

```
import numpy as np
import seaborn as sns
from sklearn.preprocessing import StandardScaler
import pandas as pd
import statsmodels
import matplotlib
plt.style.use('ggplot')
```

```
config InlineBackend.figure_format = 'retina'
%matplotlib inline
```

```
import statsmodels.api as sm
from sklearn.metrics import mean_squared_error, r2_score
from sklearn.cross_validation import cross_val_score
from sklearn.linear_model import LinearRegression, Ridge, Lasso, RidgeCV, LassoCV
from statsmodels.graphics.goregplots import qqplot
from scipy import stats
from statsmodels.stats import shapiro
```

```
file_location = 'data/Ames_Housing_Sales.csv'
```

```
df = pd.read_csv(file_location)
df = pd.DataFrame(df)
```

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2617 entries, 0 to 2616
Data columns (total 38 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   SalePrice            2617 non-null   float64
 1   LotArea              2617 non-null   float64
 2   Age                  2617 non-null   float64
 3   GrLivArea            2617 non-null   float64
 4   BasementArea         2617 non-null   float64
 5   Location             2617 non-null   int64
 6   Amenities            2617 non-null   int64
 7   RoadFrontage         2617 non-null   float64
 8   BedroomAbvGr         2617 non-null   float64
 9   Bathrooms            2617 non-null   float64
10  OverallQual          2617 non-null   float64
11  OverallCond          2617 non-null   float64
12  LotFrontage          2617 non-null   float64
13  LotArea              2617 non-null   float64
14  TwoStory_dum         2617 non-null   int64
15  FlatContour_dum      2617 non-null   int64
16  FlatRoof_dum         2617 non-null   float64
17  GarageArea           2617 non-null   float64
18  Garage_dum           2617 non-null   int64
19  CentralAir_dum       2617 non-null   int64
20  LowQualFinSF         2617 non-null   float64
21  Fireplaces           2617 non-null   int64
22  KitchenQual_Lev      2617 non-null   int64
23  Train_Lev            2617 non-null   int64
24  Zoning_3             2617 non-null   int64
25  Zoning_4             2617 non-null   int64
26  YrSold_2007          2617 non-null   int64
27  YrSold_2008          2617 non-null   int64
28  YrSold_2009          2617 non-null   int64
29  YrSold_2010          2617 non-null   int64
dtypes: float64(14), int64(16)
memory usage: 652.5 KB
```

```
# Deleting the single null value in 'GarageArea'
df = df[df['GarageArea'].isnull()]
```

```
df.isnull().values.any()
```

```
False
```

```
fig, ax = plt.subplots(figsize=(10, 10))
```

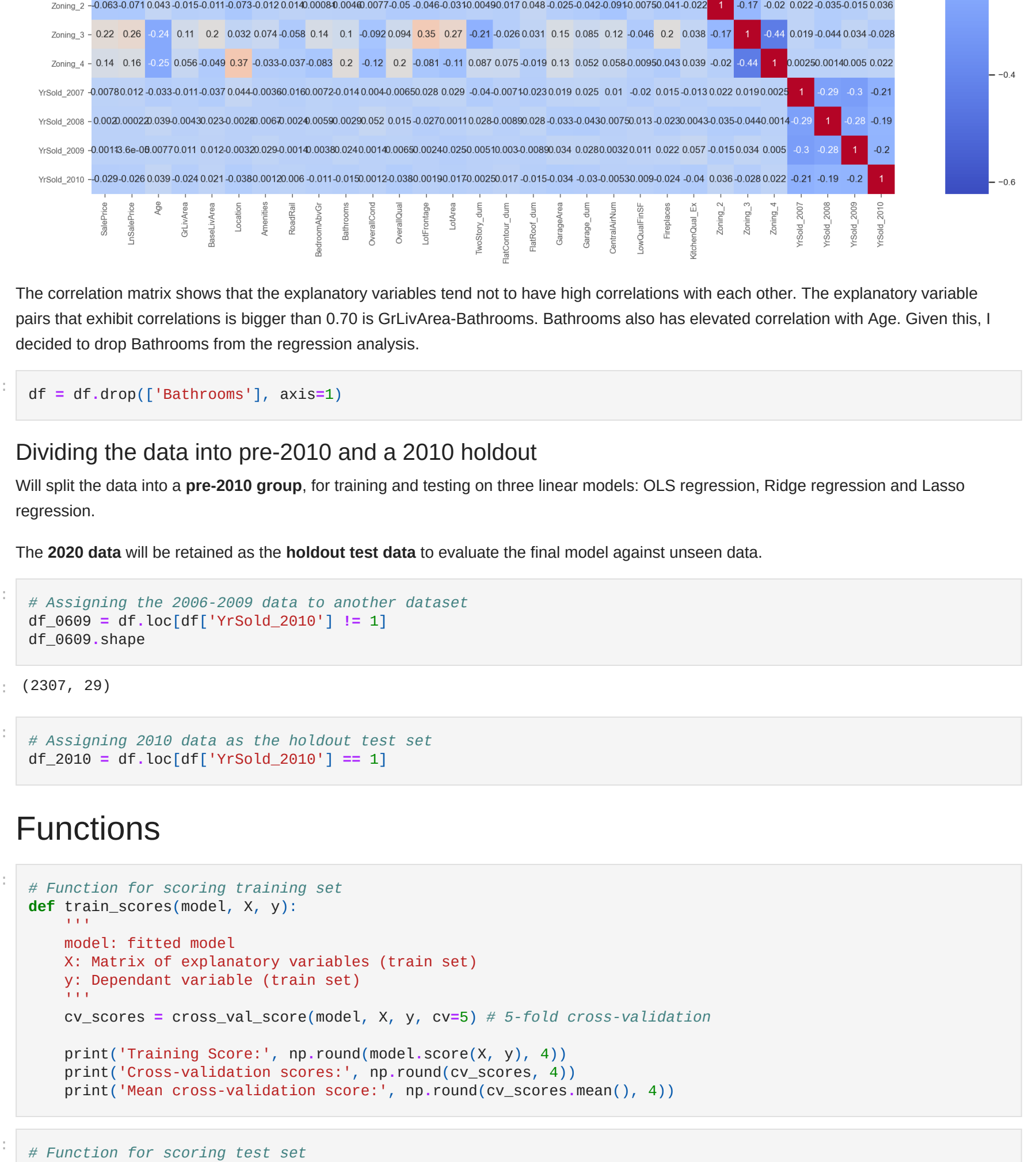
```
sns.set(font_scale=0.8)
sns.heatmap(df.corr(), annot=True, cmap='coolwarm', ax=ax)
```

