```python
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
import os
from scipy.stats import norm
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import MinMaxScaler
from sklearn.preprocessing import LabelEncoder
from sklearn.impute import SimpleImputer
from scipy import stats
from scipy.stats import norm

import warnings
warnings.filterwarnings('ignore')

%matplotlib inline
sns.set()

df_train = pd.read_csv("train.csv")
df_test = pd.read_csv("test.csv")

print(df_train.shape)
print("*"*50)
print(df_test.shape)
```

```
(1460, 81)
**************************************************
(1459, 80)
```

```python
df_train.head()
```

```
   Id  MSSubClass MSZoning  LotFrontage  LotArea Street Alley LotShape  \
0   1          60       RL         65.0     8450   Pave   NaN      Reg

1   2          20       RL         80.0     9600   Pave   NaN      Reg

2   3          60       RL         68.0    11250   Pave   NaN      IR1

3   4          70       RL         60.0     9550   Pave   NaN      IR1

4   5          60       RL         84.0    14260   Pave   NaN      IR1


   LandContour Utilities  ... PoolArea PoolQC Fence MiscFeature MiscVal
MoSold  \
0          Lvl    AllPub  ...        0    NaN   NaN         NaN       0
2
1          Lvl    AllPub  ...        0    NaN   NaN         NaN       0
```

```
5
2          Lvl    AllPub  ...      0    NaN    NaN         NaN       0
9
3          Lvl    AllPub  ...      0    NaN    NaN         NaN       0
2
4          Lvl    AllPub  ...      0    NaN    NaN         NaN       0
12

   YrSold   SaleType   SaleCondition   SalePrice
0   2008         WD          Normal      208500
1   2007         WD          Normal      181500
2   2008         WD          Normal      223500
3   2006         WD         Abnorml      140000
4   2008         WD          Normal      250000

[5 rows x 81 columns]

df_test.head()
      Id  MSSubClass  MSZoning   LotFrontage   LotArea  Street  Alley
LotShape  \
0  1461          20        RH          80.0     11622    Pave    NaN
Reg
1  1462          20        RL          81.0     14267    Pave    NaN
IR1
2  1463          60        RL          74.0     13830    Pave    NaN
IR1
3  1464          60        RL          78.0      9978    Pave    NaN
IR1
4  1465         120        RL          43.0      5005    Pave    NaN
IR1

   LandContour  Utilities   ...  ScreenPorch  PoolArea  PoolQC   Fence
MiscFeature  \
0          Lvl    AllPub   ...          120         0     NaN   MnPrv
NaN
1          Lvl    AllPub   ...            0         0     NaN     NaN
Gar2
2          Lvl    AllPub   ...            0         0     NaN   MnPrv
NaN
3          Lvl    AllPub   ...            0         0     NaN     NaN
NaN
4          HLS    AllPub   ...          144         0     NaN     NaN
NaN

   MiscVal  MoSold   YrSold   SaleType   SaleCondition
0        0       6     2010         WD          Normal
1    12500       6     2010         WD          Normal
2        0       3     2010         WD          Normal
3        0       6     2010         WD          Normal
```

```
4       0       1    2010           WD           Normal
```

[5 rows x 80 columns]

