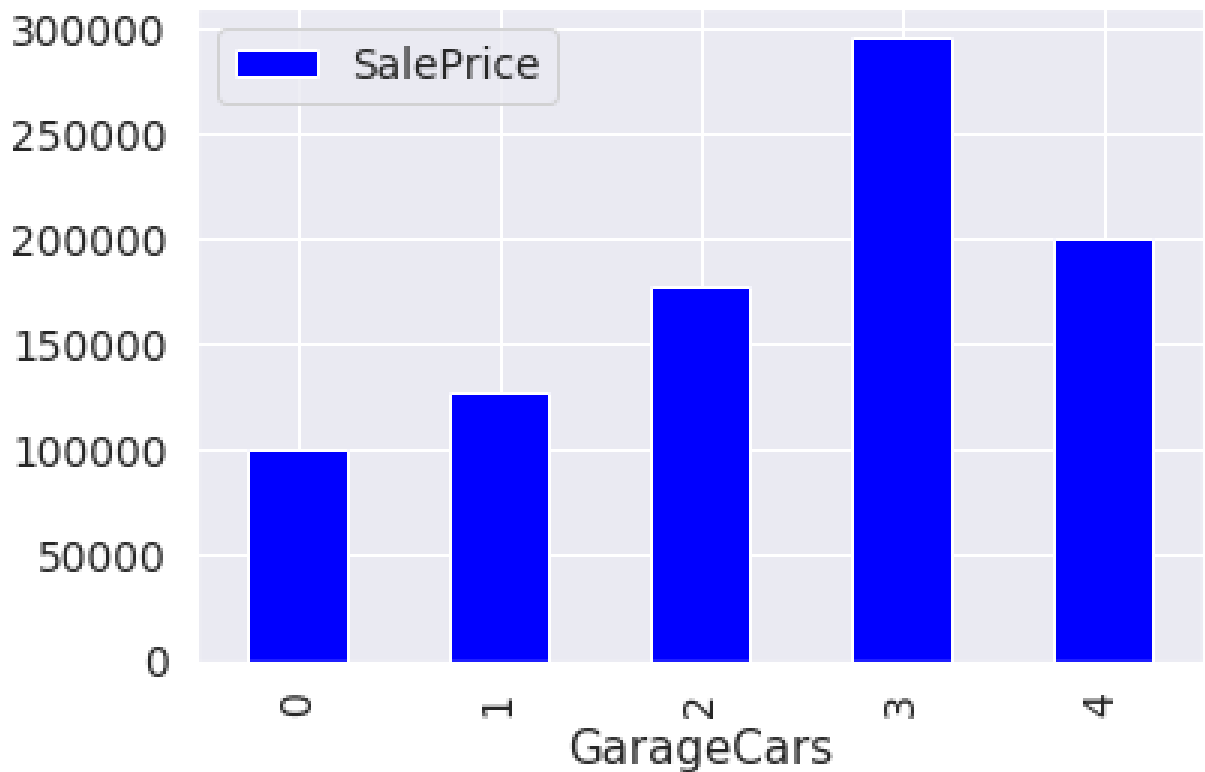


This plot shows us the SalePrice which is skewed. This means that the norm much cheaper than anywhere here in NY.

```
# TODO: code to generate Plot 2  
pivot = df1.pivot_table(index='GarageCars', values='SalePrice', ag  
pivot.plot(kind='bar', color='blue')
```



<matplotlib.axes._subplots.AxesSubplot at 0x7f5471c8a8d0



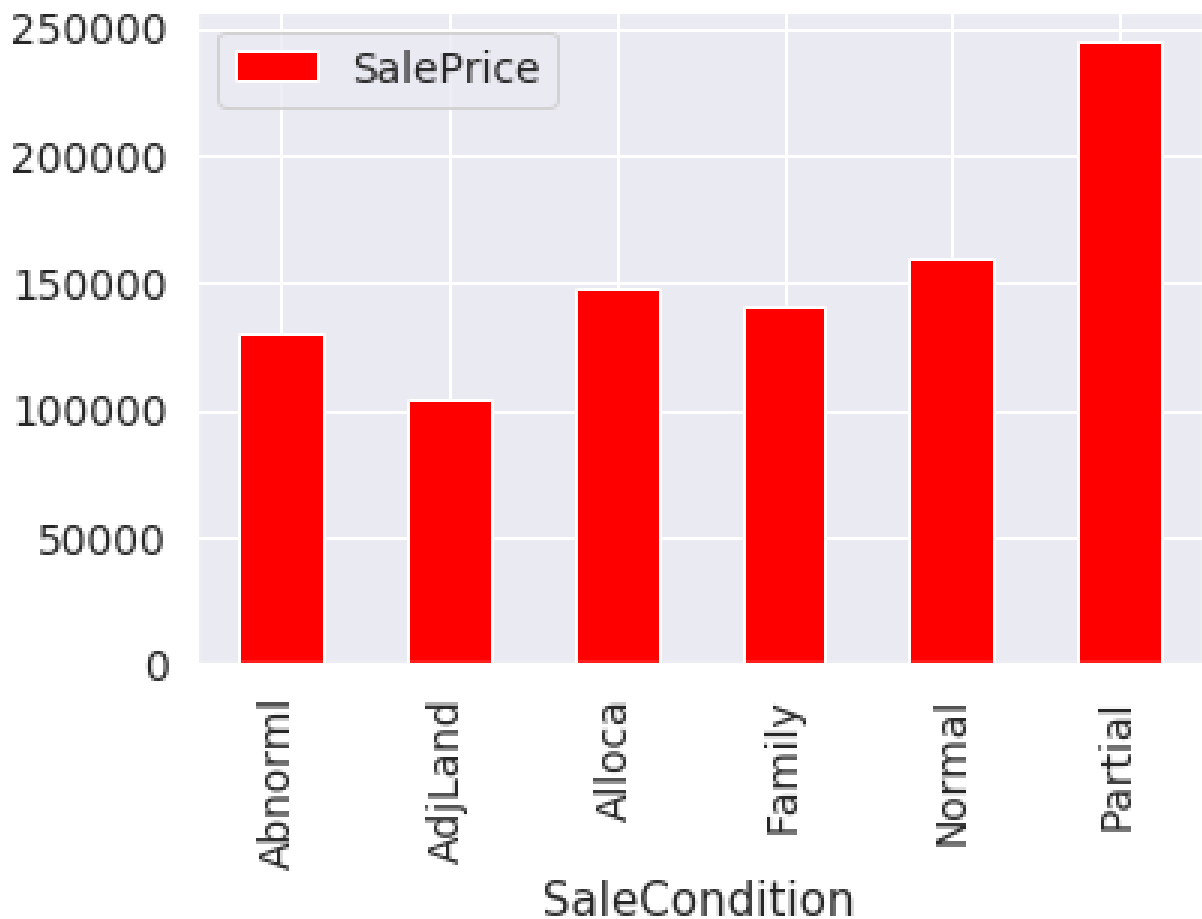
This plot tells us that a house which has place for 3 cars cost more than a p due to the fact that when the capacity of the car space increases, the living : decrease, making the place go down in price.

TODO: code to generate Plot 3

```
sp_pivot = df1.pivot_table(index='SaleCondition', values='SalePrice')
sp_pivot.plot(kind='bar', color='red')
```