```
ax.set_title('Correlation Matrix of Ames Variables', fontsize=16)
```

```
ax.set_xticklabels(ax.get_major tick labels(), fontsize=8)
```

```
ax.set_yticklabels(ax.get_major tick labels(), fontsize=8)
```

```
plt.show()
```



The correlation matrix shows that the explanatory variables tend not to have high correlations with each other. The explanatory variable pairs that exhibit correlations is bigger than 0.70 is GrLivArea-Bathrooms. Bathrooms also has elevated correlation with Age. Given this, I decided to drop Bathrooms from the regression analysis.

```
df = df.drop(['Bathrooms'], axis=1)
```

Dividing the data into pre-2010 and a 2010 holdout

Will split the data into a pre-2010 group, for training and testing on three linear models: OLS regression, Ridge regression and Lasso regression.

The 2020 data will be retained as the holdout test data to evaluate the final model against unseen data.

```
# Assigning the 2006-2009 data to another dataset
df_0609 = df.loc[df['YrSold_2010'] != 1]
df_0609.shape
```

```
(2369, 29)
```

```
# Assigning 2010 data as the holdout test set
df_2010 = df.loc[df['YrSold_2010'] == 1]
```

Functions

```
# Function for scoring training set
def train_scores(model, X, y):
    """
    model: fitted model
    X: Matrix of explanatory variables (train set)
    y: Dependent variable (train set)
    """
    cv_scores = cross_val_score(model, X, y, cv=5) # 5-fold cross-validation
    print('Training Score:', np.round(model.score(X, y), 4))
    print('Cross-validation scores:', np.round(cv_scores, 4))
    print('Mean cross-validation score:', np.round(cv_scores.mean(), 4))
```

```
# Function for scoring test set
def test_scores(model, X, y):
    """
    model: fitted model
    X: Matrix of explanatory variables (test set)
    y: Dependent variable (test set)
    """
    yhat = model.predict(X)
    print('Mean Squared Error:', np.round(metrics.mean_squared_error(y, yhat), 4))
    print('Root Mean Squared Error:', np.round(metrics.mean_squared_error(y, yhat)**0.5, 4))
```

```
# Function for plotting histogram of residuals
def resid_histogram(model, X, y, period=''):
    """
    model: fitted model
    X: Matrix of explanatory variables
    y: Dependent variable
    period: String describing data coverage period
    """
    yhat = model.predict(X)
    residuals = y - yhat
    fig, ax = plt.subplots(figsize=(10, 6))
    sns.distplot(residuals, bins=50, kde=True, ax=ax)
    plt.title('OLS Residuals, (period)', fontsize=18)
```

Regression using Ln SalePrice target on pre-2010

```
y_SP = df_0609['SalePrice']
y_LnSP = df_0609['LnSalePrice']
```

```
X = df_0609.drop(['SalePrice', 'LnSalePrice'], axis=1)
```

```
X.shape
```

```
(2369, 27)
```

```
# Train-test split of the 2006-2009 data
X_train, X_test, y_train, y_test = train_test_split(X, y_LnSP, test_size=0.3, random_state=0)
```

```
scaler = StandardScaler()
X_train = pd.DataFrame(scaler.fit_transform(X_train), columns=X_train.columns)
X_test = pd.DataFrame(scaler.transform(X_test), columns=X_test.columns)
```

```
# Fitting ordinary linear regression & getting parameter estimates
ols = LinearRegression()
ols.fit(X_train, y_train)
```