## EDA
```
df_train.describe()
                 Id    MSSubClass   LotFrontage         LotArea
OverallQual  \
count  1460.000000  1460.000000  1201.000000     1460.000000
1460.000000
mean    730.500000    56.897260    70.049958    10516.828082
6.099315
std     421.610009    42.300571    24.284752     9981.264932
1.382997
min       1.000000    20.000000    21.000000     1300.000000
1.000000
25%     365.750000    20.000000    59.000000     7553.500000
5.000000
50%     730.500000    50.000000    69.000000     9478.500000
6.000000
75%    1095.250000    70.000000    80.000000    11601.500000
7.000000
max    1460.000000   190.000000   313.000000   215245.000000
10.000000

        OverallCond     YearBuilt  YearRemodAdd    MasVnrArea
BsmtFinSF1   ...  \
count  1460.000000  1460.000000   1460.000000  1452.000000
1460.000000   ...
mean      5.575342  1971.267808   1984.865753   103.685262
443.639726   ...
std       1.112799    30.202904     20.645407   181.066207
456.098091   ...
min       1.000000  1872.000000   1950.000000     0.000000
0.000000   ...
25%       5.000000  1954.000000   1967.000000     0.000000
0.000000   ...
50%       5.000000  1973.000000   1994.000000     0.000000
383.500000   ...
75%       6.000000  2000.000000   2004.000000   166.000000
712.250000   ...
max       9.000000  2010.000000   2010.000000  1600.000000
5644.000000   ...

         WoodDeckSF   OpenPorchSF  EnclosedPorch     3SsnPorch
ScreenPorch  \
count  1460.000000  1460.000000    1460.000000  1460.000000
1460.000000
```

```
mean       94.244521      46.660274      21.954110       3.409589
15.060959
std       125.338794      66.256028      61.119149      29.317331
55.757415
min         0.000000       0.000000       0.000000       0.000000
0.000000
25%         0.000000       0.000000       0.000000       0.000000
0.000000
50%         0.000000      25.000000       0.000000       0.000000
0.000000
75%       168.000000      68.000000       0.000000       0.000000
0.000000
max       857.000000     547.000000     552.000000     508.000000
480.000000

             PoolArea        MiscVal         MoSold         YrSold
SalePrice
count   1460.000000    1460.000000    1460.000000    1460.000000
1460.000000
mean       2.758904      43.489041       6.321918    2007.815753
180921.195890
std       40.177307     496.123024       2.703626       1.328095
79442.502883
min         0.000000       0.000000       1.000000    2006.000000
34900.000000
25%         0.000000       0.000000       5.000000    2007.000000
129975.000000
50%         0.000000       0.000000       6.000000    2008.000000
163000.000000
75%         0.000000       0.000000       8.000000    2009.000000
214000.000000
max       738.000000   15500.000000      12.000000    2010.000000
755000.000000

[8 rows x 38 columns]

df_test.describe()

                Id    MSSubClass    LotFrontage         LotArea
OverallQual  \
count   1459.000000    1459.000000    1232.000000    1459.000000
1459.000000
mean    2190.000000      57.378341      68.580357    9819.161069
6.078821
std      421.321334      42.746880      22.376841    4955.517327
1.436812
min     1461.000000      20.000000      21.000000    1470.000000
1.000000
25%     1825.500000      20.000000      58.000000    7391.000000
5.000000
```

```
50%     2190.000000     50.000000     67.000000    9399.000000
6.000000
75%     2554.500000     70.000000     80.000000   11517.500000
7.000000
max     2919.000000    190.000000    200.000000   56600.000000
10.000000

        OverallCond    YearBuilt  YearRemodAdd    MasVnrArea
BsmtFinSF1  ...  \
count  1459.000000  1459.000000   1459.000000   1444.000000
1458.000000  ...
mean      5.553804  1971.357779   1983.662783    100.709141
439.203704  ...
std       1.113740    30.390071     21.130467    177.625900
455.268042  ...
min       1.000000  1879.000000   1950.000000      0.000000
0.000000  ...
25%       5.000000  1953.000000   1963.000000      0.000000
0.000000  ...
50%       5.000000  1973.000000   1992.000000      0.000000
350.500000  ...
75%       6.000000  2001.000000   2004.000000    164.000000
753.500000  ...
max       9.000000  2010.000000   2010.000000   1290.000000
4010.000000  ...

         GarageArea   WoodDeckSF   OpenPorchSF  EnclosedPorch
3SsnPorch  \
count  1458.000000  1459.000000   1459.000000    1459.000000
1459.000000
mean    472.768861    93.174777     48.313914      24.243317
1.794380
std     217.048611   127.744882     68.883364      67.227765
20.207842
min       0.000000     0.000000      0.000000       0.000000
0.000000
25%     318.000000     0.000000      0.000000       0.000000
0.000000
50%     480.000000     0.000000     28.000000       0.000000
0.000000
75%     576.000000   168.000000     72.000000       0.000000
0.000000
max    1488.000000  1424.000000    742.000000    1012.000000
360.000000

        ScreenPorch     PoolArea       MiscVal        MoSold
YrSold
count  1459.000000  1459.000000   1459.000000   1459.000000
1459.000000
mean     17.064428     1.744345     58.167923      6.104181
```

```
        2007.769705
std       56.609763     30.491646     630.806978      2.722432
1.301740
min        0.000000      0.000000       0.000000      1.000000
2006.000000
25%        0.000000      0.000000       0.000000      4.000000
2007.000000
50%        0.000000      0.000000       0.000000      6.000000
2008.000000
75%        0.000000      0.000000       0.000000      8.000000
2009.000000
max      576.000000    800.000000   17000.000000     12.000000
2010.000000

[8 rows x 37 columns]

df_train.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1460 entries, 0 to 1459
Data columns (total 81 columns):
 #   Column         Non-Null Count   Dtype
---  ------         --------------   -----
 0   Id             1460 non-null    int64
 1   MSSubClass     1460 non-null    int64
 2   MSZoning       1460 non-null    object
 3   LotFrontage    1201 non-null    float64
 4   LotArea        1460 non-null    int64
 5   Street         1460 non-null    object
 6   Alley          91 non-null      object
 7   LotShape       1460 non-null    object
 8   LandContour    1460 non-null    object
 9   Utilities      1460 non-null    object
 10  LotConfig      1460 non-null    object
 11  LandSlope      1460 non-null    object
 12  Neighborhood   1460 non-null    object
 13  Condition1     1460 non-null    object
 14  Condition2     1460 non-null    object
 15  BldgType       1460 non-null    object
 16  HouseStyle     1460 non-null    object
 17  OverallQual    1460 non-null    int64
 18  OverallCond    1460 non-null    int64
 19  YearBuilt      1460 non-null    int64
 20  YearRemodAdd   1460 non-null    int64
 21  RoofStyle      1460 non-null    object
 22  RoofMatl       1460 non-null    object
 23  Exterior1st    1460 non-null    object
 24  Exterior2nd    1460 non-null    object
 25  MasVnrType     1452 non-null    object
 26  MasVnrArea     1452 non-null    float64
```

```
27  ExterQual       1460 non-null   object
28  ExterCond       1460 non-null   object
29  Foundation      1460 non-null   object