<matplotlib.axes._subplots.AxesSubplot at 0x7f54710cfeb8

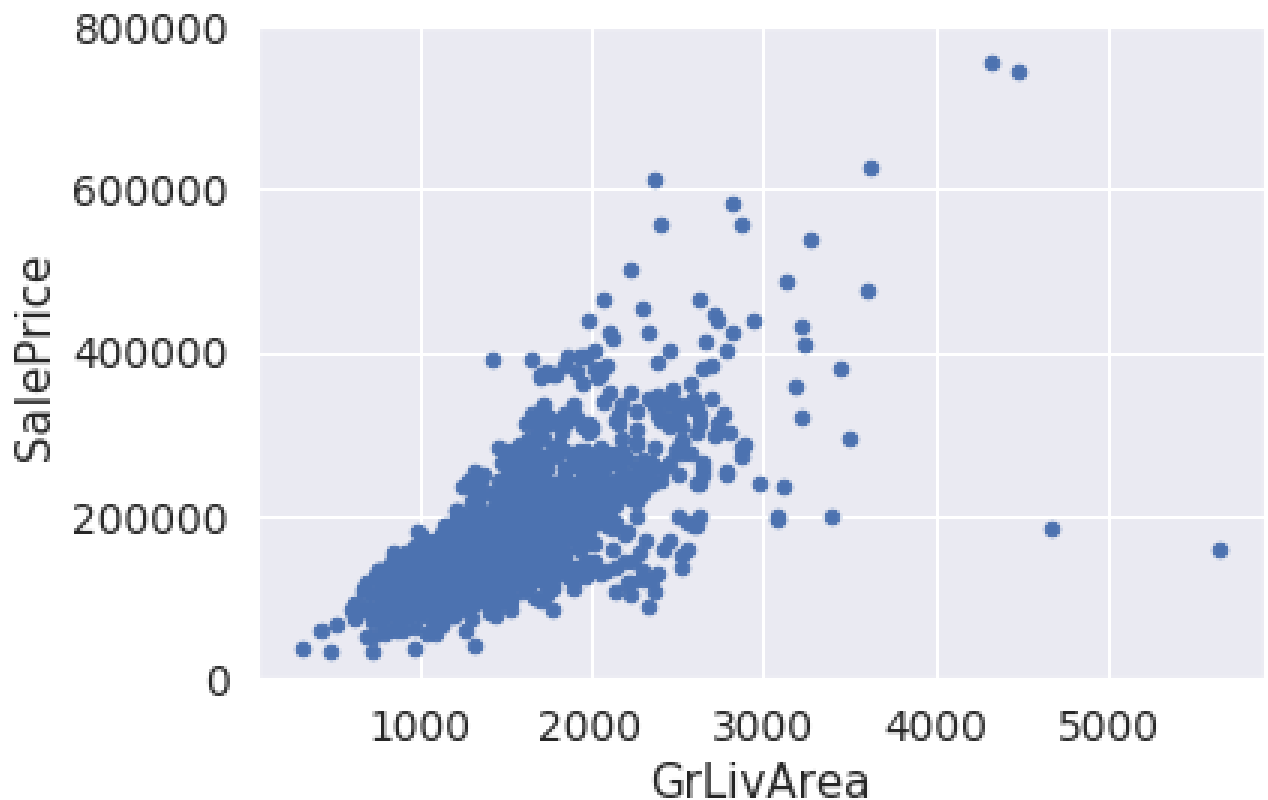


It seems that people love Partial housing. They dont want the whole thing.

TODO: code to generate Plot 4

```
data = pd.concat([df1['SalePrice'], df1['GrLivArea']], axis=1)  
data.plot.scatter(x='GrLivArea', y='SalePrice', ylim=(0,800000));
```



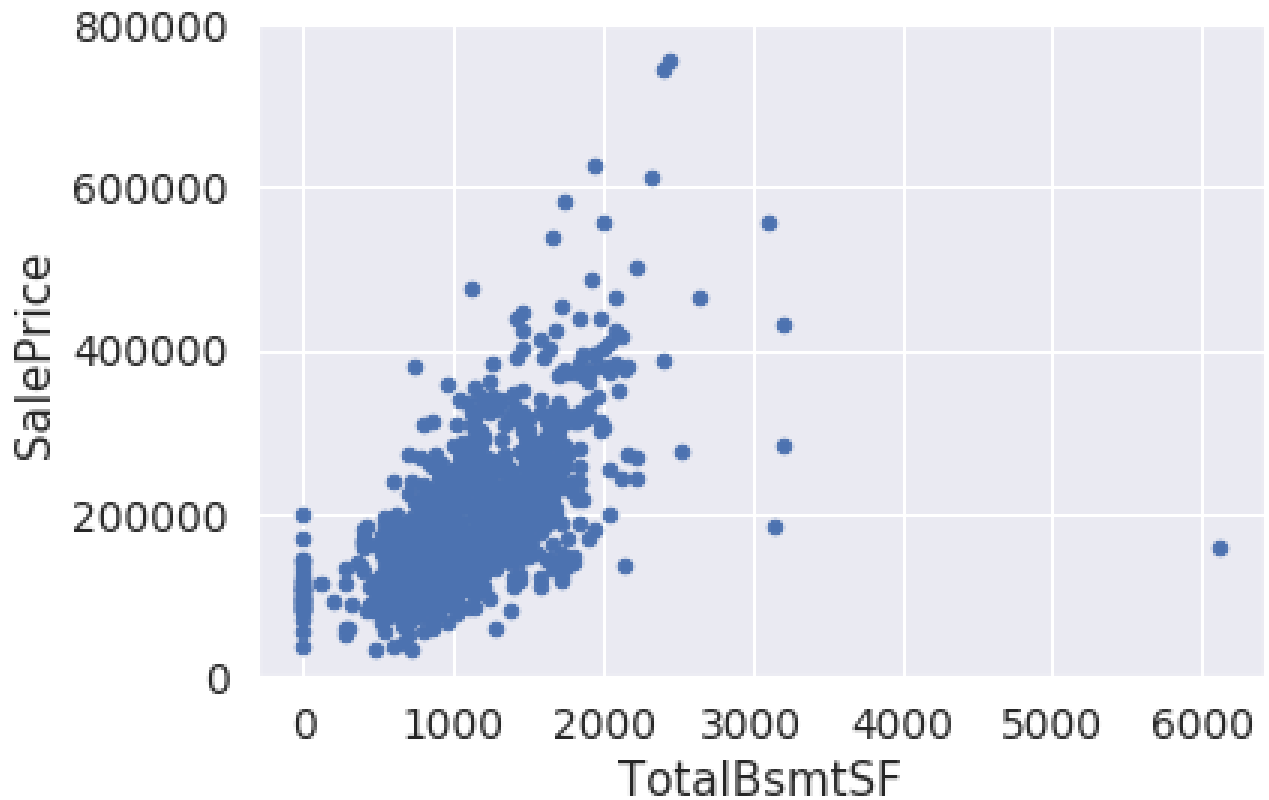


This graph shows that GrLivArea, i.e. Ground Live Area is highly correlated to actually very self explanatory. But, we also see some outliers in the data. We which seems odd and might be an outlier or faulty entry.

TODO: code to generate Plot 5

```
data = pd.concat([df1['SalePrice'], df1['TotalBsmtSF']], axis=1)
data.plot.scatter(x='TotalBsmtSF', y='SalePrice', ylim=(0,800000))
```