```
print('Intercept:', ols.intercept_)
print('Coefficients:', ols.coef_)
```

```
Intercept: 12.034995130834467
Coefficients: [-8.31451502e-02  1.59833637e-01  4.16824559e-02  3.57682462e-02
 0.51432178e-04 -0.08753798e-03 -1.08074458e-02  4.72361881e-02
 1.04875638e-01 -1.39488289e-02  2.24271387e-02 -2.49826156e-02
 8.89554266e-02 -2.01677541e-01 -2.30295166e-02  8.95397428e-03
 1.49051226e-02 -5.40926107e-03  1.49262932e-02  2.71861865e-02
 2.149784104e-02 0.61718106e-03 5.63237319e-03
 1.21379704e-03 -7.29921336e-04 0.00908090e+00]
```

```
# OLS training set scores, including CV scores
train_scores(ols, X_train, y_train)
```

```
Training Score: 0.936
Cross-validation scores: [0.8789 0.9275 0.9272 0.8881 0.9126]
Mean cross-validation score: 0.9069
```

```
# Shuffled 5-fold cross validation scores are rather similar
kf = KFold(n_splits=5, shuffle=True, random_state=1)
cv_scores_shuffled = cross_val_score(ols, X_train, y_train, cv=kf)
```

```
print('Shuffled cross validation score:', np.round(cv_scores_shuffled, 4))
print('Mean shuffled cross validation score:', np.round(cv_scores_shuffled.mean(), 4))
```

```
Shuffled cross validation score: [0.92  0.8973 0.8921 0.929  0.8926]
Mean shuffled cross validation score: 0.9082
```

```
# OLS test set score
test_scores(ols, X_test, y_test)
```

```
Test Score: 0.9142
```

```
# OLS MSE & RMSE scores
accuracy_scores(ols, X_test, y_test)
```

```
Mean Squared Error: 0.0112
Root Mean Squared Error: 0.1057
```

```
# Collect the coefficients
df_ols_coef = pd.DataFrame(ols.coef_, index=X_train.columns, columns=['Coefficients'])
df_ols_coef['Coef_abs'] = df_ols_coef['Coefficients'].abs()
```

Analysis of the OLS residuals

```
predictions_train = ols.predict(X_train)
predictions_test = ols.predict(X_test)
```

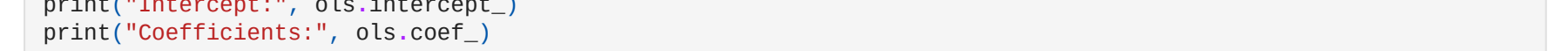
```
# Descriptive statistics of training set residuals
ols_residuals_0609 = (y_train - predictions_train)
ols_residuals_2010 = (y_test - predictions_test)
```

```
Out[28]: count      1.614000e+03
mean      -4.944025e-16
std       0.105772
min       -0.639473
25%       -0.055359
50%       0.005378
75%       0.058850
max       0.519638
Name: LnSalePrice, dtype: float64
```

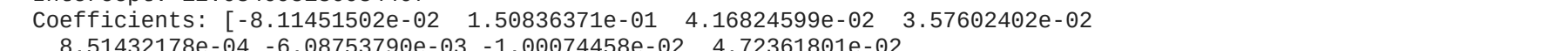
```
# Histogram of training set residuals show that they are approximately normally distributed with mean 0
# There is no indication of a left-tail, indicating that the model overpredicts the target variable
# at the very low end of 'LnSalePrice'
resid_histogram(ols, X_train, y_train, period='2006-2009 train data')
```

```
C:\Users\cam\Anaconda3\Lib\site-packages\seaborn\distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
```

OLS Residuals, 2006-2009 train data



```
from scipy import stats
stats.probplot(ols_residuals_0609, dist='norm', plot=plt)
plt.title('Quantile-Quantile Plot, 2006-2009 training set');
```



```
# Descriptive statistics of test set residuals
ols_residuals_test = (y_test - predictions_test)
ols_residuals_test.describe()
```

```
Out[32]: count      693.000000
mean      0.093343
std       0.105772
min       -0.639473
25%       -0.055359
50%       0.005378
75%       0.058850
max       0.519638
Name: LnSalePrice, dtype: float64
```

```
# Histogram of test set residuals is symmetric and approximately normal in shape
resid_histogram(ols, X_test, y_test, period='2006-2009 test data')
```

```
C:\Users\cam\Anaconda3\Lib\site-packages\seaborn\distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
```

OLS Residuals, 2006-2009 test data



```
from scipy import stats
stats.probplot(ols_residuals_test, dist='norm', plot=plt)
plt.title('Quantile-Quantile Plot, 2006-2009 test set');
```



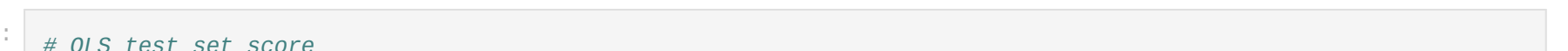
```
# Descriptive statistics of test set residuals
ols_residuals_test = (y_test - predictions_test)
ols_residuals_test.describe()
```

