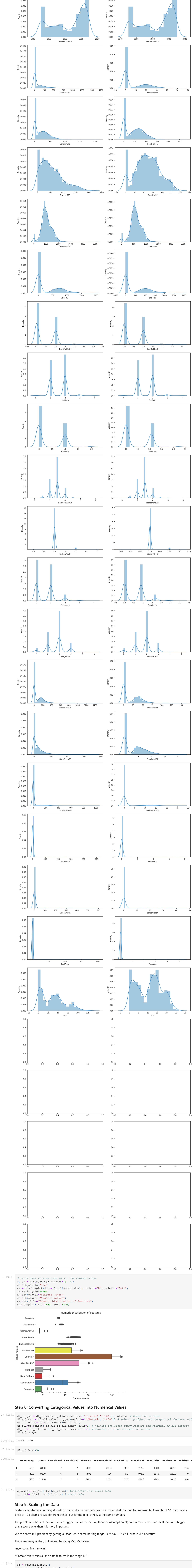


In [90]: # Normalize skewed features using a Box-Cox power transformation, we can use other techniques but am using boxcox
as it works very well on this dataset
for i in skew_index:
df_all[i] = boxcox1p(df_all[i], boxcox_normmax(df_all[i] + 1.002))

In [91]: df_all_num = df_all.select_dtypes(include=['int64', 'float64'])
skew_features = df_all_num.apply(lambda x: skew(x)).sort_values(ascending=False)
high_skew = skew_features[skew_features > 0.5]
skew_index = high_skew.index
skewness = pd.DataFrame({'Skew': high_skew})
skew_features

Out[91]: PoolArea 16.199734
3SeasonPorch 8.809104
KitchenAbvGr 3.778599
ScreenPorch 3.152066
EnclosedPorch 2.140485
MasVnrArea 0.981665
2ndFlrSF 0.887824
WoodDeckSF 0.789439
HalfBath 0.735962
BmnFlrSF 0.626280
OpenPorchSF 0.621761
Fireplaces 0.560109
OverallCond 0.381814
BmnFlrSF1 0.365664
BedroomAbvGr 0.333155
TotalBmnSF 0.273947
OverallQual 0.173972
FullBath 0.169433
OpenFrontage 0.059498
BmnUnfSF 0.052920
age -0.008823
LotArea -0.144950
GarageCars -0.251691
YearRemodAdd -0.445471
YearBuilt -0.380338
dtype: float64

In [92]: Rechecking if skewness is fixed
li_num_feats = list(num_feats)
nr_rows = 30
nr_cols = 2
fig, axs = plt.subplots(nr_rows, nr_cols, figsize=(nr_cols*7, nr_rows*4))
for r in range(nr_rows):
i = r
if i < len(li_num_feats):
sns.distplot(df_all_copy[li_num_feats[i]], ax = axs[r][0])
sns.distplot(df_all[li_num_feats[i]], ax = axs[r][1])
plt.tight_layout()



In [93]: # Let's make sure we handled all the skewed values
fig = plt.subplots(figsize=(8, 7))
ax = sns.kdeplot(log=True)
ax = sns.boxplot(data=df_all[skew_index], orient='h', palette='Set1')
ax.xaxis.grid(True)
ax.set(ylabel='Feature names')
ax.set(xlabel='Numeric values')
sns.boxplot(trim=True, label=True)

Step 8: Converting Categorical Values into Numerical Values

In [169]: df_all_num = df_all.select_dtypes(include=['float64', 'int64']).columns # Numerical columns
df_all_cat = df_all.select_dtypes(exclude=['float64', 'int64']) # selecting object and categorical features only
df_all_dummies = pd.get_dummies(df_all_cat)
df_all_concat = (df_all.concat(df_all_dummies, axis=1)) # joining converted dummy feature and original df_all dataset
df_all = df_all.drop(df_all_cat.columns, axis=1) # removing original categorical columns

Out[169]: (2919, 319)

In [171]: df_all.head(3)

Out[171]:

In [172]: x_train = df_all[~df_all['train']] # converted into train data
x_test = df_all[df_all['train']] # test data

Step 9: Scaling the Data

Scaler class: Machine learning algorithm that works on numbers does not know what that number represents. A weight of 10 grams and a price of 10 dollars are two different things, but for mode it is just the same numbers.

The problem is that if 1 feature is much bigger than other feature, then the assumption algorithm makes that since first feature is bigger than second one, then it is more important.

We can solve this problem by getting all features in same not big range. Let's say $-1 \leq x_i \leq 1$, where x_i is a feature

There are many scalers, but we will be using Min-Max scaler.

$x_{new} = \frac{x - x_{min}}{x_{max} - x_{min}}$

MinMaxScaler scales all the data features in the range [0,1]

In [175]: sc = StandardScaler()
x_train = sc.fit_transform(x_train)
x_test = sc.transform(x_test)

Step 10: Training Models on this data.

A. Linear Regression

In [123]: lr = LinearRegression().fit(x_train, y_train)

In [128]: r_sq = lr.score(x_train, y_train)
print('coefficient of determination - r-square - on training data :', r_sq)

coefficient of determination - r-square - on training data : 0.9558501439818528

In [129]: macrof1 = np.argmax(np.abs(Rs_coef))
coef = lr.coef_[macrof1]
for i in range(0, 5):
print("{}({:.02f}) <= 0.10.{}e".format(df_all.columns[macrof1[i]], coef[i]))

TotalBmnSF..... 1.5317e+10
RoomMat_ClyFile..... -2.1815e+10
2ndFlrSF..... 9.8110e-02
OverallQual..... 6.0376e-02
GarageCars..... 2.8345e-02

In [120]: from sklearn.linear_model import LinearRegression, RidgeCV, LassoCV, ElasticNetCV
from sklearn.model_selection import cross_val_score

B. Linear regression, L1 regularisation

In [121]: # Create linear regression object
ls = LassoCV()

Train the model using the training sets
ls.fit(x_train, y_train)

Out[121]: LassoCV()


```
In [224]: r_sq = ls.score(x_train, y_train)
print('coefficient of determination - rsquare - on training data :', r_sq)

coefficient of determination - rsquare - on training data : 0.939325943077308
```

```
In [225]: maxcoef = np.argsort(-np.abs(ls.coef_))
coef = ls.coef_[maxcoef]
for i in range(0, 5):
    print("{:0.25} {:< 010.4e}".format(df_all.columns[maxcoef[i]], coef[i]))

OverallQual..... 9.7745e-02
TotalBsmntSF..... 9.3286e-02
2ndFlrSF..... 5.7333e-02
Condition2_PosN..... -4.8166e-02
GarageCars..... 4.2237e-02
```

C. Linear regression, L2 regularisation!

```
In [226]: # Create linear regression object
Rr = RidgeCV()

# Train the model using the training sets
Rr.fit(x_train, y_train)
```

```
Out[226]: RidgeCV(alpha=array([ 0.1, 1. , 10. ]))
```

```
In [227]: r_sq = Rr.score(x_train, y_train)
print('coefficient of determination - rsquare - on training data :', r_sq)

coefficient of determination - rsquare - on training data : 0.9548127192732505
```

```
In [228]: maxcoef = np.argsort(-np.abs(Rr.coef_))
coef = Rr.coef_[maxcoef]
for i in range(0, 5):
    print("{:0.25} {:< 010.4e}".format(df_all.columns[maxcoef[i]], coef[i]))

TotalBsmntSF..... 9.5422e-02
2ndFlrSF..... 6.7126e-02
OverallQual..... 5.8948e-02
Condition2_PosN..... -5.2500e-02
OverallCond..... 4.6214e-02
```