



Case Study: Regression

Doojin Kim, Guilherme, Sabina

TABLE OF CONTENTS

01

INTRODUCTION

02

OBJECTIVE OF
BUSINESS CASE

03

DATASET
OVERVIEW &
METHODOLOGY

04

KEY INSIGHTS

05

ML MODEL

06

CONCLUSION

INTRODUCTION



22000

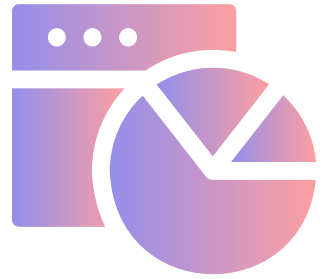
Dataset containing information on
22,000 properties



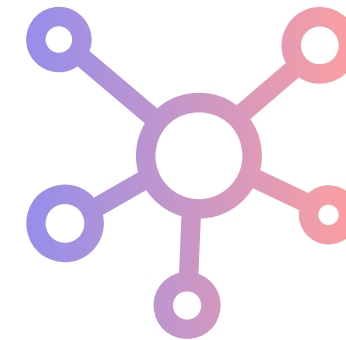
**SOLD 2014-
2015**

sold between May 2014 and May 2015

2.OBJECTIVES OF BUSINESS CASE



**Understand and
perform the
necessary EDA
steps**



**build a ML
Model that can
accurately
predict the
selling prices**



**To identify the
factors that
influence the
selling price**

3. DATASET OVERVIEW AND METHODOLOGY



ID



DATE



WATERFRONT



CONDITION



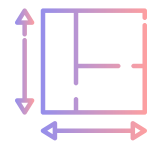
GRADE



'SQFT_ABOVE'