```
Out[32]: count      693.000000
mean      0.093343
std       0.105772
min       -0.639473
25%       -0.055359
50%       0.005378
75%       0.058850
max       0.519638
Name: LnSalePrice, dtype: float64
```

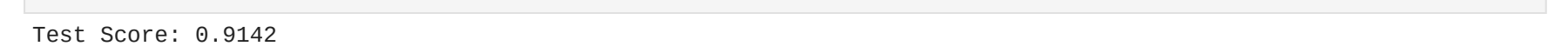
```
# Histogram of test set residuals is symmetric and approximately normal in shape
resid_histogram(ols, X_test, y_test, period='2006-2009 test data')
```

```
C:\Users\cam\Anaconda3\Lib\site-packages\seaborn\distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
```

OLS Residuals, 2006-2009 test data



```
stats.probplot(ols_residuals_test, dist='norm', plot=plt)
plt.title('Quantile-Quantile Plot, 2006-2009 test set');
```



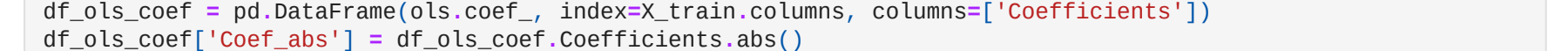
```
# Descriptive statistics of test set residuals
ols_residuals_test = (y_test - predictions_test)
ols_residuals_test.describe()
```

```
Out[32]: count      693.000000
mean      0.093343
std       0.105772
min       -0.639473
25%       -0.055359
50%       0.005378
75%       0.058850
max       0.519638
Name: LnSalePrice, dtype: float64
```

```
# Histogram of test set residuals is symmetric and approximately normal in shape
resid_histogram(ols, X_test, y_test, period='2006-2009 test data')
```

```
C:\Users\cam\Anaconda3\Lib\site-packages\seaborn\distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
```

OLS Residuals, 2006-2009 test data



```
stats.probplot(ols_residuals_test, dist='norm', plot=plt)
plt.title('Quantile-Quantile Plot, 2006-2009 test set');
```



```
# Descriptive statistics of test set residuals
ols_residuals_test = (y_test - predictions_test)
ols_residuals_test.describe()
```

```
Out[32]: count      693.000000
mean      0.093343
std       0.105772
min       -0.639473
25%       -0.055359
50%       0.005378
75%       0.058850
max       0.519638
Name: LnSalePrice, dtype: float64
```

```
# Histogram of test set residuals is symmetric and approximately normal in shape
resid_histogram(ols, X_test, y_test, period='2006-2009 test data')
```

```
C:\Users\cam\Anaconda3\Lib\site-packages\seaborn\distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
```

OLS Residuals, 2006-2009 test data



```
stats.probplot(ols_residuals_test, dist='norm', plot=plt)
plt.title('Quantile-Quantile Plot, 2006-2009 test set');
```



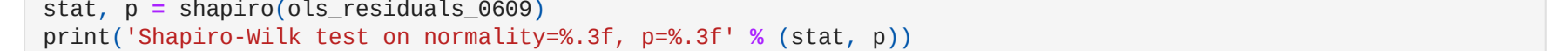
```
# Descriptive statistics of test set residuals
ols_residuals_test = (y_test - predictions_test)
ols_residuals_test.describe()
```

```
Out[32]: count      693.000000
mean      0.093343
std       0.105772
min       -0.639473
25%       -0.055359
50%       0.005378
75%       0.058850
max       0.519638
Name: LnSalePrice, dtype: float64
```

```
# Histogram of test set residuals is symmetric and approximately normal in shape
resid_histogram(ols, X_test, y_test, period='2006-2009 test data')
```

```
C:\Users\cam\Anaconda3\Lib\site-packages\seaborn\distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
```

OLS Residuals, 2006-2009 test data



```
stats.probplot(ols_residuals_test, dist='norm', plot=plt)
plt.title('Quantile-Quantile Plot, 2006-2009 test set');
```



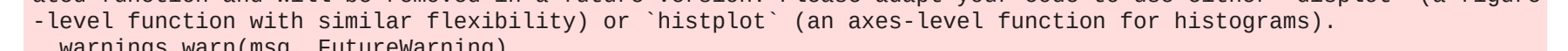
```
# Descriptive statistics of test set residuals
ols_residuals_test = (y_test - predictions_test)
ols_residuals_test.describe()
```

```
Out[32]: count      693.000000
mean      0.093343
std       0.105772
min       -0.639473
25%       -0.055359
50%       0.005378
75%       0.058850
max       0.519638
Name: LnSalePrice, dtype: float64
```

```
# Histogram of test set residuals is symmetric and approximately normal in shape
resid_histogram(ols, X_test, y_test, period='2006-2009 test data')
```

```
C:\Users\cam\Anaconda3\Lib\site-packages\seaborn\distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
```

OLS Residuals, 2006-2009 test data



