

Importing Libraries and Dataset

```
In [4]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

dataset = pd.read_csv("HousePricePrediction.xlsx - Sheet1.csv")

# Printing first 5 records of the dataset
print(dataset.head(5))
```

	Id	MSSubClass	MSZoning	LotArea	LotConfig	BldgType	OverallCond	\
0	0	60	RL	8450	Inside	1Fam	5	
1	1	20	RL	9600	FR2	1Fam	8	
2	2	60	RL	11250	Inside	1Fam	5	
3	3	70	RL	9550	Corner	1Fam	5	
4	4	60	RL	14260	FR2	1Fam	5	

	YearBuilt	YearRemodAdd	Exterior1st	BsmtFinSF2	TotalBsmtSF	SalePrice
0	2003	2003	VinylSd	0.0	856.0	208500.0
1	1976	1976	MetalSd	0.0	1262.0	181500.0
2	2001	2002	VinylSd	0.0	920.0	223500.0
3	1915	1970	Wd Sdng	0.0	756.0	140000.0
4	2000	2000	VinylSd	0.0	1145.0	250000.0

```
In [5]: dataset.shape
```

```
Out[5]: (2919, 13)
```

```
In [ ]:
```

Data Preprocessing

```
In [6]: obj = (dataset.dtypes == 'object')
object_cols = list(obj[obj].index)
print("Categorical variables:", len(object_cols))

int_ = (dataset.dtypes == 'int')
num_cols = list(int_[int_].index)
print("Integer variables:", len(num_cols))

fl = (dataset.dtypes == 'float')
fl_cols = list(fl[fl].index)
print("Float variables:", len(fl_cols))
```

```
Categorical variables: 4
Integer variables: 0
Float variables: 3
```

Exploratory Data Analysis

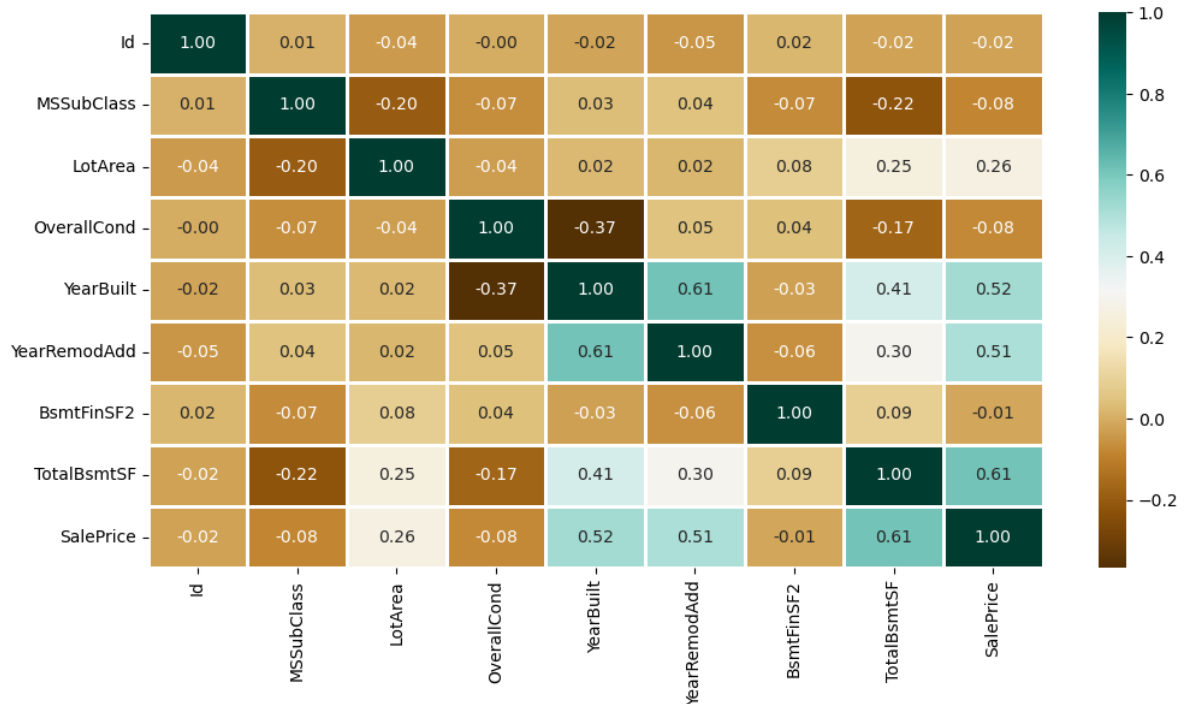
```
In [7]: plt.figure(figsize=(12, 6))
sns.heatmap(dataset.corr(),
             cmap = 'BrBG',
             fmt = '.2f',
```

```
linewidths = 2,
annot = True)
```

C:\Users\venut\AppData\Local\Temp\ipykernel_13384\3487798585.py:2: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.