30  BsmtQual        1423 non-null   object
31  BsmtCond        1423 non-null   object
32  BsmtExposure    1422 non-null   object
33  BsmtFinType1    1423 non-null   object
34  BsmtFinSF1      1460 non-null   int64
35  BsmtFinType2    1422 non-null   object
36  BsmtFinSF2      1460 non-null   int64
37  BsmtUnfSF       1460 non-null   int64
38  TotalBsmtSF     1460 non-null   int64
39  Heating         1460 non-null   object
40  HeatingQC       1460 non-null   object
41  CentralAir      1460 non-null   object
42  Electrical      1459 non-null   object
43  1stFlrSF        1460 non-null   int64
44  2ndFlrSF        1460 non-null   int64
45  LowQualFinSF    1460 non-null   int64
46  GrLivArea       1460 non-null   int64
47  BsmtFullBath    1460 non-null   int64
48  BsmtHalfBath    1460 non-null   int64
49  FullBath        1460 non-null   int64
50  HalfBath        1460 non-null   int64
51  BedroomAbvGr    1460 non-null   int64
52  KitchenAbvGr    1460 non-null   int64
53  KitchenQual     1460 non-null   object
54  TotRmsAbvGrd    1460 non-null   int64
55  Functional      1460 non-null   object
56  Fireplaces      1460 non-null   int64
57  FireplaceQu     770 non-null    object
58  GarageType      1379 non-null   object
59  GarageYrBlt     1379 non-null   float64
60  GarageFinish    1379 non-null   object
61  GarageCars      1460 non-null   int64
62  GarageArea      1460 non-null   int64
63  GarageQual      1379 non-null   object
64  GarageCond      1379 non-null   object
65  PavedDrive      1460 non-null   object
66  WoodDeckSF      1460 non-null   int64
67  OpenPorchSF     1460 non-null   int64
68  EnclosedPorch   1460 non-null   int64
69  3SsnPorch       1460 non-null   int64
70  ScreenPorch     1460 non-null   int64
71  PoolArea        1460 non-null   int64
72  PoolQC          7 non-null      object
73  Fence           281 non-null    object
74  MiscFeature     54 non-null     object
75  MiscVal         1460 non-null   int64
76  MoSold          1460 non-null   int64
```

```
 77  YrSold          1460 non-null   int64
 78  SaleType        1460 non-null   object
 79  SaleCondition   1460 non-null   object
 80  SalePrice       1460 non-null   int64
dtypes: float64(3), int64(35), object(43)
memory usage: 924.0+ KB

df_test.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1459 entries, 0 to 1458
Data columns (total 80 columns):
 #   Column        Non-Null Count  Dtype
---  ------        --------------  -----
 0   Id            1459 non-null   int64
 1   MSSubClass    1459 non-null   int64
 2   MSZoning      1455 non-null   object
 3   LotFrontage   1232 non-null   float64
 4   LotArea       1459 non-null   int64
 5   Street        1459 non-null   object
 6   Alley         107 non-null    object
 7   LotShape      1459 non-null   object
 8   LandContour   1459 non-null   object
 9   Utilities     1457 non-null   object
 10  LotConfig     1459 non-null   object
 11  LandSlope     1459 non-null   object
 12  Neighborhood  1459 non-null   object
 13  Condition1    1459 non-null   object
 14  Condition2    1459 non-null   object
 15  BldgType      1459 non-null   object
 16  HouseStyle    1459 non-null   object
 17  OverallQual   1459 non-null   int64
 18  OverallCond   1459 non-null   int64
 19  YearBuilt     1459 non-null   int64
 20  YearRemodAdd  1459 non-null   int64
 21  RoofStyle     1459 non-null   object
 22  RoofMatl      1459 non-null   object
 23  Exterior1st   1458 non-null   object
 24  Exterior2nd   1458 non-null   object
 25  MasVnrType    1443 non-null   object
 26  MasVnrArea    1444 non-null   float64
 27  ExterQual     1459 non-null   object
 28  ExterCond     1459 non-null   object
 29  Foundation    1459 non-null   object
 30  BsmtQual      1415 non-null   object
 31  BsmtCond      1414 non-null   object
 32  BsmtExposure  1415 non-null   object
 33  BsmtFinType1  1417 non-null   object
 34  BsmtFinSF1    1458 non-null   float64
 35  BsmtFinType2  1417 non-null   object
 36  BsmtFinSF2    1458 non-null   float64
```

```
37   BsmtUnfSF      1458 non-null    float64
38   TotalBsmtSF    1458 non-null    float64
39   Heating        1459 non-null    object
40   HeatingQC      1459 non-null    object
41   CentralAir     1459 non-null    object
42   Electrical     1459 non-null    object
43   1stFlrSF       1459 non-null    int64
44   2ndFlrSF       1459 non-null    int64
45   LowQualFinSF   1459 non-null    int64
46   GrLivArea      1459 non-null    int64
47   BsmtFullBath   1457 non-null    float64
48   BsmtHalfBath   1457 non-null    float64
49   FullBath       1459 non-null    int64
50   HalfBath       1459 non-null    int64
51   BedroomAbvGr   1459 non-null    int64
52   KitchenAbvGr   1459 non-null    int64
53   KitchenQual    1458 non-null    object
54   TotRmsAbvGrd   1459 non-null    int64
55   Functional     1457 non-null    object
56   Fireplaces     1459 non-null    int64
57   FireplaceQu    729 non-null     object
58   GarageType     1383 non-null    object
59   GarageYrBlt    1381 non-null    float64
60   GarageFinish   1381 non-null    object
61   GarageCars     1458 non-null    float64
62   GarageArea     1458 non-null    float64
63   GarageQual     1381 non-null    object
64   GarageCond     1381 non-null    object
65   PavedDrive     1459 non-null    object
66   WoodDeckSF     1459 non-null    int64
67   OpenPorchSF    1459 non-null    int64
68   EnclosedPorch  1459 non-null    int64
69   3SsnPorch      1459 non-null    int64
70   ScreenPorch    1459 non-null    int64
71   PoolArea       1459 non-null    int64
72   PoolQC         3 non-null       object
73   Fence          290 non-null     object
74   MiscFeature    51 non-null      object
75   MiscVal        1459 non-null    int64
76   MoSold         1459 non-null    int64
77   YrSold         1459 non-null    int64
78   SaleType       1458 non-null    object
79   SaleCondition  1459 non-null    object
dtypes: float64(11), int64(26), object(43)
memory usage: 912.0+ KB
```

```python
df_train['SalePrice'].describe()
```

```
count      1460.000000
mean     180921.195890
std       79442.502883
```
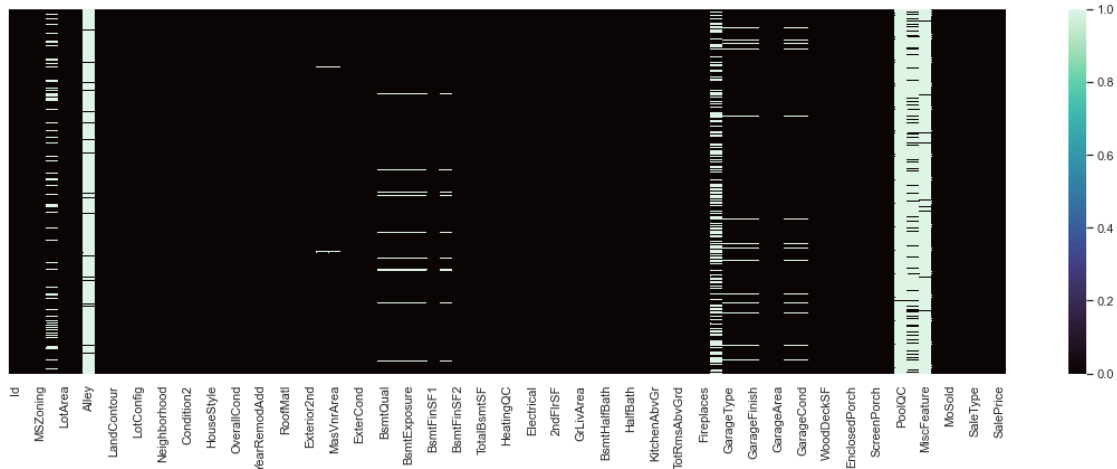
```
min          34900.000000
25%         129975.000000
50%         163000.000000
75%         214000.000000
max         755000.000000
Name: SalePrice, dtype: float64
```

```python
sns.distplot(df_train['SalePrice']);
print("Skewness: %f" % df_train['SalePrice'].skew())
print("Kurtosis: %f" % df_train['SalePrice'].kurt())
```

```
Skewness: 1.882876
Kurtosis: 6.536282
```



```python
plt.figure(figsize=(20,6))
sns.heatmap(df_train.isnull(),yticklabels=False,cbar=True,cmap='mako')
```

```
<AxesSubplot:>
```

```python
total_null = df_train.isnull().sum().sort_values(ascending=False)
#First sum and order all null values for each variable
percentage =
(df_train.isnull().sum()/df_train.isnull().count()).sort_values(ascend
ing=False) #Get the percentage
missing_data = pd.concat([total_null, percentage], axis=1,
keys=['Total', 'Percentage'])
missing_data.head(20)
```