```
stats.probplot(ols_residuals_test, dist='norm', plot=plt)
plt.title('Quantile-Quantile Plot, 2006-2009 test set');
```



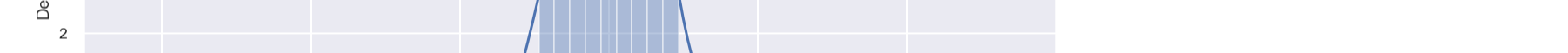
```
# Descriptive statistics of test set residuals
ols_residuals_test = (y_test - predictions_test)
ols_residuals_test.describe()
```

```
Out[32]: count      693.000000
mean      0.093343
std       0.105772
min       -0.639473
25%       -0.055359
50%       0.005378
75%       0.058850
max       0.519638
Name: LnSalePrice, dtype: float64
```

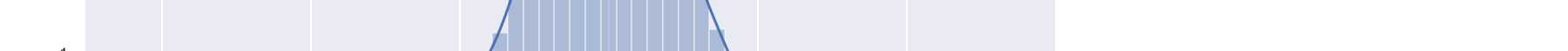
```
# Histogram of test set residuals is symmetric and approximately normal in shape
resid_histogram(ols, X_test, y_test, period='2006-2009 test data')
```

```
C:\Users\cam\Anaconda3\Lib\site-packages\seaborn\distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
```

OLS Residuals, 2006-2009 test data



```
stats.probplot(ols_residuals_test, dist='norm', plot=plt)
plt.title('Quantile-Quantile Plot, 2006-2009 test set');
```



```
# Descriptive statistics of test set residuals
ols_residuals_test = (y_test - predictions_test)
ols_residuals_test.describe()
```

```
Out[32]: count      693.000000
mean      0.093343
std       0.105772
min       -0.639473
25%       -0.055359
50%       0.005378
75%       0.058850
max       0.519638
Name: LnSalePrice, dtype: float64
```

```
# Histogram of test set residuals is symmetric and approximately normal in shape
resid_histogram(ols, X_test, y_test, period='2006-2009 test data')
```

```
C:\Users\cam\Anaconda3\Lib\site-packages\seaborn\distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
```

OLS Residuals, 2006-2009 test data



```
stats.probplot(ols_residuals_test, dist='norm', plot=plt)
plt.title('Quantile-Quantile Plot, 2006-2009 test set');
```



```
# Descriptive statistics of test set residuals
ols_residuals_test = (y_test - predictions_test)
ols_residuals_test.describe()
```

```
Out[32]: count      693.000000
mean      0.093343
std       0.105772
min       -0.639473
25%       -0.055359
50%       0.005378
75%       0.058850
max       0.519638
Name: LnSalePrice, dtype: float64
```

```
# Histogram of test set residuals is symmetric and approximately normal in shape
resid_histogram(ols, X_test, y_test, period='2006-2009 test data')
```

```
C:\Users\cam\Anaconda3\Lib\site-packages\seaborn\distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
```

OLS Residuals, 2006-2009 test data



```
stats.probplot(ols_residuals_test, dist='norm', plot=plt)
plt.title('Quantile-Quantile Plot, 2006-2009 test set');
```



```
# Descriptive statistics of test set residuals
ols_residuals_test = (y_test - predictions_test)
ols_residuals_test.describe()
```

```
Out[32]: count      693.000000
mean      0.093343
std       0.105772
min       -0.639473
25%       -0.055359
50%       0.005378
75%       0.058850
max       0.519638
Name: LnSalePrice, dtype: float64
```

```
# Histogram of test set residuals is symmetric and approximately normal in shape
resid_histogram(ols, X_test, y_test, period='2006-2009 test data')
```

```
C:\Users\cam\Anaconda3\Lib\site-packages\seaborn\distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
```

OLS Residuals, 2006-2009 test data



```
stats.probplot(ols_residuals_test, dist='norm', plot=plt)
plt.title('Quantile-Quantile Plot, 2006-2009 test set');
```



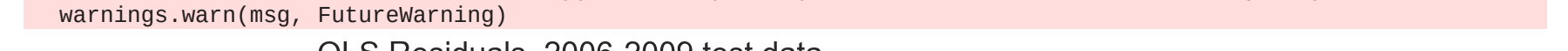
```
# Descriptive statistics of test set residuals
ols_residuals_test = (y_test - predictions_test)
ols_residuals_test.describe()
```

```
Out[32]: count      693.000000
mean      0.093343
std       0.105772
min       -0.639473
25%       -0.055359
50%       0.005378
75%       0.058850
max       0.519638
Name: LnSalePrice, dtype: float64
```

```
# Histogram of test set residuals is symmetric and approximately normal in shape
resid_histogram(ols, X_test, y_test, period='2006-2009 test data')
```

```
C:\Users\cam\Anaconda3\Lib\site-packages\seaborn\distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
```

OLS Residuals, 2006-2009 test data



```
stats.probplot(ols_residuals_test, dist='norm', plot=plt)
plt.title('Quantile-Quantile Plot, 2006-2009 test set');
```



