```
In [85]:    1  #loading required libraries
            2  import numpy as np
            3  import pandas as pd
            4  import matplotlib.pyplot as plt
            5  from sklearn.linear_model import LinearRegression
            6  from sklearn.metrics import mean_squared_error
            7  from sklearn.impute import SimpleImputer
            8  from sklearn.ensemble import HistGradientBoostingRegressor
            9  import matplotlib.pyplot as plt
           10  import seaborn as sns
           11  from sklearn.model_selection import cross_val_score
           12  import seaborn as sns
           13  import statsmodels.api as sm
           14
```
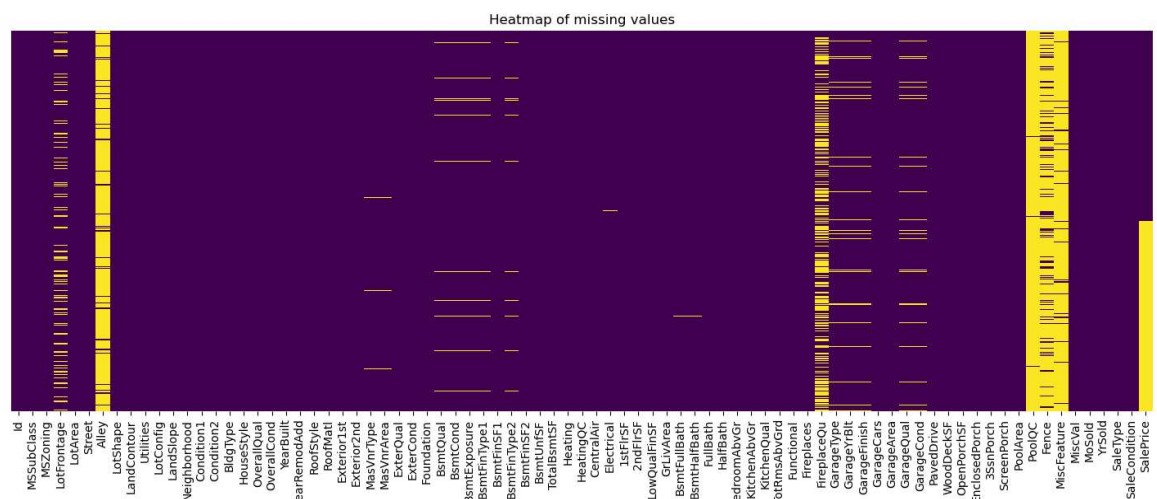
```
In [86]:    1  train_data = pd.read_csv('train.csv')
            2  test_data=pd.read_csv('test.csv')
            3  data=pd.concat([train_data,test_data])
            4  plt.figure(figsize=(18,6))
            5  plt.title('Heatmap of missing values')
            6  sns.heatmap(data.isnull(),yticklabels=False,cbar=False,cmap='viridis')
            7
```

Out[86]:  <Axes: title={'center': 'Heatmap of missing values'}>

```python
In [87]:   1  # separate numerical and categorical columns
           2  numerical_cols = train_data.select_dtypes(include=['int64', 'float64'])
           3  categorical_cols = train_data.select_dtypes(include=['object']).columns
           4
           5  # replace missing values with mean for numerical columns
           6  imputer = SimpleImputer(strategy='mean')
           7  train_data[numerical_cols] = imputer.fit_transform(train_data[numerical_
           8
           9  # replace missing values with mode for categorical columns
          10  imputer = SimpleImputer(strategy='most_frequent')
          11  train_data[categorical_cols] = imputer.fit_transform(train_data[categori
          12
          13  # verify that missing values are replaced
          14  print(train_data.isnull().sum())
          15
          16
          17
```

```
Id               0
MSSubClass       0
MSZoning         0
LotFrontage      0
LotArea          0
                ..