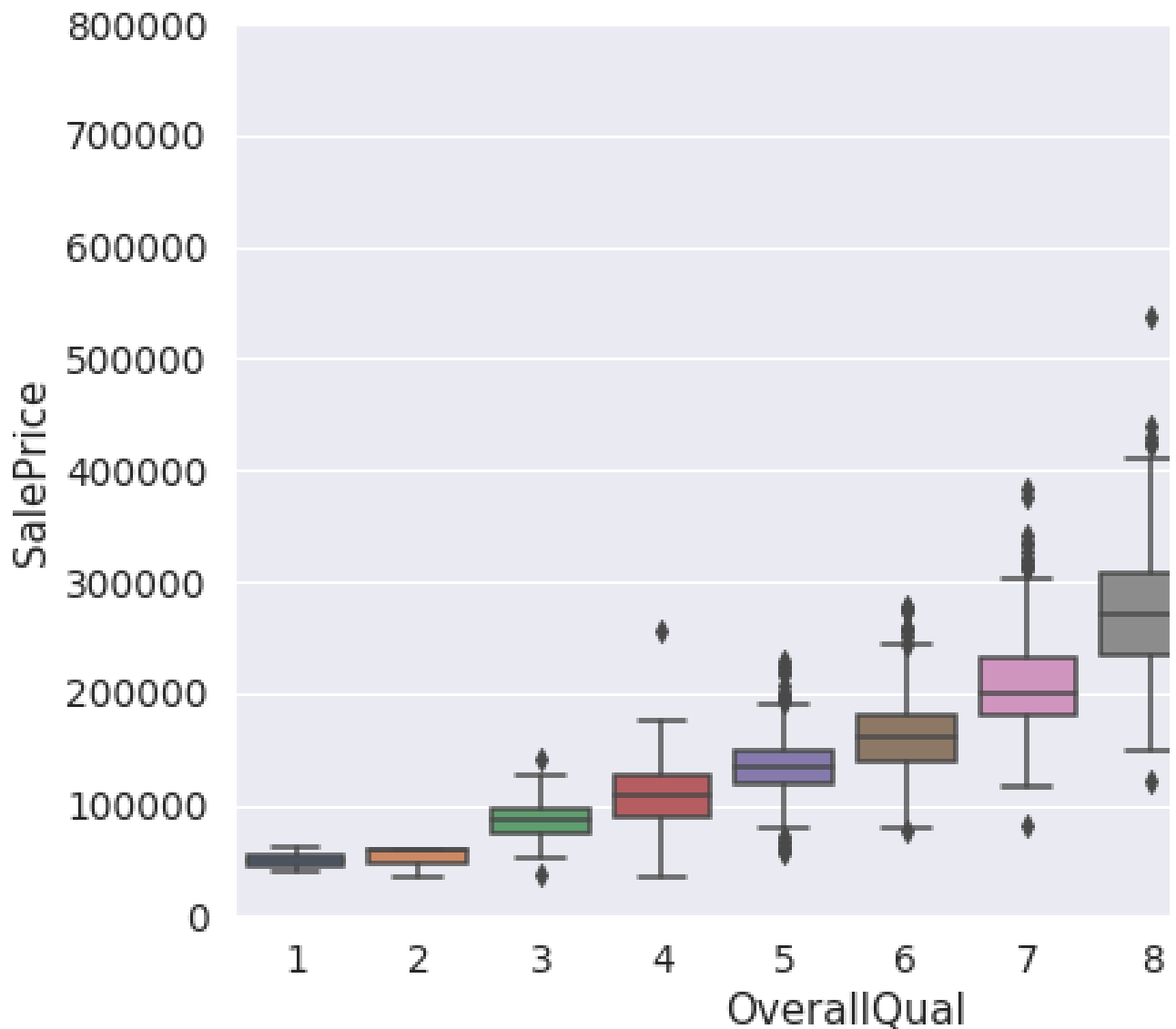




This is also a very intuitive graph and shows that the more the area, more is

```
var = 'OverallQual'
data = pd.concat([df1['SalePrice'], df1[var]], axis=1)
f, ax = plt.subplots(figsize=(8, 6))
fig = sns.boxplot(x=var, y="SalePrice", data=data)
fig.axis(ymin=0, ymax=800000);
```





This again shows a linear relation. This essentially means that as the quality intuitive and is also inherently true.

▼ Data Cleaning

```
test = pd.read_csv('/kaggle/input/house-prices-advanced-regressor
ntrain = df1.shape[0]
ntest = test.shape[0]
y_train = df1.SalePrice.values
all_data = pd.concat((df1, test), sort=False).reset_index(drop=True)
```

```
all_data.drop(['SalePrice'], axis=1, inplace=True)  
print("all_data size is : {}".format(all_data.shape))
```



all_data size is : (2919, 80)

```
total = all_data.isnull().sum().sort_values(ascending=False)  
percent = (all_data.isnull().sum()/all_data.isnull().count()).sort  
missing_data = pd.concat([total, percent], axis=1, keys=['Total',  
missing_data.head(35)
```



	Total	Percent
PoolQC	2909	0.996574
MiscFeature	2814	0.964029
Alley	2721	0.932169
Fence	2348	0.804385
FireplaceQu	1420	0.486468
LotFrontage	486	0.166495
GarageCond	159	0.054471
GarageQual	159	0.054471
GarageYrBlt	159	0.054471
GarageFinish	159	0.054471
GarageType	157	0.053786
BsmtCond	82	0.028092
BsmtExposure	82	0.028092
BsmtQual	81	0.027749
BsmtFinType2	80	0.027407
BsmtFinType1	79	0.027064
MasVnrType	24	0.008222
MasVnrArea	23	0.007879
MSZoning	4	0.001370
BsmtHalfBath	2	0.000685

Utilities	2	0.000685
Functional	2	0.000685
BsmtFullBath	2	0.000685
BsmtFinSF1	1	0.000343
Exterior1st	1	0.000343
Exterior2nd	1	0.000343
BsmtFinSF2	1	0.000343
BsmtUnfSF	1	0.000343
TotalBsmtSF	1	0.000343
SaleType	1	0.000343
Electrical	1	0.000343
KitchenQual	1	0.000343
GarageArea	1	0.000343
GarageCars	1	0.000343
HouseStyle	0	0.000000

Here, we see that PoolQC is the mostly missing data which might be correct house. Similar for Alley, MiscFetaure, FireplaceQu. We can see that these feat Garage values have same missing values and hence that can also be delete have 1 missing values which might mean that there is 1 row that has all the

```
all_data = all_data.drop((missing_data[missing_data['Total'] > 1])
all_data = all_data.drop((missing_data[missing_data['Total'] == 1])
all_data.isnull().sum().max()
```



0

```
train = all_data[:ntrain]
test = all_data[ntrain:]
```

```
train['Neighborhood']
```



```
0      CollgCr
1      Veenker
2      CollgCr
3      Crawfor
4      NoRidge
...
1455    Gilbert
1456      NWAmes
1457    Crawfor
1458      NAmes
1459    Edwards
Name: Neighborhood, Length: 1460, dtype: object
```

▼ Part 3 - Handcrafted Scoring Function

```
# TODO: code for scoring function
df2 = train
f = df2.groupby(['GrLivArea']).mean()
# g = f.rank(1, 'SalePrice')
g1 = df1.groupby('Neighborhood').mean()['SalePrice'].rank(ascending=
g2 = df1.groupby('GrLivArea').mean()['SalePrice'].rank(ascending=
g3 = df1.groupby('OverallQual').mean()['SalePrice'].rank(ascending=
g4 = df1.groupby('GarageCars').mean()['SalePrice'].rank(ascending=
g5 = df1.groupby('GarageArea').mean()['SalePrice'].rank(ascending=
g6 = df1.groupby('TotalBsmtSF').mean()['SalePrice'].rank(ascending=
g7 = df1.groupby('1stFlrSF').mean()['SalePrice'].rank(ascending=Tr
g8 = df1.groupby('FullBath').mean()['SalePrice'].rank(ascending=Tr
```

```

g9 = df1.groupby('TotRmsAbvGrd').mean()['SalePrice'].rank(ascending=
g10 = df1.groupby('YearBuilt').mean()['SalePrice'].rank(ascending=

t1 = g1.to_dict()
t2 = g2.to_dict()
t3 = g3.to_dict()
t4 = g4.to_dict()
t5 = g5.to_dict()
t6 = g6.to_dict()
t7 = g7.to_dict()
t8 = g8.to_dict()
t9 = g9.to_dict()
t10 = g10.to_dict()

train1 = pd.DataFrame()

cols = ['Neighborhood', 'OverallQual', 'GrLivArea', 'GarageCars', 'Gar
        'TotRmsAbvGrd', 'YearBuilt']
for x in cols:
    train1[x] = df1[x]
train1.head()

```



	Neighborhood	OverallQual	GrLivArea	GarageCars	Gar
0	CollgCr	7	1710	2	
1	Veenker	6	1262	2	
2	CollgCr	7	1786	2	
3	Crawfor	7	1717	3	
4	NoRidge	8	2198	3	

```

t90 = pd.DataFrame()
def score_fn():
    c = train1['Neighborhood'].count()
    lst = list(range(c))
    print(train1)
    for i in lst:
        j = int(i)
        t = train1.iloc[j]['Neighborhood']

```



```

    f = t1[t]
    train1.set_value(j, 'Neighborhood', f)
t90 = train1
print(t90)
train1['Sum'] = train1.sum(axis=1)
#train1['Neighborhood']
t3 = train1.nlargest(10, ['Sum'])
t5 = train1.nsmallest(10, ['Sum'])
for index, row in train1.iterrows():
    x1 = row['Neighborhood']
    for item in g1.items():
        if item[1] == x1:
            t3.set_value(index, 'Neighborhood', item[0])

for index, row in train1.iterrows():
    x1 = row['Neighborhood']
    for item in g1.items():
        if item[1] == x1:
            t5.set_value(index, 'Neighborhood', item[0])
t4 = t3.nlargest(10, ['Sum'])
t6 = t5.nlargest(10, ['Sum'])
return t4,t6

```

What is the ten most desirable houses?

```

t4,t6 = score_fn()
t4

```



	Neighborhood	OverallQual	GrLivArea	GarageCars	G
0	CollgCr	7	1710	2	
1	Veenker	6	1262	2	
2	CollgCr	7	1786	2	
3	Crawfor	7	1717	3	
4	NoRidge	8	2198	3	
...	
1455	Gilbert	6	1647	2	
1456	NWAmes	6	2073	2	
1457	Crawfor	7	2340	1	
1458	NAmes	5	1078	1	
1459	Edwards	5	1256	1	