```
In [69]: lasso_mod = LassoCV(lphs=np.logspace(-4, 4, 10), cv=5)
lasso_mod.fit(X_train, y_train)

print('Best Lasso alpha:', lasso_mod.alpha_)
```

Best Lasso alpha: 0.000774263682681127

```
In [70]: train_scores = lasso_mod.X_train, y_train)
train_scores(lasso_mod.X_train, y_train)
```

Training Score: 0.9159  
Cross-validation scores: [0.9182 0.9249 0.9074 0.8966 0.9098]  
Mean cross-validation score: 0.9098  
Test Score: 0.9013

Using the 2010 holdout data as the test set, we see that the **R-squared drops marginally to near-0.90 in the test set** across the three linear models from the 0.91+ achieved in the 2006-2009 training set. The scores again appear **consistent across the OLS, Ridge and Lasso regressions**.

The **OLS residuals** of the 2010 data has a distribution with a very slight positive mean value and a longer left-tail, but it looks approximately normal.

Given this, I will proceed with just the OLS regression model for **statistical inference**.

# Regression model using full 2006-2010 data

```
In [71]: y_lnSP = df['LnSalePrice']
```

```
In [72]: X_fin = df.drop(['SalePrice', 'LnSalePrice'], axis=1)
```

```
In [73]: X_fin.shape
```

Out[73]: (2616, 27)

```
In [74]: X_fin = pd.DataFrame(Scaler().fit_transform(X_fin), columns=X_fin.columns)
```

```
In [75]: ols_fit(X_fin, y_lnSP)
```

Out[75]: LinearRegression()

```
In [76]: # Model scores on the full 2006-2010 data
test_scores(ols_fit(X_fin, y_lnSP))
accuracy_scores(ols_fit(X_fin, y_lnSP))
```

Test Score: 0.9148  
Mean Squared Error: 0.0113  
Root Mean Squared Error: 0.1062

```
In [77]: df_lnSP_coef = pd.DataFrame(ols_fit.coef_, index=X_fin.columns,
                                columns=['Coefficients'])
df_lnSP_coef['Coeff_abs'] = df_lnSP_coef.Coefficients.abs()
```

```
In [78]: # Descriptive statistics of the residuals
predictions = ols_fit.predict(X_fin)
error_term = (y_lnSP - predictions)
error_term.describe()
```

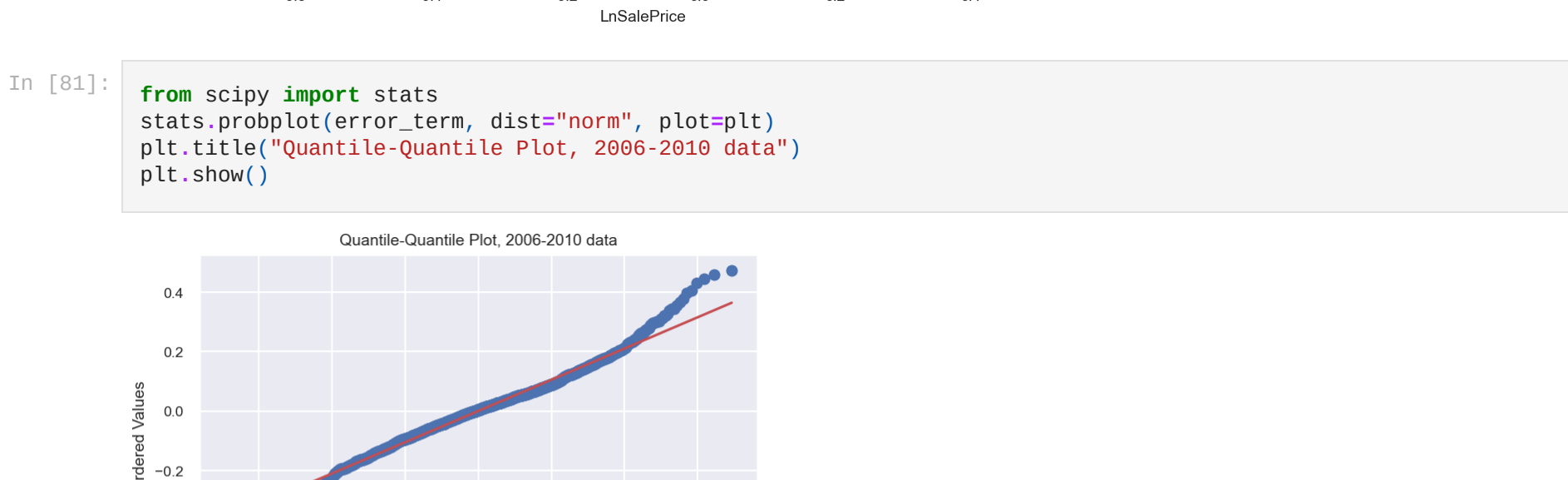
Out[78]: count 2.616000e+03  
mean -0.85955e-16  
std 1.062177e-01  
min -6.37565e-01  
25% -6.01957e-02  
50% 4.038420e-03  
75% 5.920534e-02  
max 4.715965e-01  
Name: LnSalePrice, dtype: float64