MoSold           0
YrSold           0
SaleType         0
SaleCondition    0
SalePrice        0
Length: 81, dtype: int64
```

```python
In [88]:   1  # Select the predictor and target variable (for both training and testi
           2  x_train = train_data[['LotArea', 'BedroomAbvGr', 'BsmtFullBath','BsmtHa]
           3  y_train = train_data['SalePrice']
           4
           5  x_test=test_data[['LotArea', 'BedroomAbvGr', 'BsmtFullBath','BsmtHalfBat
```

```
In [89]:    1  # separate numerical and categorical columns for test data
            2  numerical_cols2 = x_test.select_dtypes(include=['int64', 'float64']).col
            3
            4  # replace missing values with mean for numerical columns
            5  imputer = SimpleImputer(strategy='mean')
            6  x_test[numerical_cols2] = imputer.fit_transform(x_test[numerical_cols2]
            7
            8  # verify that missing values are replaced
            9  print(x_test.isnull().sum())
```

```
LotArea          0
BedroomAbvGr     0
BsmtFullBath     0
BsmtHalfBath     0
FullBath         0
HalfBath         0
dtype: int64

C:\Users\Lenovo\AppData\Local\Temp\ipykernel_8988\281113641.py:6: SettingWi
thCopyWarning:
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: https://pandas.pydata.org/pandas-doc
s/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://p
andas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-vi
ew-versus-a-copy)
  x_test[numerical_cols2] = imputer.fit_transform(x_test[numerical_cols2])
```

```
In [90]:    1  # Initializing a linear regression model and fit this model to train se
            2  model = LinearRegression()
            3  print(model)
            4  model.fit(x_train, y_train)
            5
```

```
LinearRegression()
```

```
Out[90]:   ▼ LinearRegression
           LinearRegression()
```

```
In [91]:    1  # Making predictions on the test set using the fitted model
            2
            3  y_pred = model.predict(x_test)
            4  print(y_pred)
            5
```
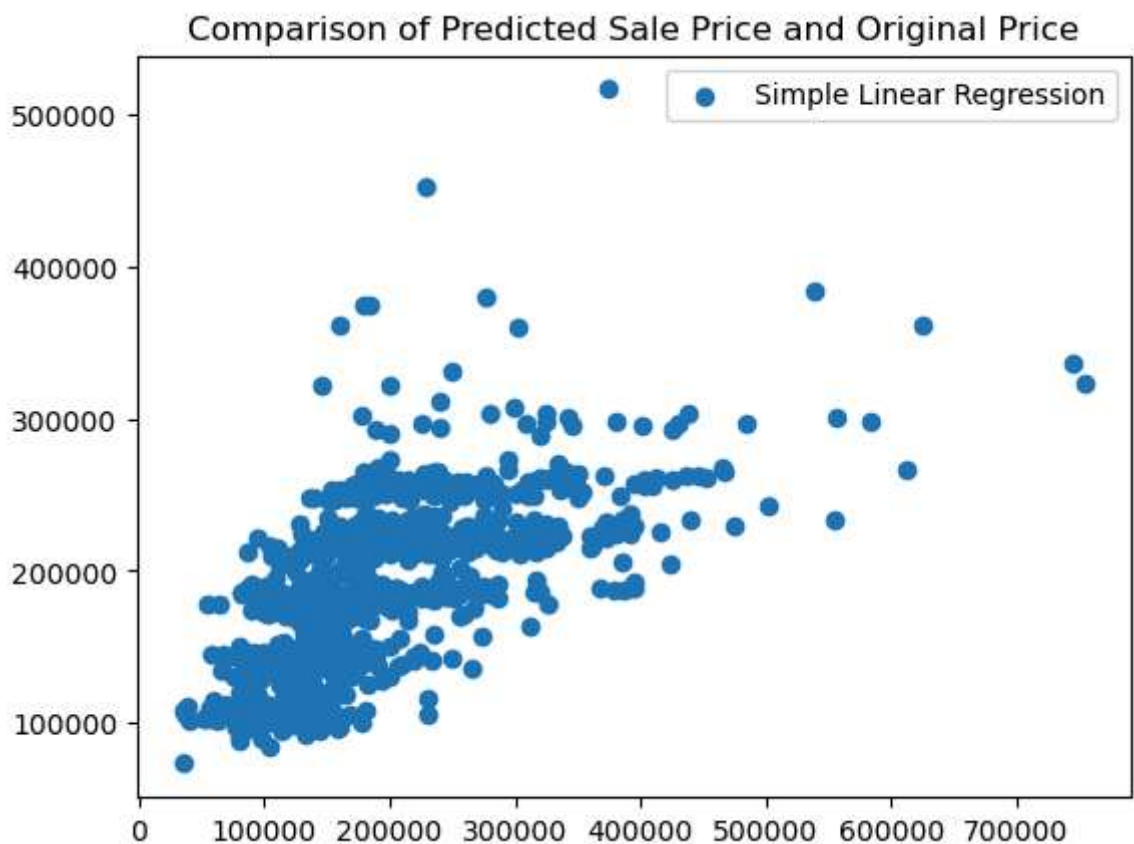
```
[112606.119798    145550.48551853 224695.03092123 ... 146582.41859758
 120117.45460539 219383.04414268]
```

```
In [92]:    1  #instead of simple linear regression we are using HistGradientBoosting
            2  #missing values are assigned to the left or right child accordingly1. Th
```

```
In [81]:   1  train_data = pd.read_csv('train.csv')
           2  test_data=pd.read_csv('test.csv')
           3
           4  x_train = train_data[['LotArea', 'BedroomAbvGr', 'BsmtFullBath','BsmtHa
           5  y_train = train_data['SalePrice']
           6
           7  x_test=test_data[['LotArea', 'BedroomAbvGr', 'BsmtFullBath','BsmtHalfBat
           8  # Initializing a linear regression model and fit this model to train se
           9  histgrad_model =  HistGradientBoostingRegressor().fit(x_train, y_train)
          10
          11  # Making predictions on the test set using the fitted model
          12  y_pred = histgrad_model.predict(x_test)
          13  print("predicted saleprices;",y_pred)
          14
```

predicted saleprices; [129151.88642225 178516.40056678 253716.44635059 ...