|  | Total | Percentage |
|---|---|---|
| PoolQC | 1453 | 0.995205 |
| MiscFeature | 1406 | 0.963014 |
| Alley | 1369 | 0.937671 |
| Fence | 1179 | 0.807534 |
| FireplaceQu | 690 | 0.472603 |
| LotFrontage | 259 | 0.177397 |
| GarageYrBlt | 81 | 0.055479 |
| GarageCond | 81 | 0.055479 |
| GarageType | 81 | 0.055479 |
| 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 |
| MasVnrArea | 8 | 0.005479 |
| MasVnrType | 8 | 0.005479 |
| Electrical | 1 | 0.000685 |
| Id | 0 | 0.000000 |

```python
df_train = df_train.drop((missing_data[missing_data["Percentage"] >
0.05]).index,1) #Drop All Var. with null values > 1

df_train.isnull().sum()
```

```
Id                  0
MSSubClass          0
MSZoning            0
LotArea             0
Street              0
                   ..
MoSold              0
YrSold              0
SaleType            0
SaleCondition       0
SalePrice           0
Length: 70, dtype: int64
```

```python
num_col=df_train._get_numeric_data().columns.tolist()
num_col
```

```
['Id',
 'MSSubClass',
 'LotArea',
 'OverallQual',
 'OverallCond',
 'YearBuilt',
 'YearRemodAdd',
 'MasVnrArea',
 'BsmtFinSF1',
 'BsmtFinSF2',
 'BsmtUnfSF',
 'TotalBsmtSF',
 '1stFlrSF',
 '2ndFlrSF',
 'LowQualFinSF',
 'GrLivArea',
 'BsmtFullBath',
 'BsmtHalfBath',
 'FullBath',
 'HalfBath',
 'BedroomAbvGr',
 'KitchenAbvGr',
 'TotRmsAbvGrd',
 'Fireplaces',
 'GarageCars',
 'GarageArea',
 'WoodDeckSF',
 'OpenPorchSF',
 'EnclosedPorch',
 '3SsnPorch',
 'ScreenPorch',
 'PoolArea',
 'MiscVal',
 'MoSold',
```

```python
 'YrSold',
 'SalePrice']

cat_col=set(df_train.columns)-set(num_col)
cat_col

{'BldgType',
 'BsmtCond',
 'BsmtExposure',
 'BsmtFinType1',
 'BsmtFinType2',
 'BsmtQual',
 'CentralAir',
 'Condition1',
 'Condition2',
 'Electrical',
 'ExterCond',
 'ExterQual',
 'Exterior1st',
 'Exterior2nd',
 'Foundation',
 'Functional',
 'Heating',
 'HeatingQC',
 'HouseStyle',
 'KitchenQual',
 'LandContour',
 'LandSlope',
 'LotConfig',
 'LotShape',
 'MSZoning',
 'MasVnrType',
 'Neighborhood',
 'PavedDrive',
 'RoofMatl',
 'RoofStyle',
 'SaleCondition',
 'SaleType',
 'Street',
 'Utilities'}

for col in num_col:
    df_train[col].fillna(0, inplace=True)

for col in cat_col:

    df_train[col].fillna('None', inplace=True)

## NA Check: Verify that we covered all 'NAs' in our data
print(f'Number of NAs in train df: {sum(df_train.isnull().sum())}')
```

```
Number of NAs in train df: 0

plt.figure(figsize=(20,6))
sns.heatmap(df_train.isnull(),yticklabels=False,cbar=True,cmap='mako')

<AxesSubplot:>
```



