# Data Cleaning and EDA

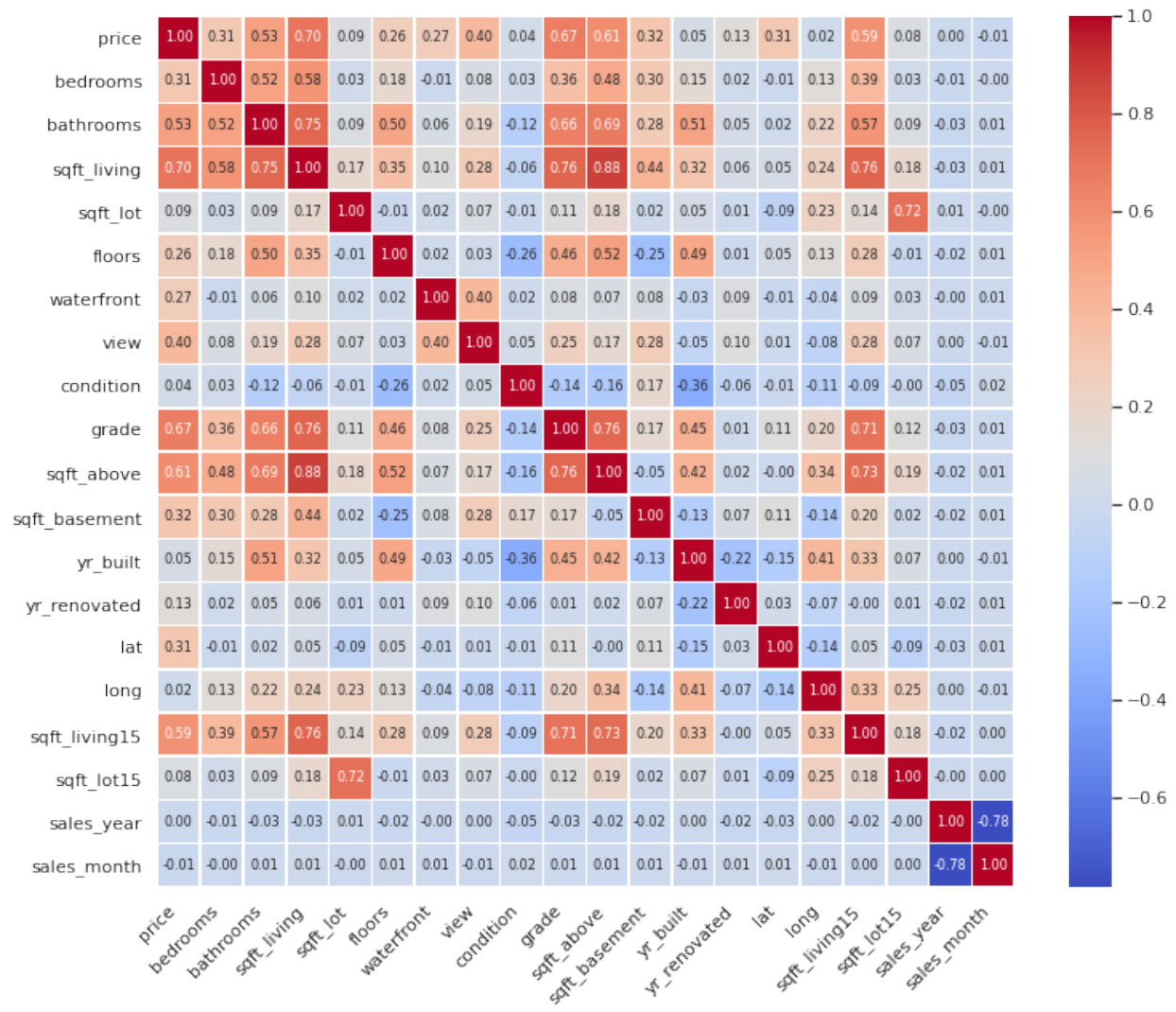| | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors | waterfront | view | condition | grade | sqft_above | sqft_basement | yr_built | yr_renovated | lat | long | sqft_living15 | sqft_lot15 | sales_year | sales_month |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| price | 1.00 | 0.31 | 0.53 | 0.70 | 0.09 | 0.26 | 0.27 | 0.40 | 0.04 | 0.67 | 0.61 | 0.32 | 0.05 | 0.13 | 0.31 | 0.02 | 0.59 | 0.08 | 0.00 | -0.01 |
| bedrooms | 0.31 | 1.00 | 0.52 | 0.58 | 0.03 | 0.18 | -0.01 | 0.08 | 0.03 | 0.36 | 0.48 | 0.30 | 0.15 | 0.02 | -0.01 | 0.13 | 0.39 | 0.03 | -0.01 | -0.00 |
| bathrooms | 0.53 | 0.52 | 1.00 | 0.75 | 0.09 | 0.50 | 0.06 | 0.19 | -0.12 | 0.66 | 0.69 | 0.28 | 0.51 | 0.05 | 0.02 | 0.22 | 0.57 | 0.09 | -0.03 | 0.01 |
| sqft_living | 0.70 | 0.58 | 0.75 | 1.00 | 0.17 | 0.35 | 0.10 | 0.28 | -0.06 | 0.76 | 0.88 | 0.44 | 0.32 | 0.06 | 0.05 | 0.24 | 0.76 | 0.18 | -0.03 | 0.01 |
| sqft_lot | 0.09 | 0.03 | 0.09 | 0.17 | 1.00 | -0.01 | 0.02 | 0.07 | -0.01 | 0.11 | 0.18 | 0.02 | 0.05 | 0.01 | -0.09 | 0.23 | 0.14 | 0.72 | 0.01 | -0.00 |
| floors | 0.26 | 0.18 | 0.50 | 0.35 | -0.01 | 1.00 | 0.02 | 0.03 | -0.26 | 0.46 | 0.52 | -0.25 | 0.49 | 0.01 | 0.05 | 0.13 | 0.28 | -0.01 | -0.02 | 0.01 |
| waterfront | 0.27 | -0.01 | 0.06 | 0.10 | 0.02 | 0.02 | 1.00 | 0.40 | 0.02 | 0.08 | 0.07 | 0.08 | -0.03 | 0.09 | -0.01 | -0.04 | 0.09 | 0.03 | -0.00 | 0.01 |
| view | 0.40 | 0.08 | 0.19 | 0.28 | 0.07 | 0.03 | 0.40 | 1.00 | 0.05 | 0.25 | 0.17 | 0.28 | -0.05 | 0.10 | 0.01 | -0.08 | 0.28 | 0.07 | 0.00 | -0.01 |
| condition | 0.04 | 0.03 | -0.12 | -0.06 | -0.01 | -0.26 | 0.02 | 0.05 | 1.00 | -0.14 | -0.16 | 0.17 | -0.36 | -0.06 | -0.01 | -0.11 | -0.09 | -0.00 | -0.05 | 0.02 |
| grade | 0.67 | 0.36 | 0.66 | 0.76 | 0.11 | 0.46 | 0.08 | 0.25 | -0.14 | 1.00 | 0.76 | 0.17 | 0.45 | 0.01 | 0.11 | 0.20 | 0.71 | 0.12 | -0.03 | 0.01 |