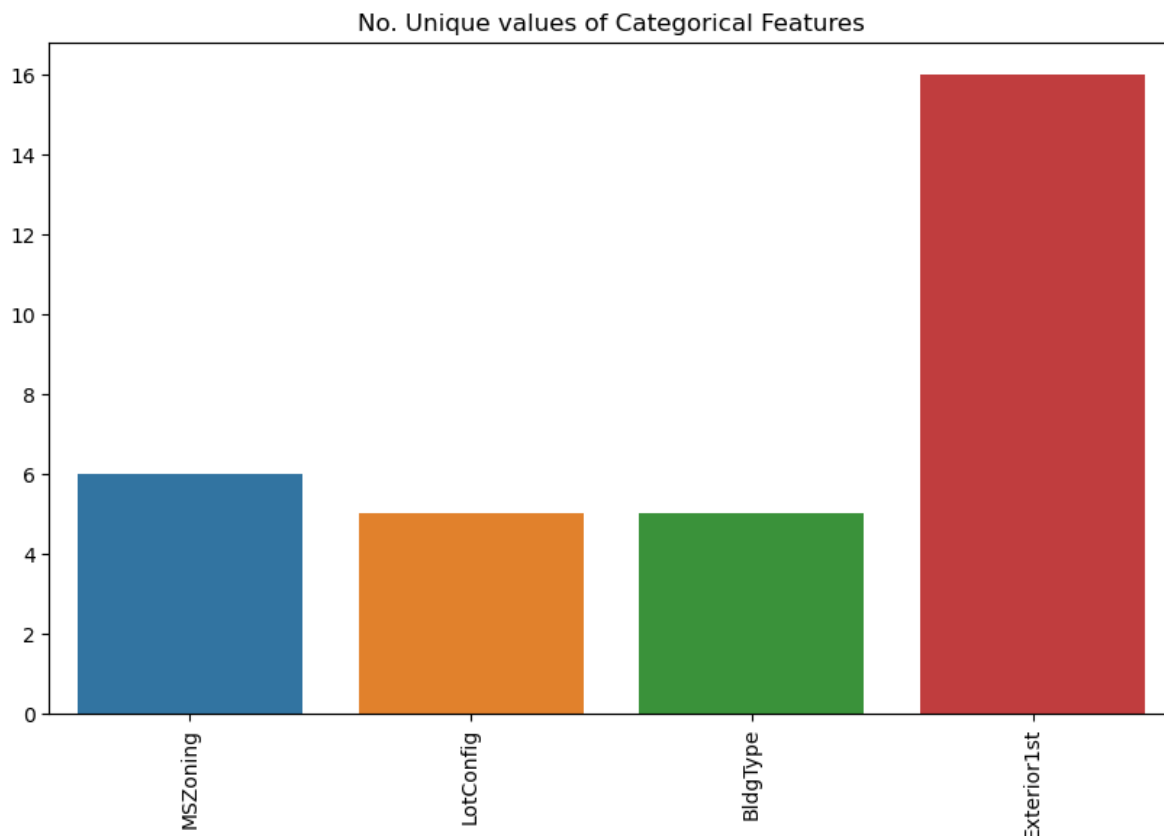
```
sns.heatmap(dataset.corr(),
```

Out[7]: <Axes: >



```
In [8]: unique_values = []
for col in object_cols:
    unique_values.append(dataset[col].unique().size)
plt.figure(figsize=(10,6))
plt.title('No. Unique values of Categorical Features')
plt.xticks(rotation=90)
sns.barplot(x=object_cols,y=unique_values)
```