## Investigate potential features & outliers

Below, We can see a few of the highest correlating predictors of SalePrice. Based on these features, it is obvious that usable square footage cumulatively amounts to the highest correlation to SalePrice (GrLivArea, TotalBsmtSF, 1stFlrSF, GarageArea). Other discrete and categorical variables (OverallQual, GarageCars, FullBath, TotRmsAdvGrd) influence the dependent variable as well.

```
corr_mat = df_train.corr().SalePrice.sort_values(ascending=False)
corr_mat.head(10)

SalePrice      1.000000
OverallQual    0.795774
GrLivArea      0.734968
TotalBsmtSF    0.651153
GarageCars     0.641047
1stFlrSF       0.631530
GarageArea     0.629217
FullBath       0.562165
TotRmsAbvGrd   0.537769
YearBuilt      0.523608
Name: SalePrice, dtype: float64
```

Below we can see the distribution of a few of these variables and assess how outliers may impact the data.
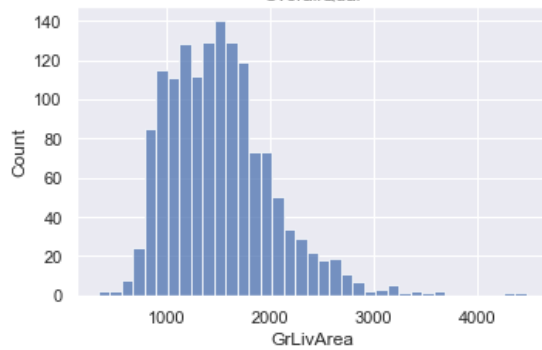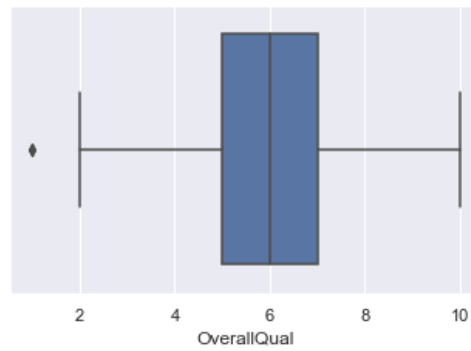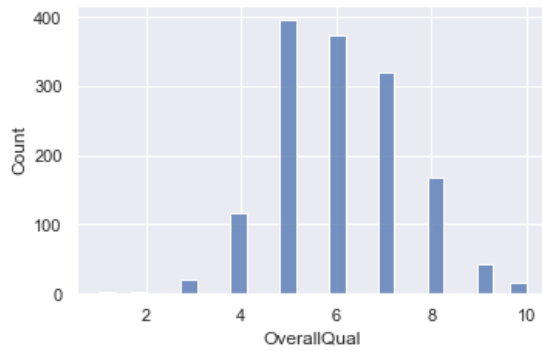
```
cor_features = ['OverallQual', 'GrLivArea', 'TotalBsmtSF',
'GarageCars', '1stFlrSF', 'YearBuilt' ]
```

```python
n = len(cor_features)

fig = plt.figure(figsize=(6*2, 4*n))
# add 2 graph for each column variable
gs = fig.add_gridspec(n, 2)
ax = [[fig.add_subplot(gs[i, j]) for j in range(2)] for i in range(n)]

for i in range(n):
    sns.histplot(x=cor_features[i], data=df_train, ax=ax[i][0])
    sns.boxplot(x=cor_features[i], data=df_train, ax=ax[i][1])

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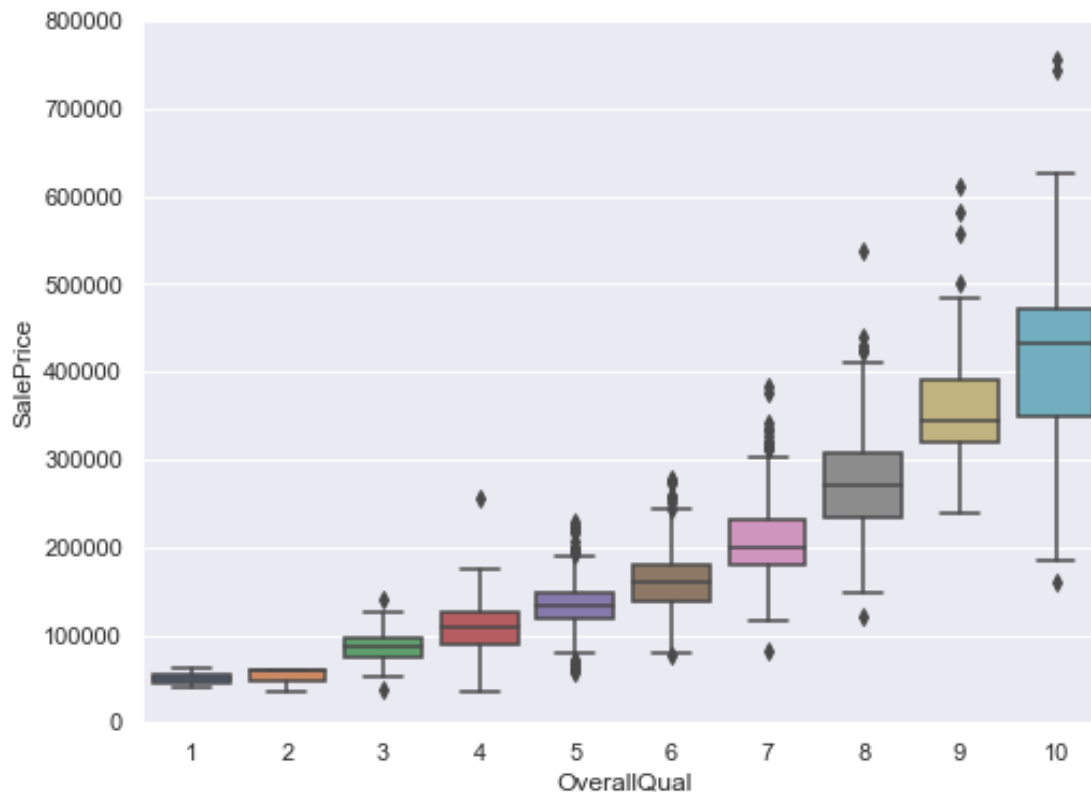
```python
# OverallQual and SalePrice
data = pd.concat([df_train['SalePrice'], df_train['OverallQual']], axis=1)
f, ax = plt.subplots(figsize=(8, 6))
fig = sns.boxplot(x='OverallQual', y="SalePrice", data=data)
fig.axis(ymin=0, ymax=800000);
```
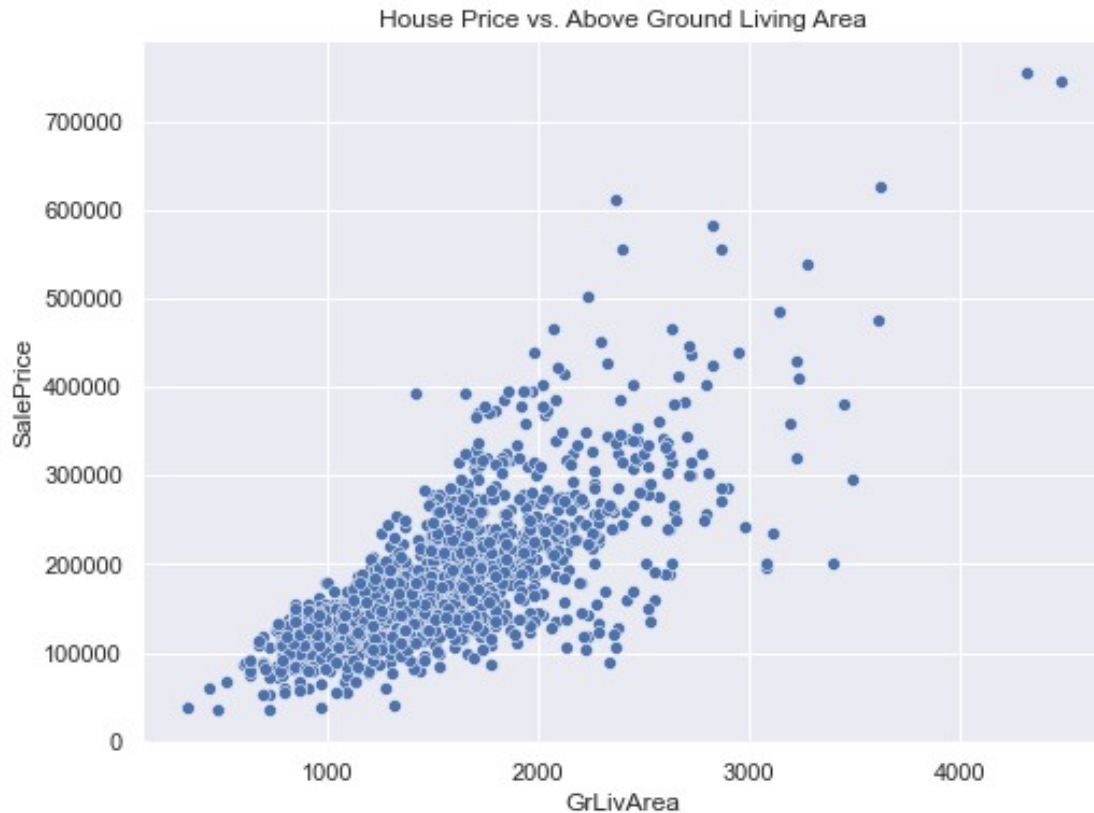