```
In [79]: # Residuals are approximately normal, given the skew and kurtosis
print("Skew:", error_term.skew())
print("Kurtosis:", error_term.kurtosis())
stat, p = shapiro(error_term)
print('Shapiro-Wilk test on normality=%3f, p=%3f' % (stat, p))
```

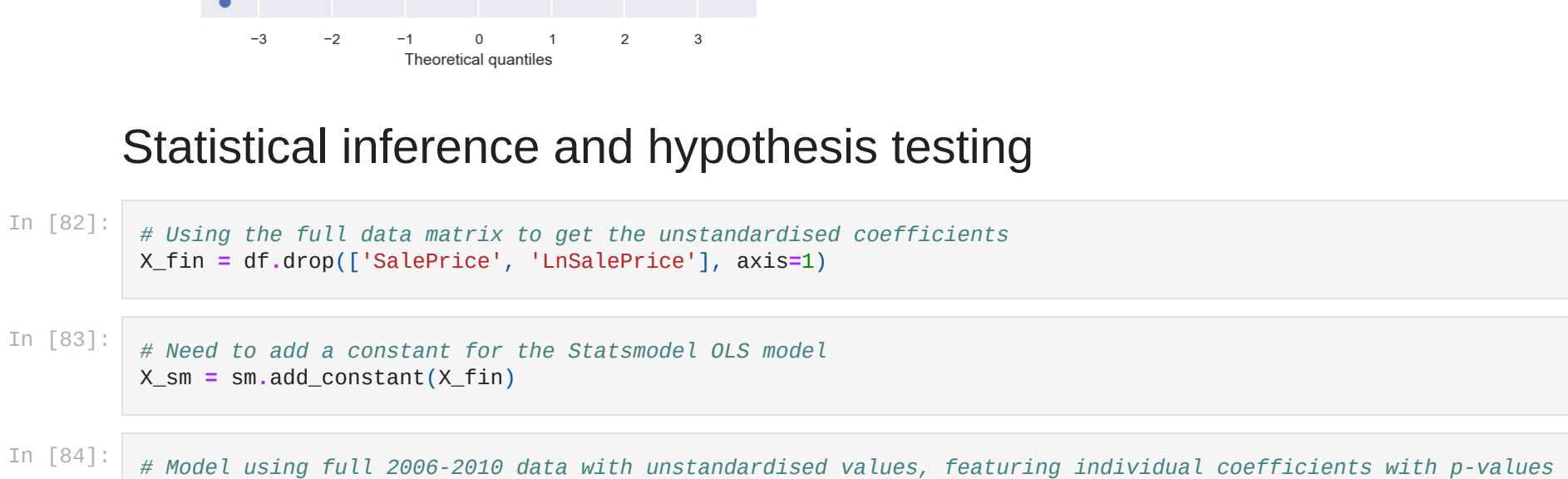
Skew: -0.17607125393577763  
Kurtosis: 2.310713114664527  
Shapiro-Wilk test on normality=0.978, p=0.000

```
In [80]: resid_histgram(ols_fit, X_fin, y_lnSP, period='full 2006-2010 data')
```

C:\Users\camb7\anaconda3\lib\site-packages\seaborn\distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).



```
In [81]: from scipy import stats
stats.probplot(error_term, dist='norm', plot=plt)
plt.title("Quantile-Quantile Plot, 2006-2010 data")
plt.show()
```



# Statistical inference and hypothesis testing

```
In [82]: # Using the full data matrix to get the unstandardised coefficients
X_fin = df.drop(['SalePrice', 'LnSalePrice'], axis=1)
```

```
In [83]: # Need to add a constant for the Statsmodel OLS model
X_fin = sm.add_constant(X_fin)
```

```
In [84]: # Model using full 2006-2010 data with unstandardised values, featuring individual coefficients with p-values
model = sm.OLS(y_lnSP, X_fin)
results = model.fit()
print(results.summary())
```

```
OLS Regression Results
=====
Dep. Variable:  LnSalePrice  R-squared:  0.915
Model:  OLS  Adj. R-squared:  0.914
Method:  Least Squares  F-statistic:  1029.
Date:  Sun, 23 May 2021  Prob (F-statistic):  0.00
Time:  18:14:04  Log-Likelihood:  2154.3
No. Observations:  2616  AIC:  -4253.
Df Residuals:  2588  BIC:  -4089.
Df Model:  27
Covariance Type:  nonrobust