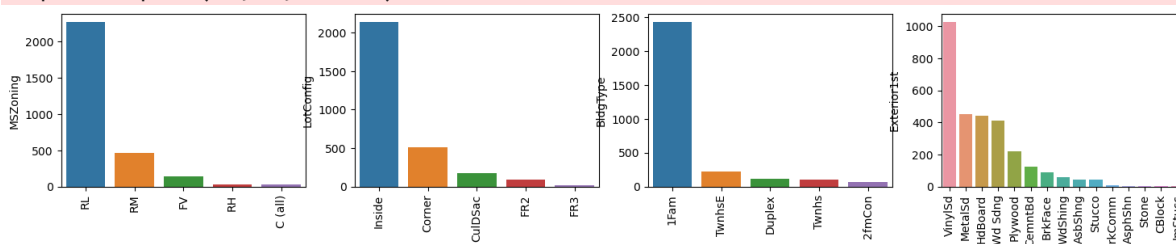
Out[8]: <Axes: title={'center': 'No. Unique values of Categorical Features'}>



```
In [9]: plt.figure(figsize=(18, 36))
plt.title('Categorical Features: Distribution')
plt.xticks(rotation=90)
index = 1

for col in object_cols:
    y = dataset[col].value_counts()
    plt.subplot(11, 4, index)
    plt.xticks(rotation=90)
    sns.barplot(x=list(y.index), y=y)
    index += 1
```

C:\Users\venut\AppData\Local\Temp\ipykernel_13384\2166598984.py:8: MatplotlibDeprecationWarning: Auto-removal of overlapping axes is deprecated since 3.6 and will be removed two minor releases later; explicitly call ax.remove() as needed.



Data Cleaning

```
In [10]: dataset.drop(['Id'],
axis=1,
inplace=True)
```

```
In [11]: dataset['SalePrice'] = dataset['SalePrice'].fillna(
dataset['SalePrice'].mean())
```

```
In [12]: new_dataset = dataset.dropna()
```

```
In [13]: new_dataset.isnull().sum()
```

```
Out[13]: MSSubClass      0
         MSZoning        0
         LotArea         0
         LotConfig       0
         BldgType        0
         OverallCond     0
         YearBuilt       0
         YearRemodAdd    0
         Exterior1st     0
         BsmtFinSF2      0
         TotalBsmtSF     0
         SalePrice       0
         dtype: int64
```

OneHotEncoder – For Label categorical features

```
In [14]: from sklearn.preprocessing import OneHotEncoder
```

```
s = (new_dataset.dtypes == 'object')
object_cols = list(s[s].index)
print("Categorical variables:")
print(object_cols)
print('No. of. categorical features: ',
      len(object_cols))
```

```
Categorical variables:
['MSZoning', 'LotConfig', 'BldgType', 'Exterior1st']
No. of. categorical features: 4
```

```
In [23]: OH_encoder = OneHotEncoder(sparse=False, handle_unknown='ignore')
         OH_cols = pd.DataFrame(OH_encoder.fit_transform(new_dataset[object_cols]))
         OH_cols.index = new_dataset.index
         OH_cols.columns = OH_encoder.get_feature_names_out()
         df_final = new_dataset.drop(object_cols, axis=1)
         df_final = pd.concat([df_final, OH_cols], axis=1)
```

```
C:\Users\venut\anaconda3\lib\site-packages\sklearn\preprocessing\_encoders.py:828:
FutureWarning: `sparse` was renamed to `sparse_output` in version 1.2 and will be
removed in 1.4. `sparse_output` is ignored unless you leave `sparse` to its default
value.
  warnings.warn(
```

```
In [18]: df_final.
```

```
Out[18]: MSSubClass      0
          LotArea        0
          OverallCond    0
          YearBuilt      0
          YearRemodAdd    0
          BsmtFinSF2     0
          TotalBsmtSF    0
          SalePrice      0
          0              0
          1              0
          2              0
          3              0
          4              0
          5              0
          6              0
          7              0
          8              0
          9              0
          10             0
          11             0
          12             0
          13             0
          14             0
          15             0
          16             0
          17             0
          18             0
          19             0
          20             0
          21             0
          22             0
          23             0
          24             0
          25             0
          26             0
          27             0
          28             0
          29             0
dtype: int64
```

Splitting Dataset into Training and Testing

```
In [24]: from sklearn.metrics import mean_absolute_error
          from sklearn.model_selection import train_test_split

          X = df_final.drop(['SalePrice'], axis=1)
          Y = df_final['SalePrice']

          # Split the training set into
          # training and validation set
          X_train, X_valid, Y_train, Y_valid = train_test_split(
              X, Y, train_size=0.8, test_size=0.2, random_state=0)
```

SVM – Support vector Machine

```
In [25]: from sklearn import svm
          from sklearn.svm import SVC
          from sklearn.metrics import mean_absolute_percentage_error
```

```
model_SVR = svm.SVR()
model_SVR.fit(X_train,Y_train)
Y_pred = model_SVR.predict(X_valid)

print(mean_absolute_percentage_error(Y_valid, Y_pred))

0.1870512931870423
```