| sqft_above | 0.61 | 0.48 | 0.69 | 0.88 | 0.18 | 0.52 | 0.07 | 0.17 | -0.16 | 0.76 | 1.00 | -0.05 | 0.42 | 0.02 | -0.00 | 0.34 | 0.73 | 0.19 | -0.02 | 0.01 |
| sqft_basement | 0.32 | 0.30 | 0.28 | 0.44 | 0.02 | -0.25 | 0.08 | 0.28 | 0.17 | 0.17 | -0.05 | 1.00 | -0.13 | 0.07 | 0.11 | -0.14 | 0.20 | 0.02 | -0.02 | 0.01 |
| yr_built | 0.05 | 0.15 | 0.51 | 0.32 | 0.05 | 0.49 | -0.03 | -0.05 | -0.36 | 0.45 | 0.42 | -0.13 | 1.00 | -0.22 | -0.15 | 0.41 | 0.33 | 0.07 | 0.00 | -0.01 |
| yr_renovated | 0.13 | 0.02 | 0.05 | 0.06 | 0.01 | 0.01 | 0.09 | 0.10 | -0.06 | 0.01 | 0.02 | 0.07 | -0.22 | 1.00 | 0.03 | -0.07 | -0.00 | 0.01 | -0.02 | 0.01 |
| lat | 0.31 | -0.01 | 0.02 | 0.05 | -0.09 | 0.05 | -0.01 | 0.01 | -0.01 | 0.11 | -0.00 | 0.11 | -0.15 | 0.03 | 1.00 | -0.14 | 0.05 | -0.09 | -0.03 | 0.01 |
| long | 0.02 | 0.13 | 0.22 | 0.24 | 0.23 | 0.13 | -0.04 | -0.08 | -0.11 | 0.20 | 0.34 | -0.14 | 0.41 | -0.07 | -0.14 | 1.00 | 0.33 | 0.25 | 0.00 | -0.01 |
| sqft_living15 | 0.59 | 0.39 | 0.57 | 0.76 | 0.14 | 0.28 | 0.09 | 0.28 | -0.09 | 0.71 | 0.73 | 0.20 | 0.33 | -0.00 | 0.05 | 0.33 | 1.00 | 0.18 | -0.02 | 0.00 |
| sqft_lot15 | 0.08 | 0.03 | 0.09 | 0.18 | 0.72 | -0.01 | 0.03 | 0.07 | -0.00 | 0.12 | 0.19 | 0.02 | 0.07 | 0.01 | -0.09 | 0.25 | 0.18 | 1.00 | -0.00 | 0.00 |
| sales_year | 0.00 | -0.01 | -0.03 | -0.03 | 0.01 | -0.02 | -0.00 | 0.00 | -0.05 | -0.03 | -0.02 | -0.02 | 0.00 | -0.02 | -0.03 | 0.00 | -0.02 | -0.00 | 1.00 | -0.78 |
| sales_month | -0.01 | -0.00 | 0.01 | 0.01 | -0.00 | 0.01 | 0.01 | -0.01 | 0.02 | 0.01 | 0.01 | 0.01 | -0.01 | 0.01 | 0.01 | -0.01 | 0.00 | 0.00 | -0.78 | 1.00 |

```
from sklearn.model_selection import train_test_split

# Split the data into training (80%) and testing (20%) while keeping 'price' in both
X_train, X_test = train_test_split(df, test_size=0.2, random_state=42)  # Set random_state for reproducibility

# Print the number of rows in each dataset
print("Training set size:", len(X_train))
print("Test set size:", len(X_test))
```