```python
# GrLivArea and SalesPrice
sns.set_style('darkgrid')
plt.figure(figsize=(8, 6))
sns.scatterplot(x='GrLivArea', y='SalePrice', data=df_train)
title = plt.title('House Price vs. Above Ground Living Area')
```

House Price vs. Above Ground Living Area

The scatter plot above reveals a few outliers where a larger living area is recorded with a low sale price. These outliers can be removed to ensure they do not influence future models.

```python
# Clean df_train (GrLiveArea)
outlier = df_train[(df_train.GrLivArea > 4000) & (df_train.SalePrice <
200000)].index
df_train.drop(outlier, axis=0, inplace=True)

# TotalBsmtSF and SalesPrice
sns.set_style('darkgrid')
plt.figure(figsize=(8, 6))
sns.scatterplot(x='TotalBsmtSF', y='SalePrice', data=df_train)
title = plt.title('House Price vs. Basement (sqft)')
```

House Price vs. Basement

```
# 1stFlrSF and SalesPrice
sns.set_style('darkgrid')
plt.figure(figsize=(8, 6))
sns.scatterplot(x='1stFlrSF', y='SalePrice', data=df_train)
title = plt.title('House Price vs. First Floor (sqft)')
```

House Price vs. First Floor (sqft)

## Feature Creation

Feature creation is likely to be a useful approach to finding more potent predictors in this data set. Based on the list of high correlating variables, it is apparent that features representing usable square feet are strong predictors and can be merged to create a stronger predictive feature. Additionally, the current dataframe seems to categorically discriminate based on above or below ground features. Combining some of high correlation variable, both above and below ground, may yield an overall stronger predictor. Finally, YearBuilt showed up a on the bottom of the correlation list with a comparatively low correlation. However, it remains an interesting feature to explore given some obvious and real world implications. Ideally, it would be nice to see in depth how larger renovations might impact the value of older homes. However, the data makes it difficult to define what renovation may have occurred.

Potentially interesting new predictors include:

-Total Square Feet of living Space (Below and Above ground)

-Total Number of Bathrooms (Below and Above Ground)

-Age of House when sold

```python
# Total Square Feet Column
df_train['TotalSqft'] = df_train['TotalBsmtSF'] + df_train['1stFlrSF']
+ df_train['2ndFlrSF']

# Total Bathrooms Column
df_train['TotalBath'] = df_train['FullBath'] +
df_train['BsmtFullBath'] + 0.5*(df_train['HalfBath'] +
df_train['BsmtHalfBath'])

# Age of House
df_train['HouseAge'] = df_train['YrSold'] - df_train['YearBuilt']

# Check for new columns
df_train.head()
```

```
    Id  MSSubClass MSZoning  LotArea Street LotShape LandContour
Utilities  \
0   1          60       RL     8450   Pave      Reg         Lvl
AllPub
1   2          20       RL     9600   Pave      Reg         Lvl
AllPub
2   3          60       RL    11250   Pave      IR1         Lvl
AllPub
3   4          70       RL     9550   Pave      IR1         Lvl
AllPub
4   5          60       RL    14260   Pave      IR1         Lvl
AllPub

  LotConfig LandSlope  ... PoolArea MiscVal MoSold YrSold SaleType  \
0    Inside       Gtl  ...        0       0      2   2008       WD
1       FR2       Gtl  ...        0       0      5   2007       WD
2    Inside       Gtl  ...        0       0      9   2008       WD
3    Corner       Gtl  ...        0       0      2   2006       WD
4       FR2       Gtl  ...        0       0     12   2008       WD

  SaleCondition  SalePrice  TotalSqft  HouseAge TotalBath
0        Normal     208500       2566         5       3.5
1        Normal     181500       2524        31       2.5
2        Normal     223500       2706         7       3.5
3       Abnorml     140000       2473        91       2.0
4        Normal     250000       3343         8       3.5