199879.47712159
 137533.8702215   208269.17549794]

```
In [82]:   1  plt.figure()
           2  plt.title('Comparison of Predicted Sale Price and Original Price')
           3  plt.scatter(y_train, model.predict(x_train), label='Simple Linear Regres
           4  plt.legend()
           5  plt.show()
           6
```

```
In [93]:  1  plt.figure()
          2  plt.title('Comparison of Predicted Sale Price and Original Price')
          3  plt.scatter(y_train, histgrad_model.predict(x_train), label='HistGradie
          4  plt.legend()
          5  plt.show()
          6
```

### Comparison of Predicted Sale Price and Original Price



```
In [96]:  1  #preparing the submission data
          2  y_pred = model.predict(x_test)
          3  sub_data=pd.DataFrame()
          4  sub_data['Id']=test_data['Id']
          5  sub_data['SalePrice']=y_pred
          6  sub_data
          7  sub_data.to_excel('submission.xlsx', index=False)
          8
```

```python
###LINEAR REGRESSION USING ORDINARY LEAST SQUARE METHOD

data = pd.read_csv('train.csv')

# Select the relevant features: square footage, number of bedrooms, and
X1 = data[['LotArea', 'BedroomAbvGr', 'BsmtFullBath','BsmtHalfBath','Ful
y1 = data['SalePrice']  # Target variable (house price)

# Add a constant term for the intercept in the regression model
X1 = sm.add_constant(X1)

# Fit the linear regression model
model = sm.OLS(y1, X1).fit()
model
# Print the summary of the model
print(model.summary())
```

```
                            OLS Regression Results
==================================================================================
Dep. Variable:                 SalePrice   R-squared:                       0.459
Model:                               OLS   Adj. R-squared:                  0.457
Method:                    Least Squares   F-statistic:                     205.9
Date:                   Tue, 04 Jun 2024   Prob (F-statistic):           4.28e-190
Time:                           03:01:15   Log-Likelihood:                -18095.
No. Observations:                   1460   AIC:                         3.620e+04
Df Residuals:                       1453   BIC:                         3.624e+04
Df Model:                              6
Covariance Type:               nonrobust
==================================================================================
                   coef     std err          t      P>|t|      [0.025      0.975]
----------------------------------------------------------------------------------
const           3.178e+04    6555.902      4.847      0.000    1.89e+04    4.46e+04
LotArea            1.2639       0.158      7.981      0.000       0.953       1.575
BedroomAbvGr   -6777.5551    2087.330     -3.247      0.001   -1.09e+04   -2683.053
BsmtFullBath    3.694e+04    3074.859     12.014      0.000    3.09e+04     4.3e+04
BsmtHalfBath    1.578e+04    6532.302      2.416      0.016    2967.754    2.86e+04
FullBath          7.97e+04    3014.745     26.436      0.000    7.38e+04    8.56e+04
HalfBath        3.638e+04    3134.939     11.604      0.000    3.02e+04    4.25e+04
==================================================================================
Omnibus:                      539.540   Durbin-Watson:                   1.975
Prob(Omnibus):                  0.000   Jarque-Bera (JB):             3595.355
Skew:                           1.564   Prob(JB):                         0.00
Kurtosis:                      10.023   Cond. No.                     6.70e+04
==================================================================================

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 6.7e+04. This might indicate that there are
strong multicollinearity or other numerical problems.
```

```
In [ ]: 1
```