	TotalBsmntSF	1stFlrSF	FullBath	TotRmsAbvGrd	Yea
0	856	856	2	8	
1	1262	1262	2	6	
2	920	920	2	6	
3	756	961	1	7	
4	1145	1145	2	9	
...	
1455	953	953	2	7	
1456	1542	2073	2	7	
1457	1152	1188	2	9	
1458	1078	1078	1	5	
1459	1256	1256	1	6	

[1460 rows x 10 columns]

```
/opt/conda/lib/python3.6/site-packages/ipykernel_launcher
# Remove the CWD from sys.path while we load stuff.
```

	Neighborhood	OverallQual	GrLivArea	GarageCars	G
0	17	7	1710	2	
1	21	6	1262	2	
2	17	7	1786	2	
3	18	7	1717	3	
4	25	8	2198	3	
...	
1455	15	6	1647	2	
1456	14	6	2073	2	

1457	18	7	2340	1
1458	11	5	1078	1
1459	5	5	1256	1

	TotalBsmfSF	1stFlrSF	FullBath	TotRmsAbvGrd	Yea
0	856	856	2	8	
1	1262	1262	2	6	
2	920	920	2	6	
3	756	961	1	7	
4	1145	1145	2	9	
...	
1455	953	953	2	7	
1456	1542	2073	2	7	
1457	1152	1188	2	9	
1458	1078	1078	1	5	
1459	1256	1256	1	6	

[1460 rows x 10 columns]

/opt/conda/lib/python3.6/site-packages/ipykernel_launcher
/opt/conda/lib/python3.6/site-packages/ipykernel_launcher

	Neighborhood	OverallQual	GrLivArea	GarageCars
1298	Edwards	10.0	5642.0	2.0
523	Edwards	10.0	4676.0	3.0
496	NoRidge	8.0	3228.0	2.0
1182	NoRidge	10.0	4476.0	3.0
691	NoRidge	10.0	4316.0	3.0
1373	NoRidge	10.0	2633.0	3.0
440	NridgHt	10.0	2402.0	3.0
1169	NoRidge	10.0	3627.0	3.0
224	NridgHt	10.0	2392.0	3.0

What is the ten least desirable houses?

t90



t6



	Neighborhood	OverallQual	GrLivArea	GarageCars
125	IDOTRR	6.0	754.0	0.0
29	BrkSide	4.0	520.0	1.0
916	IDOTRR	2.0	480.0	1.0
1321	BrkSide	3.0	720.0	1.0
710	BrkSide	3.0	729.0	0.0
528	Edwards	4.0	605.0	0.0
1218	BrkSide	4.0	912.0	0.0
705	IDOTRR	4.0	1092.0	0.0
1100	SWISU	2.0	438.0	1.0
533	BrkSide	1.0	334.0	0.0

I have selected items that have the highest correlation with the SalePrice. To values and made a score for each one of them. The neighborhood has been has been given the highest value. The rest being integer values were added house of area 3000>2000 and hence will give more weight to overall score. :

▼ Part 4 - Pairwise Distance Function

```
# TODO: code for distance function
cols = ['Neighborhood', 'OverallQual', 'GrLivArea', '1stFlrSF', 'FullBath',
        'TotRmsAbvGrd', 'YearBuilt']
for x in cols:
    train1[x] = train[x]
train1.head()
```



	Neighborhood	OverallQual	GrLivArea	GarageCars	Gar
0	CollgCr	7	1710	2	
1	Veenker	6	1262	2	
2	CollgCr	7	1786	2	
3	Crawfor	7	1717	3	
4	NoRidge	8	2198	3	

```
train['Neighborhood'].unique()
```



```
array(['CollgCr', 'Veenker', 'Crawfor', 'NoRidge', 'Mitic',
      'NWAmes', 'OldTown', 'BrkSide', 'Sawyer', 'NridgH',
      'SawyerW', 'IDOTRR', 'MeadowV', 'Edwards', 'Timbe',
      'StoneBr', 'ClearCr', 'NPkVill', 'Blmngtn', 'BrDa
```

```
from sklearn.preprocessing import LabelEncoder
train3 = train
cols = ('Street', 'LotShape', 'LandContour', 'LotConfig', 'LandSlope',
       'HouseStyle', 'RoofStyle', 'RoofMatl', 'ExterQual', 'ExterCond',
       'PavedDrive', 'SaleCondition')
```

```
# process columns, apply LabelEncoder to categorical features
for c in cols:
    lbl = LabelEncoder()
    lbl.fit(list(train3[c].values))
    train3[c] = lbl.transform(list(train3[c].values))
```

```
# shape
print('Shape all_data: {}'.format(train3.shape))
```



Shape all_data: (1460, 46)

/opt/conda/lib/python3.6/site-packages/ipykernel_launcher.py:13:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/10min.html#copy-on-write>
This is added back by InteractiveShellApp.init_path()

```
#train90 = train1
def pairwise(a,b):
    #print(np.linalg.norm(a-b))
    return np.linalg.norm(a-b)
```

```
train90 = train1
```

```
train90.rename(columns={ train90.columns[0]: "Neighborhood" }, inplace=True)
c = train90['Neighborhood'].count()
lst = list(range(c))
print(train90)
for i in lst:
    j = int(i)
    t = train90.iloc[j]['Neighborhood']
    f = t1[t]
    train90.set_value(j, 'Neighborhood', f)
```

train90



```
/opt/conda/lib/python3.6/site-packages/ipykernel_launcher
```

```
# Remove the CWD from sys.path while we load stuff.
```

	Neighborhood	OverallQual	GrLivArea	GarageCars	G
0	CollgCr	7	1710	2	
1	Veenker	6	1262	2	
2	CollgCr	7	1786	2	
3	Crawfor	7	1717	3	
4	NoRidge	8	2198	3	
...	
1455	Gilbert	6	1647	2	
1456	NWAmes	6	2073	2	
1457	Crawfor	7	2340	1	
1458	NAmes	5	1078	1	
1459	Edwards	5	1256	1	

	TotalBsmntSF	1stFlrSF	FullBath	TotRmsAbvGrd	Yea
0	856	856	2	8	
1	1262	1262	2	6	
2	920	920	2	6	
3	756	961	1	7	
4	1145	1145	2	9	
...	
1455	953	953	2	7	
1456	1542	2073	2	7	
1457	1152	1188	2	9	
1458	1078	1078	1	5	
1459	1256	1256	1	6	

```
[1460 rows x 11 columns]
```

	Neighborhood	OverallQual	GrLivArea	GarageCars
0	17	7	1710	2
1	21	6	1262	2
2	17	7	1786	2
3	18	7	1717	3


```

train3 = train
from sklearn.preprocessing import LabelEncoder
cols = ('Street', 'LotShape', 'LandContour', 'LotConfig', 'LandSlope',
        'HouseStyle', 'RoofStyle', 'RoofMatl', 'ExterQual', 'ExterCond',
        'PavedDrive', 'SaleCondition')

# process columns, apply LabelEncoder to categorical features
for c in cols:
    lbl = LabelEncoder()
    lbl.fit(list(train3[c].values))
    train3[c] = lbl.transform(list(train3[c].values))

# shape
print('Shape all_data: {}'.format(train3.shape))

```



Shape all_data: (1460, 46)

/opt/conda/lib/python3.6/site-packages/ipykernel_launcher.py:13:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/10min.html#copy-on-write>
This is added back by InteractiveShellApp.init_path()

How well does the distance function work? When does it do well/badly?

```
train3 = train
```

▼ Part 5 - Clustering

```

from sklearn.cluster import DBSCAN
clustering = DBSCAN(eps=280, min_samples=2, metric = pairwise).fit(train3)
clust = clustering.labels_

```



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

	principal component 1	principal component 2	target
0	-235.736040	355.674021	0
1	44.762207	-340.774393	1
2	53.978297	354.994340	0
3	-210.277845	384.773220	1

```
type(clust)
```