Training set size: 17290
Test set size: 4323

```
                        OLS Regression Results
================================================================================
Dep. Variable:                  price   R-squared:                       0.492
Model:                            OLS   Adj. R-squared:                  0.492
Method:                 Least Squares   F-statistic:                 1.677e+04
Date:                Sat, 08 Feb 2025   Prob (F-statistic):               0.00
Time:                        21:36:33   Log-Likelihood:            -2.3995e+05
No. Observations:               17290   AIC:                         4.799e+05
Df Residuals:                   17288   BIC:                         4.799e+05
Df Model:                           1
Covariance Type:            nonrobust
================================================================================
                 coef    std err          t      P>|t|      [0.025      0.975]
--------------------------------------------------------------------------------
Intercept     -4.2e+04   4886.778     -8.594      0.000   -5.16e+04   -3.24e+04
sqft_living   279.5548      2.159    129.496      0.000     275.323     283.786
================================================================================
Omnibus:                    11990.495   Durbin-Watson:                   2.030
Prob(Omnibus):                  0.000   Jarque-Bera (JB):           483410.340
Skew:                           2.835   Prob(JB):                         0.00
Kurtosis:                      28.276   Cond. No.                     5.65e+03
================================================================================

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 5.65e+03. This might indicate that there are
strong multicollinearity or other numerical problems.
Adjusted R-squared: 0.4923544744403926
```

List:

sqft_living: 0.4923
grade: 0.4451
sqft_above: 0.3856
bathrooms: 0.3338
sqft_living15: 0.3249
view: 0.1877
sqft_basement: 0.1721
lat: 0.1124
waterfront: 0.1056
floors: 0.0895
yr_built: 0.0521
sqft_lot: 0.0087
sqft_lot15: 0.0079
yr_renovated: 0.0032
long: 0.0028
condition: 0.0019

The top 3 predictors based on the adj R squared:  sqft_living , grade, and  sqft_above.  I conducted simple linear regression for each predictor using statsmodels.ols(), with price as the dependent variable. After fitting the models, I extracted the Adjusted R-squared values for each predictor to evaluate their performance. Finally, I ranked the predictors in descending order based on their Adjusted R-squared

values to determine which variables had the strongest relationship with price.   Yes, the correlation matrix analysis identified sqft_living as the best guess predictor because it had the highest correlation with price. After validating this through Adjusted R-squared values from linear regression, sqft_living remains the strongest predictor, confirming its significant relationship with price.