=====
coef    std err          t    P>|t|    [0.025    0.975]
-----
const    10.5360      0.024    445.575      0.000    10.492    10.585
Age      -0.0027      0.000    -23.277      0.000    -0.003    -0.002
GrLivArea    0.0003    0.35e-06    35.585      0.000    0.000    0.000
BaselLivArea    9.71e-05    5.44e-06    17.875      0.000    8.65e-05    0.000
Location    0.0333      0.003    11.353      0.000    0.028    0.039
Amenities    0.0146      0.015    -0.961    0.337    -0.044    0.015
RoadRail    -0.0410      0.009    -4.616      0.000    -0.058    -0.024
BedroomAbvGr    -0.0117      0.004    -3.267      0.001    -0.019    -0.005
OverallCond    0.0457      0.002    19.081      0.000    0.041    0.050
OverallQual    0.0812      0.003    30.263      0.000    0.076    0.086
LotFrontage    0.0008      0.000    0.112      0.900    -0.001    0.001
LotArea    3.189e-05    3.57e-07    8.915      0.000    2.48e-05    3.89e-05
TwoStory_dum    -0.0458      0.006    -7.563      0.000    -0.058    -0.034
FlatContour_dum    -0.0204      0.007    -2.737      0.006    -0.035    -0.006
FlatRoof_dum    0.0336      0.027    1.263      0.207    -0.019    0.086
GarageArea    0.0001    1.59e-05    0.228      0.880    9.94e-05    0.000
Garage_dum    0.0205      0.012    1.713      0.087    -0.003    0.044
CentralAirNum    0.0680      0.010    6.560      0.000    0.043    0.088
LowQualFinSF    -0.0001    5.37e-05    -2.082      0.037    -0.000    -3.87e-05
Fireplaces    0.0261      0.004    6.485      0.000    0.018    0.034
KitchenQual_Ex    0.0715      0.010      7.123      0.000    0.052    0.091
Zoning_2    -0.0370      0.025    -1.521    0.128    -0.086    0.011
Zoning_3    0.0497      0.007      7.037      0.000    0.036    0.063
Zoning_4    0.0565      0.012      4.586      0.000    0.032    0.081
YrSold_2007    0.0035      0.006      0.563      0.573    -0.009    0.016
YrSold_2008    0.0065      0.007      1.006      0.314    -0.006    0.019
YrSold_2009    -0.0064      0.006    -0.996      0.322    -0.019    0.006
YrSold_2010    0.0156      0.008      2.035      0.042    0.001    0.031

=====
Omnibus:  158.447  Durbin-Watson:  1.002
Prob(Omnibus):  0.000  Jarque-Bera (JB):  592.121
Skew:  -0.176  Prob(JB):  2.65e-129
Kurtosis:  5.394  Cond. No.  1.55e+05
=====

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 1.55e+05. This might indicate that there are strong multicollinearity or other numerical problems.
```

The majority of the variables in the OLS model are statistically significant at the 5% level, as indicated by the individual t-statistics and p-values. The exceptions are 'RoadRail' and 'LowQualFinSF'.

```
In [85]: ols_fit(X_fin, y_lnSP)
unscaled_coef = ols_fit.coef_
```

```
In [86]: import math
transformed_coef = []
for i in unscaled_coef:
    j = math.exp(i)
    transformed_coef.append(j)
print(transformed_coef)
```

[0.997252408705468, 1.0093329291585632, 1.00009718923651, 1.0338247623381687, 0.9855426176035118, 0.9598214705320343, 0.98836928093296, 1.0467486352268127, 1.0845439294576424, 1.000751771469628, 1.0000831846894666, 0.9952580332454207, 0.9797936893162256, 1.0341265404179930, 1.000130531339037, 1.0206625199609974, 1.070345859587567, 0.9989559869822011, 1.0264156398065571, 1.074077907916657, 0.963052462805285, 1.050908491836821, 1.0581614678622424, 1.0035490335757256, 1.0665624160230066, 0.9936692254476065, 1.015710639164118]

```
In [87]: coef_effect = [(i - 1)*df.SalePrice.mean() for i in transformed_coef]
```

```
In [88]: var_impact = pd.DataFrame(data=[X_fin.columns, coef_effect]).T
```

```
In [89]: var_impact.columns = ["variable", "1-unit change"]
```

```
In [90]: var_impact
```

Out[90]:

	variable	1-unit change
0	Age	-487.364
1	GrLivArea	50.5493
2	BaselLivArea	17.3802
3	Location	6050.06
4	Amenities	-2585.92
5	RoadRail	-2186.52
6	BedroomAbvGr	-2080.32
7	OverallCond	8361.69
8	OverallQual	15121.9
9	LotFrontage	134.465
10	LotArea	0.566629
11	TwoStory_dum	-9002.65
12	FlatContour_dum	-6113.73
13	FlatRoof_dum	6104.04
14	GarageArea	23.3475
15	Garage_dum	3695.8
16	CentralAirNum	12502.4
17	LowQualFinSF	-25.7589
18	Fireplaces	4724.83
19	KitchenQual_Ex	13249.9
20	Zoning_2	-6908.61
21	Zoning_3	9105.74
22	Zoning_4	10403
23	YrSold_2007	634.797
24	YrSold_2008	1173.79
25	YrSold_2009	-1132.35
26	YrSold_2010	2810.08

The DataFrame above shows the impact of a **1-unit change** in the value the respective variables on the **average sale price**. For example, a house that is **one unit** (1 year in this case) more in "Age" will cause the average sale price to drop by **\$487** all else being equal. A house that has a location score that is **one rank** higher will cause the average sale price to rise **\$6,050** all else being equal.