[5 rows x 73 columns]
```
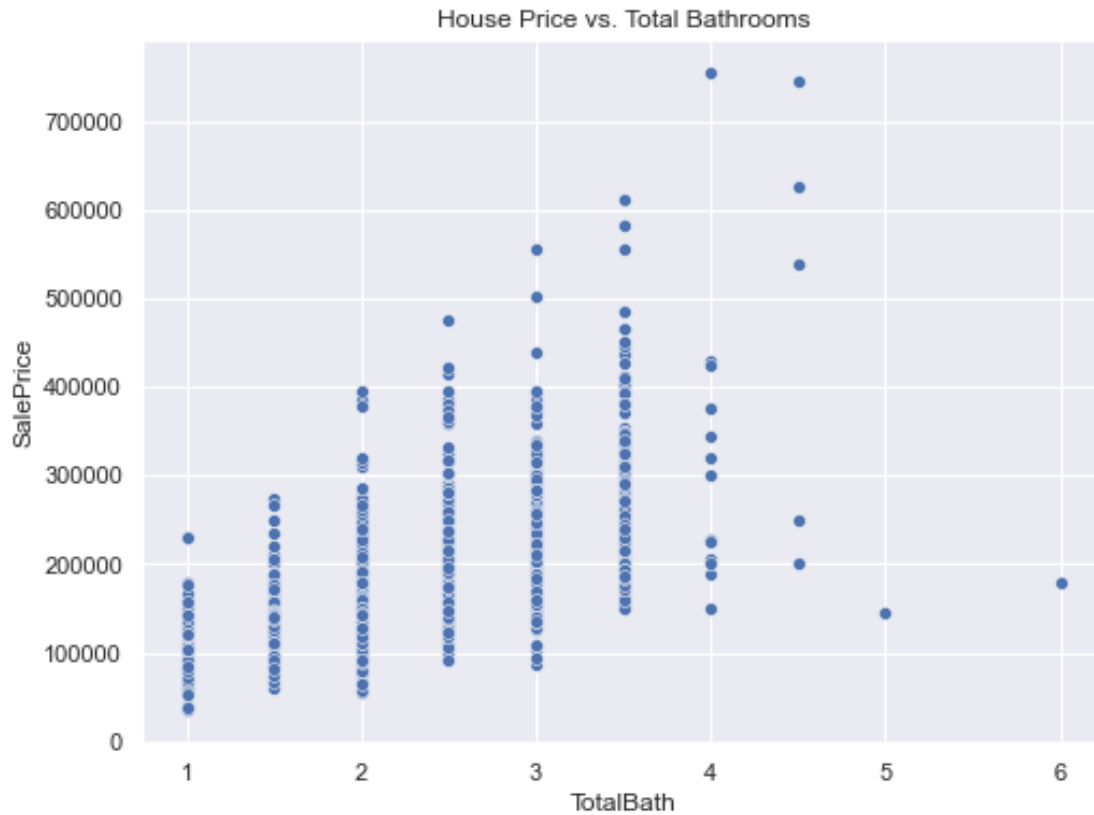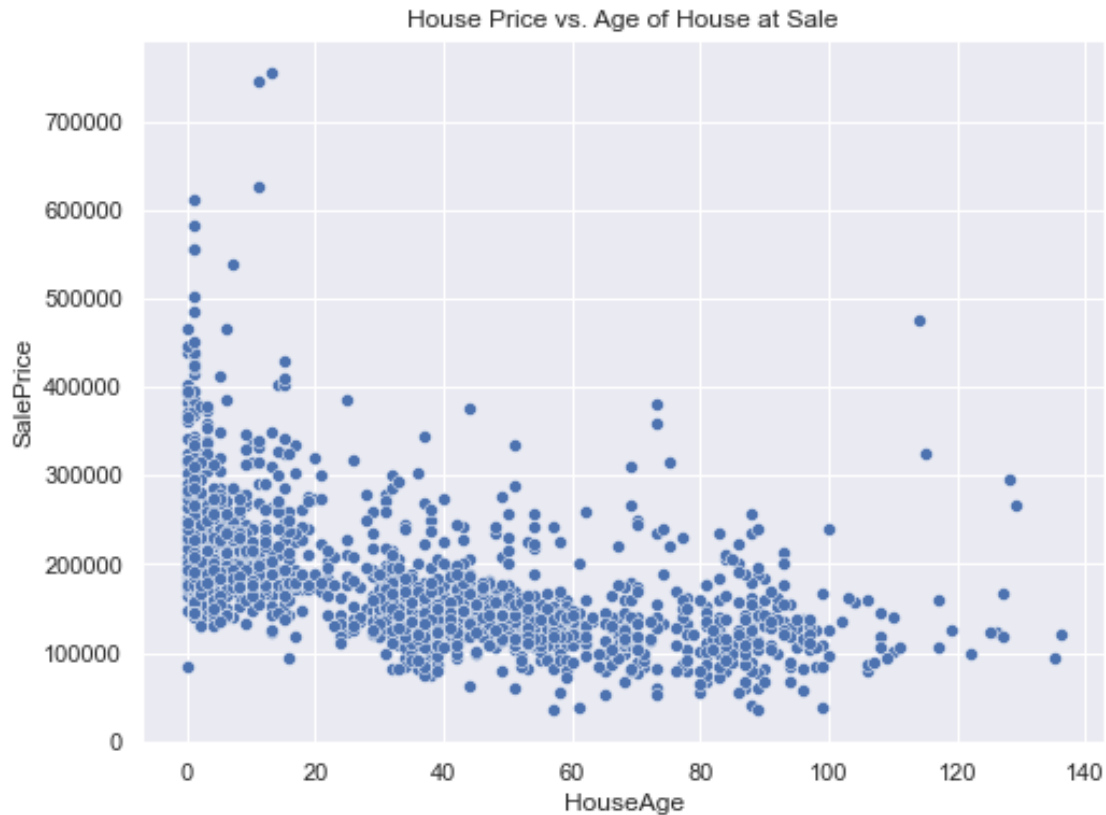
```python
# TotalSqft and SalesPrice
sns.set_style('darkgrid')
plt.figure(figsize=(8, 6))
sns.scatterplot(x='TotalSqft', y='SalePrice', data=df_train)
title = plt.title('House Price vs. Total Living Space')
```

House Price vs. Total Living Space

```
# TotalBath and SalesPrice
sns.set_style('darkgrid')
plt.figure(figsize=(8, 6))
sns.scatterplot(x='TotalBath', y='SalePrice', data=df_train)
title = plt.title('House Price vs. Total Bathrooms')
```

House Price vs. Total Bathrooms

```
# HouseAge and SalesPrice
sns.set_style('darkgrid')
plt.figure(figsize=(8, 6))
sns.scatterplot(x='HouseAge', y='SalePrice', data=df_train)
title = plt.title('House Price vs. Age of House at Sale')
```

House Price vs. Age of House at Sale

```
corr_mat2 = df_train.corr().SalePrice.sort_values(ascending=False)
corr_mat2.head(10)
```

```
SalePrice      1.000000
TotalSqft      0.832877
OverallQual    0.795774
GrLivArea      0.734968
TotalBsmtSF    0.651153
GarageCars     0.641047
TotalBath      0.635896
1stFlrSF       0.631530
GarageArea     0.629217
FullBath       0.562165
Name: SalePrice, dtype: float64
```

## Standard Scaling and Min-Max

The process of scaling is important for normalizing the data for a future model. We can see how the data of the chosen variables will normalize through both min-max and standard scaling methods.