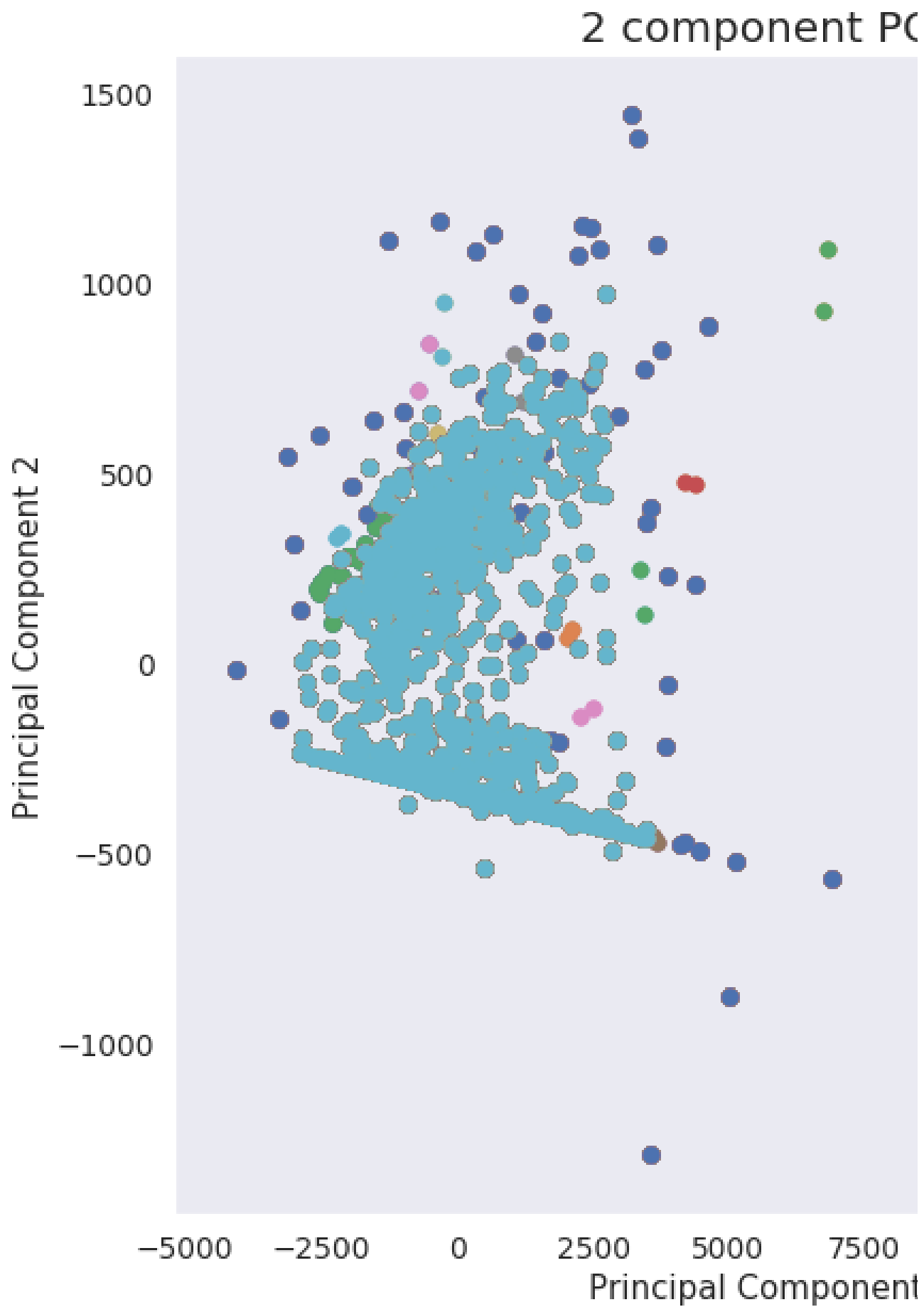


```
numpy.ndarray
```

```
fig = plt.figure(figsize = (10,10))
ax = fig.add_subplot(1,1,1)
ax.set_xlabel('Principal Component 1', fontsize = 15)
ax.set_ylabel('Principal Component 2', fontsize = 15)
ax.set_title('2 component PCA', fontsize = 20)

targets = clust
#colors = ['r', 'g', 'b']
for target in targets:
    indicesToKeep = finalDf['target'] == target
    ax.scatter(finalDf.loc[indicesToKeep, 'principal component 1'],
               finalDf.loc[indicesToKeep, 'principal component 2'],
               s = 50)
ax.legend(clust1[0].unique())
ax.grid()
```





As can be seen from the figure, the clusters are based on the distance and used PCA to decrease the number of components and hence view it in a 2D

▼ Part 6 - Linear Regression

```
from sklearn.preprocessing import LabelEncoder
cols = ('Street', 'LotShape', 'LandContour', 'LotConfig', 'LandSlope',
        'HouseStyle', 'RoofStyle', 'RoofMatl', 'ExterQual', 'ExterCond',
        'PavedDrive', 'SaleCondition')

# process columns, apply LabelEncoder to categorical features
for c in cols:
    lbl = LabelEncoder()
    lbl.fit(list(test[c].values))
    test[c] = lbl.transform(list(test[c].values))
```

```
# shape
print('Shape all_data: {}'.format(test.shape))
```



Shape all_data: (1459, 46)
 /opt/conda/lib/python3.6/site-packages/ipykernel_launcher
 A value is trying to be set on a copy of a slice from a
 Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: <http://pandas.pyda>
 # Remove the CWD from sys.path while we load stuff.

```
# TODO: code for linear regression
train3 = train
from sklearn.linear_model import LinearRegression
model1 = LinearRegression()
model1.fit(train3,y_train)
predict1 = model1.predict(train3)
predict = model1.predict(test)
from sklearn.metrics import mean_squared_error
from math import sqrt
rms = sqrt(mean_squared_error(y_train, predict1))

coeff_df = pd.DataFrame(model1.coef_, train3.columns, columns=['C
```

y_train.shape



(1460,)

rms



34302.96425845283

coeff_df



	Coefficient
Id	-1.566255
MSSubClass	-60.696718
LotArea	0.442976
Street	29152.318379
LotShape	-1126.725301
LandContour	3081.166099
LotConfig	3.040278
LandSlope	10423.086793
Neighborhood	619.297191
Condition1	-1055.019868
Condition2	-8359.690266
BldgType	-3289.645623
HouseStyle	-434.117981
OverallQual	15521.969041
OverallCond	5388.537397
YearBuilt	469.921419
YearRemodAdd	30.138755
RoofStyle	2615.396384
RoofMatl	3398.248784
ExterQual	-17192.187405

ExterCond	986.247012
Foundation	371.686816
Heating	-1179.238785
HeatingQC	-904.423645
CentralAir	-3488.791818
1stFlrSF	37.144754
2ndFlrSF	19.588481
LowQualFinSF	-26.787338
GrLivArea	29.945898
FullBath	-312.409458
HalfBath	-867.346692
BedroomAbvGr	-8654.721697
KitchenAbvGr	-12488.826981
TotRmsAbvGrd	3824.038355
Fireplaces	4639.057311
PavedDrive	1543.345223
WoodDeckSF	30.795121
OpenPorchSF	-10.273531
EnclosedPorch	10.708368
3SsnPorch	19.840339
ScreenPorch	62.689786

PoolArea	-18.343098
MiscVal	-1.594904
MoSold	-220.891501
YrSold	-1154.313793
SaleCondition	2736.791333

How well/badly does it work? Which are the most important variables?

It gives a RMSE of 34302. We can see that it gives a lot of weightage to Street KitchenAbvGrd and also ExterQual. Here, the theta denote the amount of weight

▼ Part 7 - External Dataset

TODO: code to import external dataset and test

```
import random
```

```
lst = []
for i in range(1460):
    lst.append(random.randint(1,5))
i = pd.DataFrame(lst)
i.rename(columns={ i.columns[0]: "Schools" }, inplace = True)
i
# random.randint(3,4)
```



Schools	
0	2
1	4
2	4
3	3
4	4
...	...
1455	3
1456	5
1457	5

```

train4 = train
len(train4.columns)
train4.shape
ext1 = np.ones(1460)
ext1 = ext1*3
ext1 = np.transpose(ext1)
ext1 = np.reshape(ext1, (1460,1))
ext1
ext1 = pd.DataFrame(ext1)
#ext1
#ext1.rename(columns={ ext1.columns[0]: "Schools" }, inplace = True)
# train4 = pd.concat([train4,ext1],axis=1)
train4 = pd.concat([train4,i],axis=1)
y_train = pd.DataFrame(y_train)
y_train.rename(columns={ y_train.columns[0]: "Y" }, inplace = True)
train4 = pd.concat([train4,y_train],axis=1)
from scipy.stats import pearsonr
corr, _ = pearsonr(train4['Schools'], train4['Y'])
corr
#train4['Schools'].isnull().sum()

```



0.030610303498988183

```
train['Neighborhood'].unique()
```



```
array([ 5, 24,  6, 15, 11, 21, 14, 17,  3, 19, 16, 12,  2,
        8, 22,  4, 13,  0,  2, 18,  1])
```

The dataset chosen here in of how many schools are present in the vicinity schools around the area, there are a limited choices they have and hence a still it shows some correlation if not 0. Under exact and correct circumstance might lead to a better prediction