```
x = df_train[['TotalSqft', 'OverallQual', 'TotalBsmtSF',
'1stFlrSF']].values
y = df_train['SalePrice'].values
```

```
fig, ax = plt.subplots(figsize=(12, 4))

ax.hist(x[:,0]);
ax.hist(x[:,1]);
ax.hist(x[:,2]);
ax.hist(x[:,3]);
```



```
fig, ax = plt.subplots(ncols=4, figsize=(24, 8))

ax[0].scatter(x[:,0], y);
ax[1].scatter(x[:,1], y);
ax[2].scatter(x[:,2], y);
ax[3].scatter(x[:,3], y);

plt.show()
```
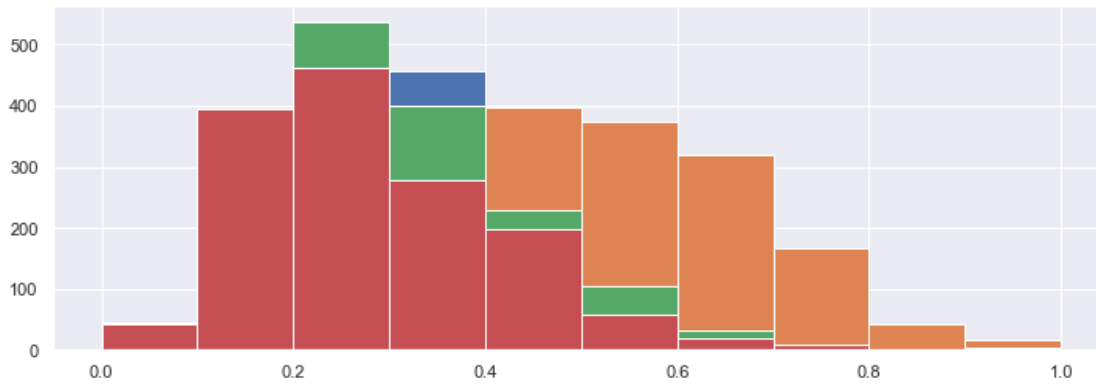


```
fig, ax = plt.subplots(figsize=(12, 4))

scaler = MinMaxScaler()
x_minmax = scaler.fit_transform(x)

ax.hist(x_minmax [:,0]);
ax.hist(x_minmax [:,1]);
ax.hist(x_minmax [:,2]);
ax.hist(x_minmax [:,3]);
```
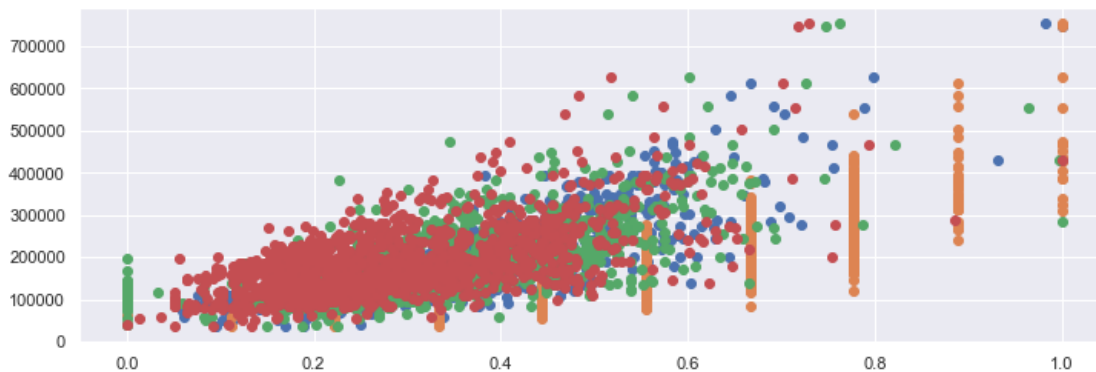
```python
fig, ax = plt.subplots(figsize=(12, 4))

scaler = MinMaxScaler()
x_minmax = scaler.fit_transform(x)

ax.scatter(x_minmax [:,0], y);
ax.scatter(x_minmax [:,1], y);
ax.scatter(x_minmax [:,2], y);
ax.scatter(x_minmax [:,3], y);
```
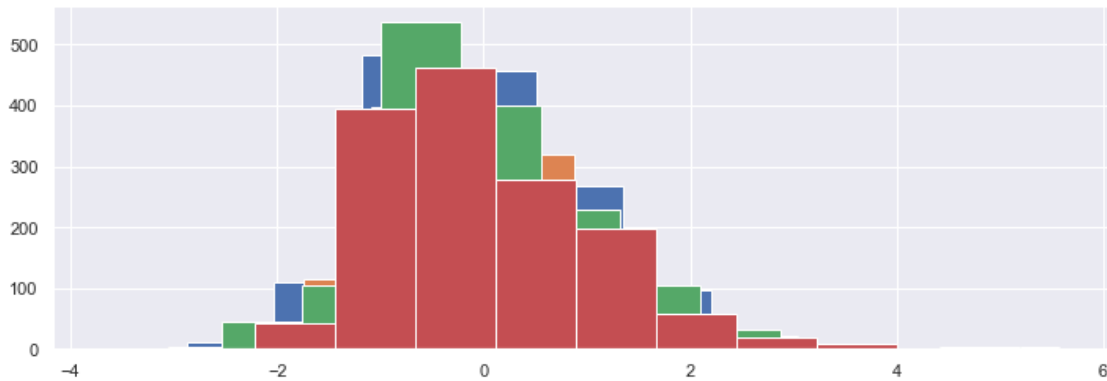


```python
fig, ax = plt.subplots(figsize=(12, 4))

scaler = StandardScaler()
x_std = scaler.fit_transform(x)

ax.hist(x_std[:,0]);
ax.hist(x_std[:,1]);
ax.hist(x_std[:,2]);
ax.hist(x_std[:,3]);
```

```
fig, ax = plt.subplots(figsize=(12, 4))

scaler = StandardScaler()
x_std = scaler.fit_transform(x)

ax.scatter(x_std[:,0], y);
ax.scatter(x_std[:,1], y);
ax.scatter(x_std[:,2], y);
ax.scatter(x_std[:,3], y);
```



## Conclusion

Conclusions drawn from performing EDA on this data set remain somewhat obvious in nature. Features that add value to the price of a home often include the most practical/usable features. Livable and usable space (sqft) drastically impacts the home value, while discrete and categorical variables regarding those livable/usable spaces add additional information to their utility and thus additional value (i.e. OverallQual, GarageCars, TotalBath). A notable point of interest that still warrants exploration is in the age of the house. This remains an interesting variable in relation to SalePrice because it hold a relatively meaningful level of correlation, yet it does not inherently communicate anything about the predictive variables regarding quality or livable/usable space. Ultimately, I believe it may be wise to consider more advanced feature creation surrounding house age and possible renovations before implementing a ML model.