▼ Part 8 - Permutation Test

```
# TODO: code for all permutation tests
train4 = pd.DataFrame()
train4['OverallQual'] = train['OverallQual']
train5 = pd.DataFrame()
train5['GrLivArea'] = train['GrLivArea']
train6 = pd.DataFrame()
train6['Street'] = train['Street']
train7 = pd.DataFrame()
train7['KitchenAbvGr'] = train['KitchenAbvGr']
train8 = pd.DataFrame()
train8['SaleCondition'] = train['SaleCondition']
train9 = pd.DataFrame()
train9['YrSold'] = train['YrSold']
train10 = pd.DataFrame()
train10['LotArea'] = train['LotArea']
train11 = pd.DataFrame()
train11['CentralAir'] = train['CentralAir']
train12 = pd.DataFrame()
train12['FullBath'] = train['FullBath']
train13 = pd.DataFrame()
train13['HalfBath'] = train['HalfBath']
```

```
model2 = LinearRegression()  
model3 = LinearRegression()  
model4 = LinearRegression()  
model5 = LinearRegression()  
model6 = LinearRegression()  
model7 = LinearRegression()  
model8 = LinearRegression()  
model9 = LinearRegression()  
model10 = LinearRegression()  
model11 = LinearRegression()
```

```
model2.fit(train4,y_train)  
model3.fit(train5,y_train)  
model4.fit(train6,y_train)  
model5.fit(train7,y_train)  
model6.fit(train8,y_train)  
model7.fit(train9,y_train)  
model8.fit(train10,y_train)  
model9.fit(train11,y_train)  
model10.fit(train12,y_train)  
model11.fit(train13,y_train)
```

```
test1 = pd.DataFrame()  
test1['OverallQual'] = test['OverallQual']  
test2 = pd.DataFrame()  
test2['GrLivArea'] = test['GrLivArea']  
test3 = pd.DataFrame()  
test3['Street'] = test['Street']  
test3 = pd.get_dummies(test3)  
test4 = pd.DataFrame()  
test4['KitchenAbvGr'] = test['KitchenAbvGr']  
test5 = pd.DataFrame()  
test5['SaleCondition'] = test['SaleCondition']  
test5 = pd.get_dummies(test5)  
test6 = pd.DataFrame()  
test6['YrSold'] = test['YrSold']  
test7 = pd.DataFrame()  
test7['LotArea'] = test['LotArea']  
test8 = pd.DataFrame()  
test8['CentralAir'] = test['CentralAir']  
test9 = pd.DataFrame()  
test9['FullBath'] = test['FullBath']  
test10 = pd.DataFrame()  
test10['HalfBath'] = test['HalfBath']
```

```
#test5 = pd.get_dummies(test5)
```

```
predicted_prices2 = model2.predict(test1)
predicted_prices3 = model3.predict(test2)
predicted_prices4 = model4.predict(test3)
predicted_prices5 = model5.predict(test4)
predicted_prices6 = model6.predict(test5)
predicted_prices7 = model7.predict(test6)
predicted_prices8 = model8.predict(test7)
predicted_prices9 = model9.predict(test8)
predicted_prices10 = model10.predict(test9)
predicted_prices11 = model11.predict(test10)
```

```
y_train_com = y_train.sample(100,random_state=1)
```

```
predicted_prices2_comp = pd.DataFrame(predicted_prices2).sample(100)
predicted_prices3_comp = pd.DataFrame(predicted_prices3).sample(100)
predicted_prices4_comp = pd.DataFrame(predicted_prices4).sample(100)
predicted_prices5_comp = pd.DataFrame(predicted_prices5).sample(100)
predicted_prices6_comp = pd.DataFrame(predicted_prices6).sample(100)
predicted_prices7_comp = pd.DataFrame(predicted_prices7).sample(100)
predicted_prices8_comp = pd.DataFrame(predicted_prices8).sample(100)
predicted_prices9_comp = pd.DataFrame(predicted_prices9).sample(100)
predicted_prices10_comp = pd.DataFrame(predicted_prices10).sample(100)
predicted_prices11_comp = pd.DataFrame(predicted_prices11).sample(100)
```

```
from sklearn.metrics import mean_squared_error
#y_train_com = pd.DataFrame(y_train_com).fillna('0')
predicted_prices2_comp = predicted_prices2_comp.reset_index()
#y_train_com = y_train_com.reset_index()
np.reshape(y_train_com, (100,1))
predicted_prices2_comp.rename(columns={ predicted_prices2_comp.columns[0]: 'predicted_prices21_comp'})
predicted_prices21_comp = pd.DataFrame()
predicted_prices21_comp['SalesPrice'] = predicted_prices2_comp['SalesPrice']
print("MSE1",mean_squared_error(np.log(y_train_com), np.log(predicted_prices21_comp)))
print("MSE2",mean_squared_error(np.log(y_train_com), np.log(predicted_prices21_comp)))
print("MSE3",mean_squared_error(np.log(y_train_com), np.log(predicted_prices21_comp)))
print("MSE4",mean_squared_error(np.log(y_train_com), np.log(predicted_prices21_comp)))
print("MSE5",mean_squared_error(np.log(y_train_com), np.log(predicted_prices21_comp)))
print("MSE6",mean_squared_error(np.log(y_train_com), np.log(predicted_prices21_comp)))
print("MSE7",mean_squared_error(np.log(y_train_com), np.log(predicted_prices21_comp)))
print("MSE8",mean_squared_error(np.log(y_train_com), np.log(predicted_prices21_comp)))
print("MSE9",mean_squared_error(np.log(y_train_com), np.log(predicted_prices21_comp)))
print("MSE10",mean_squared_error(np.log(y_train_com), np.log(predicted_prices21_comp)))
```

```

# We will look at the predicted prices to ensure we have something
#print(predicted_prices2)
score_dict={}
pvalue_dict={}
pvalue_lst = []
from sklearn.model_selection import permutation_test_score
score1,perm_score, pvalue= permutation_test_score(model2, train4,
print("OverallQual", score1)
t=train4.columns.values
pvalue_dict[str(t)]=pvalue
score2,perm_score, pvalue= permutation_test_score(model3, train5,
print("GrLiveArea", score2)
t=train5.columns.values
pvalue_dict[str(t)]=pvalue
score3,perm_score, pvalue= permutation_test_score(model4, train6,
print("Steet", score3)
t=train6.columns.values
pvalue_dict[str(t)]=pvalue
score4,perm_score, pvalue= permutation_test_score(model5, train7,
print("KitchenAbvGrd", score4)
t=train7.columns.values
pvalue_dict[str(t)]=pvalue
score5,perm_score, pvalue= permutation_test_score(model6, train8,
t=train8.columns.values
print("SaleCondition",score5)
pvalue_dict[str(t)]=pvalue
#pvalue_dict

```



```
MSE1 0.29654644179500783
MSE2 0.2536328196775782
MSE3 0.17759915876877005
MSE4 0.18886329487491202
MSE5 0.19390747652286883
MSE6 0.1752988431039135
MSE7 0.17458088838859515
MSE8 0.1959718844328004
MSE9 0.2147375052725084
MSE10 0.18619685832443314
/opt/conda/lib/python3.6/site-packages/sklearn/model_sel
warnings.warn(CV_WARNING, FutureWarning)
OverallQual 0.6247257478644835
/opt/conda/lib/python3.6/site-packages/sklearn/model_sel
warnings.warn(CV_WARNING, FutureWarning)
GrLiveArea 0.4924726967831159
/opt/conda/lib/python3.6/site-packages/sklearn/model sel
```

predicted_prices21_comp



	SalesPrice
0	176408.736044
1	221844.538637
2	221844.538637
3	221844.538637
4	312716.143823
...	...
95	130972.933451

```
y_train2 = pd.DataFrame(y_train)
y_train2.isnull().sum()
```



```
Y      0
dtype: int64
```

```
95    130972.933451
```

```
train4.isnull().sum()
```



```
OverallQual      0
dtype: int64
```

```
import random
c1=0
for i in range(100):
#     print(train4.isnull().sum())
#     print(y_train1.isnull().sum())
    train40 = pd.concat([train4, y_train2], axis=1)
#     print(train40.isnull().sum())
    train41 = train40.sample(1000)
    train41.rename(columns={ train41.columns[1]: "SalePrice" }, inplace=True)
    y_train1 = train41['SalePrice']
    train41 = train41.drop('OverallQual',axis=1)
    score,perm_score, pvalue= permutation_test_score(model2, train41,
```



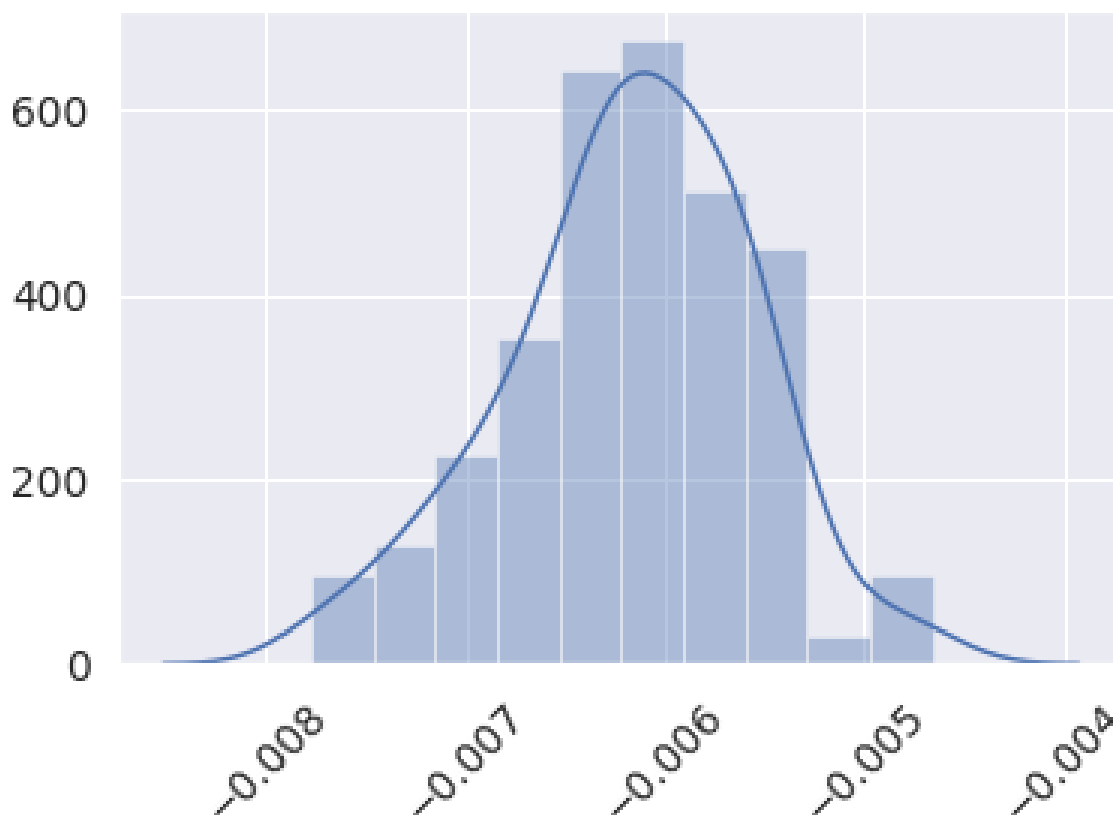
```
    c1=c1+perm_score  
c1=c1/100  
#c1  
s = sns.distplot(c1)  
plt.xticks(rotation=45)
```



1 . / 1 7 0 1 / . 1 2 1 0 . 1 / 1 7 / 1 7

<https://colab.research.google.com/drive/1LLI98XWtN1zqx8uZusjASoqDjxdCHhVg?authuser=1#scrollTo=8xEjXVNqily6&printMode=true>


```
warnings.warn(CV_WARNING, FutureWarning)
/opt/conda/lib/python3.6/site-packages/sklearn/model_sel
warnings.warn(CV_WARNING, FutureWarning)
/opt/conda/lib/python3.6/site-packages/sklearn/model_sel
warnings.warn(CV_WARNING, FutureWarning)
/opt/conda/lib/python3.6/site-packages/sklearn/model_sel
warnings.warn(CV_WARNING, FutureWarning)
/opt/conda/lib/python3.6/site-packages/sklearn/model_sel
warnings.warn(CV_WARNING, FutureWarning)
(array([-0.009, -0.008, -0.007, -0.006, -0.005, -0.004,
<a list of 7 Text xticklabel objects>)
```



```
train40.isnull().sum()
```



```
OverallQual    0
Y              0
dtype: int64
```

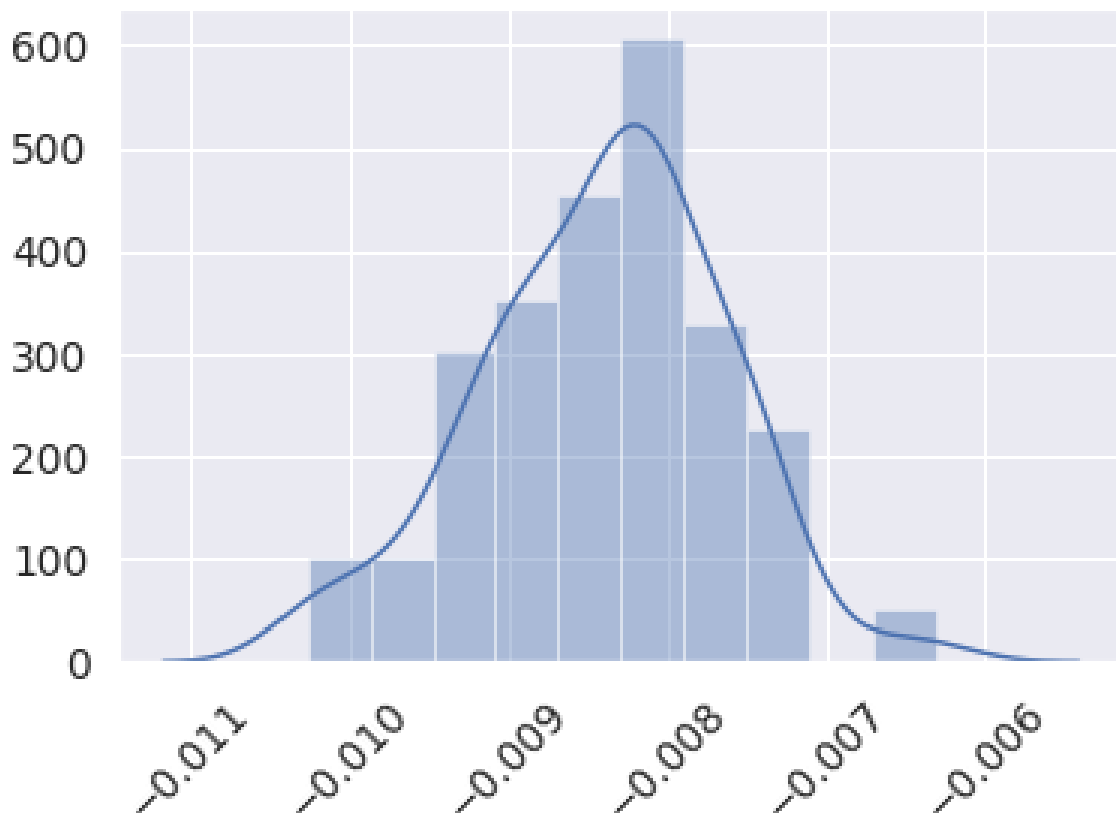
```
c2=0
for i in range(100):
    train50 = pd.concat([train5, y_train2], axis=1)
    train51 = train50.sample(750)
    train51.rename(columns={ train51.columns[1]: "SalePrice" }, inplace=True)
    y_train1 = train51['SalePrice']
    train51 = train51.drop('GrLivArea',axis=1)
    score,perm_score, pvalue= permutation_test_score(model2, train51, y_train1)
    c2=c2+perm_score
c2=c2/100
c2
s = sns.distplot(c2)
plt.xticks(rotation=45)
```



1 . / 1 7 0 1 / . 1 2 1 0 . 1 / 1 7 / 1 7

<https://colab.research.google.com/drive/1LLI98XWtN1zqx8uZusjASoqDjxdCHhVg?authuser=1#scrollTo=8xEjXVNqily6&printMode=true>


```
warnings.warn(CV_WARNING, FutureWarning)
/opt/conda/lib/python3.6/site-packages/sklearn/model_sel
warnings.warn(CV_WARNING, FutureWarning)
/opt/conda/lib/python3.6/site-packages/sklearn/model_sel
warnings.warn(CV_WARNING, FutureWarning)
/opt/conda/lib/python3.6/site-packages/sklearn/model_sel
warnings.warn(CV_WARNING, FutureWarning)
/opt/conda/lib/python3.6/site-packages/sklearn/model_sel
warnings.warn(CV_WARNING, FutureWarning)
(array([-0.012, -0.011, -0.01 , -0.009, -0.008, -0.007,
      <a list of 8 Text xticklabel objects>)
```



```
c3=0
for i in range(100):
    train60 = pd.concat([train6, y_train2], axis=1)
    train61 = train60.sample(1000)
    train61.rename(columns={ train61.columns[1]: "SalePrice" }, inplace=True)
    y_train1 = train61['SalePrice']
    train61 = train61.drop('Street',axis=1)
    score,perm_score, pvalue= permutation_test_score(model2, train61, y_train1)
    c3=c3+perm_score
```

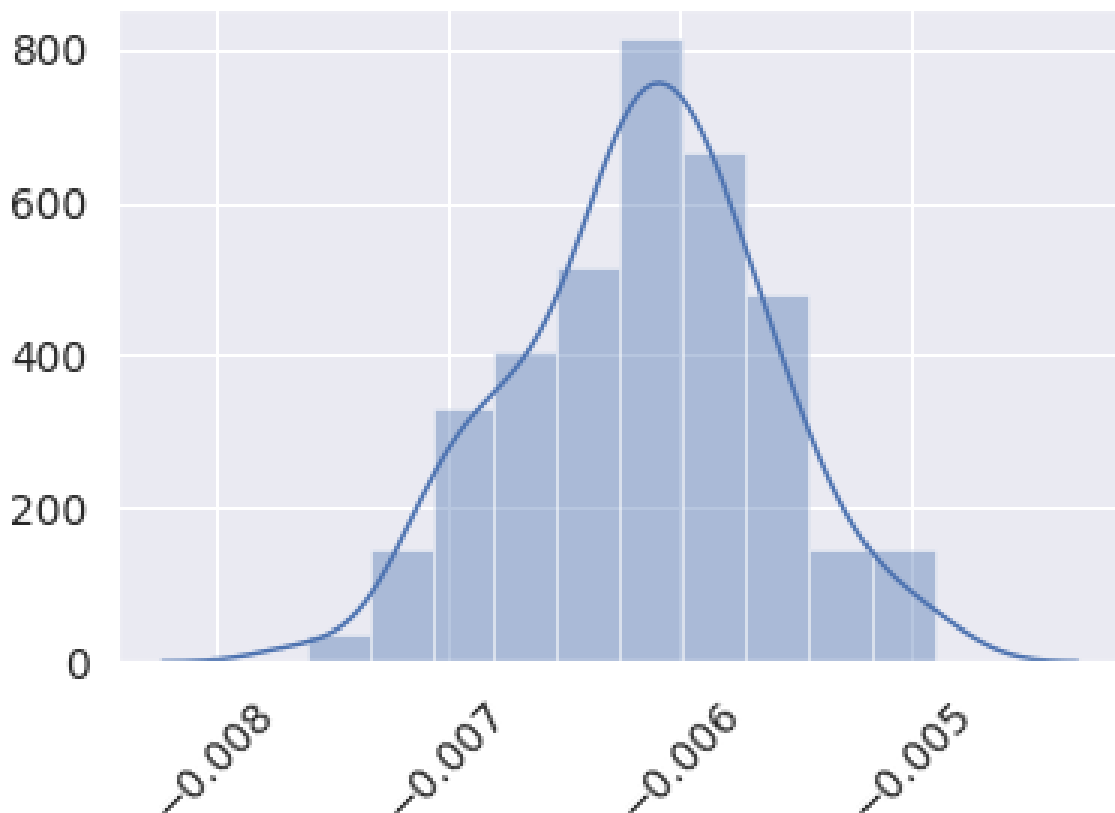
```
c3=c3/100  
c3  
s = sns.distplot(c3)  
plt.xticks(rotation=45)
```



1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100 101 102 103 104 105 106 107 108 109 110 111 112 113 114 115 116 117 118 119 120 121 122 123 124 125 126 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143 144 145 146 147 148 149 150 151 152 153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168 169 170 171 172 173 174 175 176 177 178 179 180 181 182 183 184 185 186 187 188 189 190 191 192 193 194 195 196 197 198 199 200 201 202 203 204 205 206 207 208 209 210 211 212 213 214 215 216 217 218 219 220 221 222 223 224 225 226 227 228 229 230 231 232 233 234 235 236 237 238 239 240 241 242 243 244 245 246 247 248 249 250 251 252 253 254 255 256 257 258 259 260 261 262 263 264 265 266 267 268 269 270 271 272 273 274 275 276 277 278 279 280 281 282 283 284 285 286 287 288 289 290 291 292 293 294 295 296 297 298 299 300 301 302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 378 379 380 381 382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404 405 406 407 408 409 410 411 412 413 414 415 416 417 418 419 420 421 422 423 424 425 426 427 428 429 430 431 432 433 434 435 436 437 438 439 440 441 442 443 444 445 446 447 448 449 450 451 452 453 454 455 456 457 458 459 460 461 462 463 464 465 466 467 468 469 470 471 472 473 474 475 476 477 478 479 480 481 482 483 484 485 486 487 488 489 490 491 492 493 494 495 496 497 498 499 500 501 502 503 504 505 506 507 508 509 510 511 512 513 514 515 516 517 518 519 520 521 522 523 524 525 526 527 528 529 530 531 532 533 534 535 536 537 538 539 540 541 542 543 544 545 546 547 548 549 550 551 552 553 554 555 556 557 558 559 560 561 562 563 564 565 566 567 568 569 570 571 572 573 574 575 576 577 578 579 580 581 582 583 584 585 586 587 588 589 590 591 592 593 594 595 596 597 598 599 600 601 602 603 604 605 606 607 608 609 610 611 612 613 614 615 616 617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 639 640 641 642 643 644 645 646 647 648 649 650 651 652 653 654 655 656 657 658 659 660 661 662 663 664 665 666 667 668 669 670 671 672 673 674 675 676 677 678 679 680 681 682 683 684 685 686 687 688 689 690 691 692 693 694 695 696 697 698 699 700 701 702 703 704 705 706 707 708 709 710 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771 772 773 774 775 776 777 778 779 780 781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808 809 810 811 812 813 814 815 816 817 818 819 820 821 822 823 824 825 826 827 828 829 830 831 832 833 834 835 836 837 838 839 840 841 842 843 844 845 846 847 848 849 850 851 852 853 854 855 856 857 858 859 860 861 862 863 864 865 866 867 868 869 870 871 872 873 874 875 876 877 878 879 880 881 882 883 884 885 886 887 888 889 890 891 892 893 894 895 896 897 898 899 900 901 902 903 904 905 906 907 908 909 910 911 912 913 914 915 916 917 918 919 920 921 922 923 924 925 926 927 928 929 930 931 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949 950 951 952 953 954 955 956 957 958 959 960 961 962 963 964 965 966 967 968 969 970 971 972 973 974 975 976 977 978 979 980 981 982 983 984 985 986 987 988 989 990 991 992 993 994 995 996 997 998 999 1000 1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 1020 1021 1022 1023 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 1040 1

<https://colab.research.google.com/drive/1LLI98XWtN1zqx8uZusjASoqDjxdCHhVg?authuser=1#scrollTo=8xEjXVNqily6&printMode=true>


```
warnings.warn(CV_WARNING, FutureWarning)
/opt/conda/lib/python3.6/site-packages/sklearn/model_sel
warnings.warn(CV_WARNING, FutureWarning)
/opt/conda/lib/python3.6/site-packages/sklearn/model_sel
warnings.warn(CV_WARNING, FutureWarning)
/opt/conda/lib/python3.6/site-packages/sklearn/model_sel
warnings.warn(CV_WARNING, FutureWarning)
/opt/conda/lib/python3.6/site-packages/sklearn/model_sel
warnings.warn(CV_WARNING, FutureWarning)
(array([-0.009, -0.008, -0.007, -0.006, -0.005, -0.004])
<a list of 6 Text xticklabel objects>)
```



```
c4=0
for i in range(100):
    train70 = pd.concat([train7, y_train2], axis=1)
    train71 = train70.sample(1000)
    train71.rename(columns={ train71.columns[1]: "SalePrice" }, inplace=True)
    y_train1 = train71['SalePrice']
    train71 = train71.drop('KitchenAbvGr',axis=1)
    score,perm_score, pvalue= permutation_test_score(model2, train71, y_train1)
    c4=c4+perm_score
```

```
c4=c4/100  
c4  
s = sns.distplot(c4)  
plt.xticks(rotation=45)
```