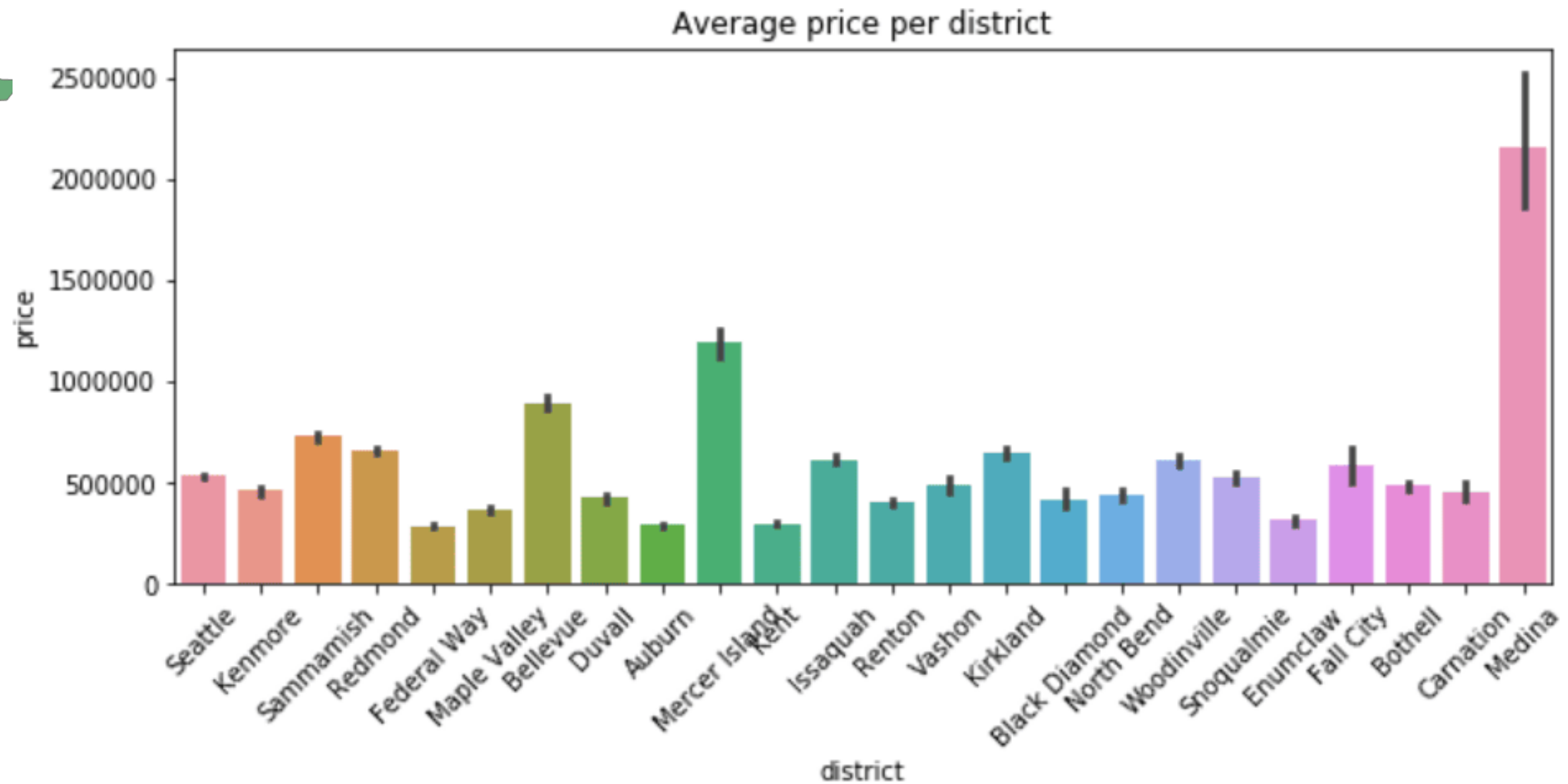
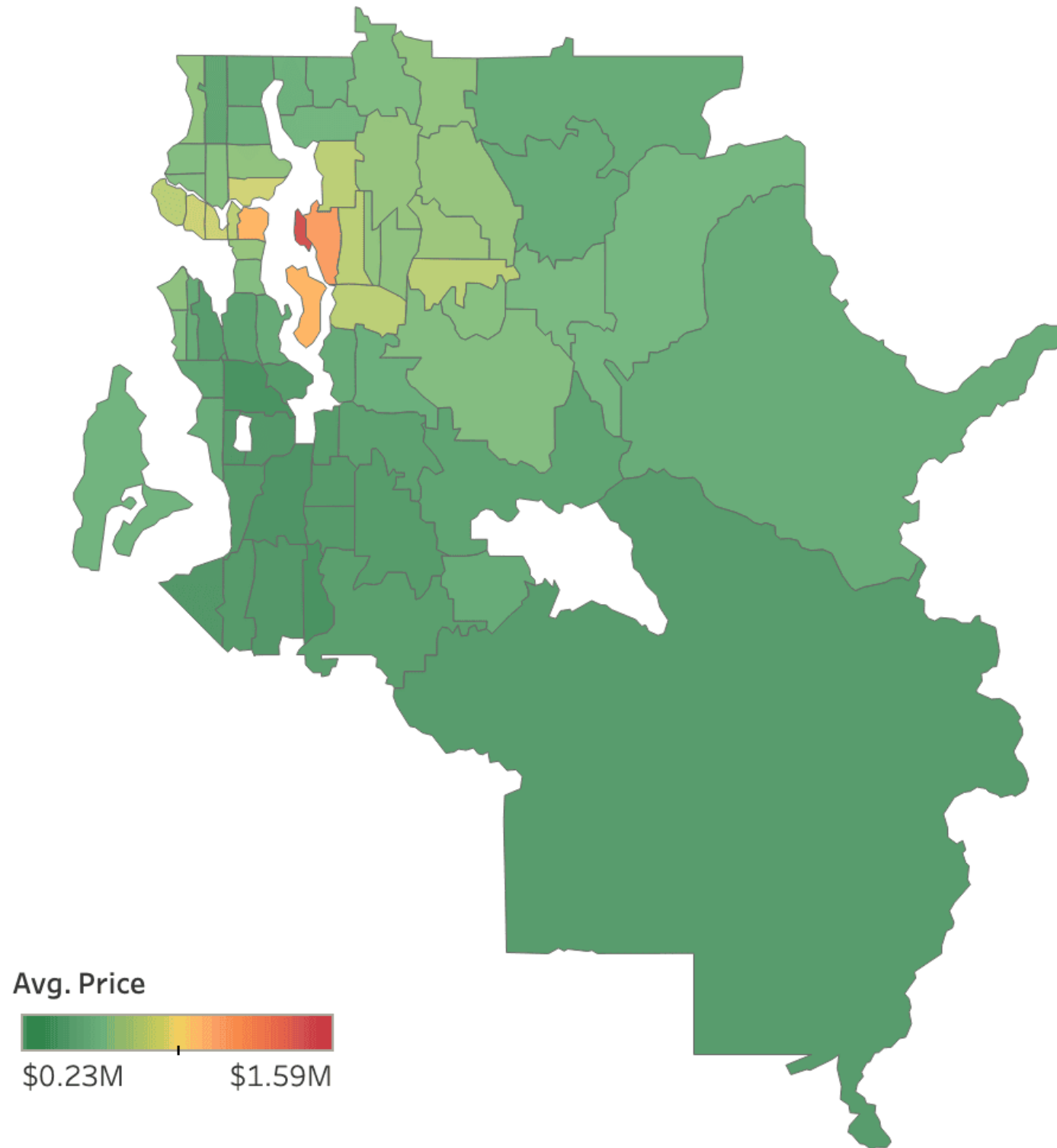
'SQFT_LIVING
15'



'SQFT_LOT1
5'

LIBRARIES	PANDAS	MATPLOTLIB	SEABORN	SCIKIT.LEARN
DATA CLEANING	DEALING WITH NULL VALUES	DROP COLUMNS	HANDLING OUTLIERS	
EDA	MY SQL	PYTHON	TABLEAU	
DATA MODELLING	PREDICTION MODELS	MODEL VALIDATION	MODEL IMPROOVEMENT	

4.1 DISTRICT

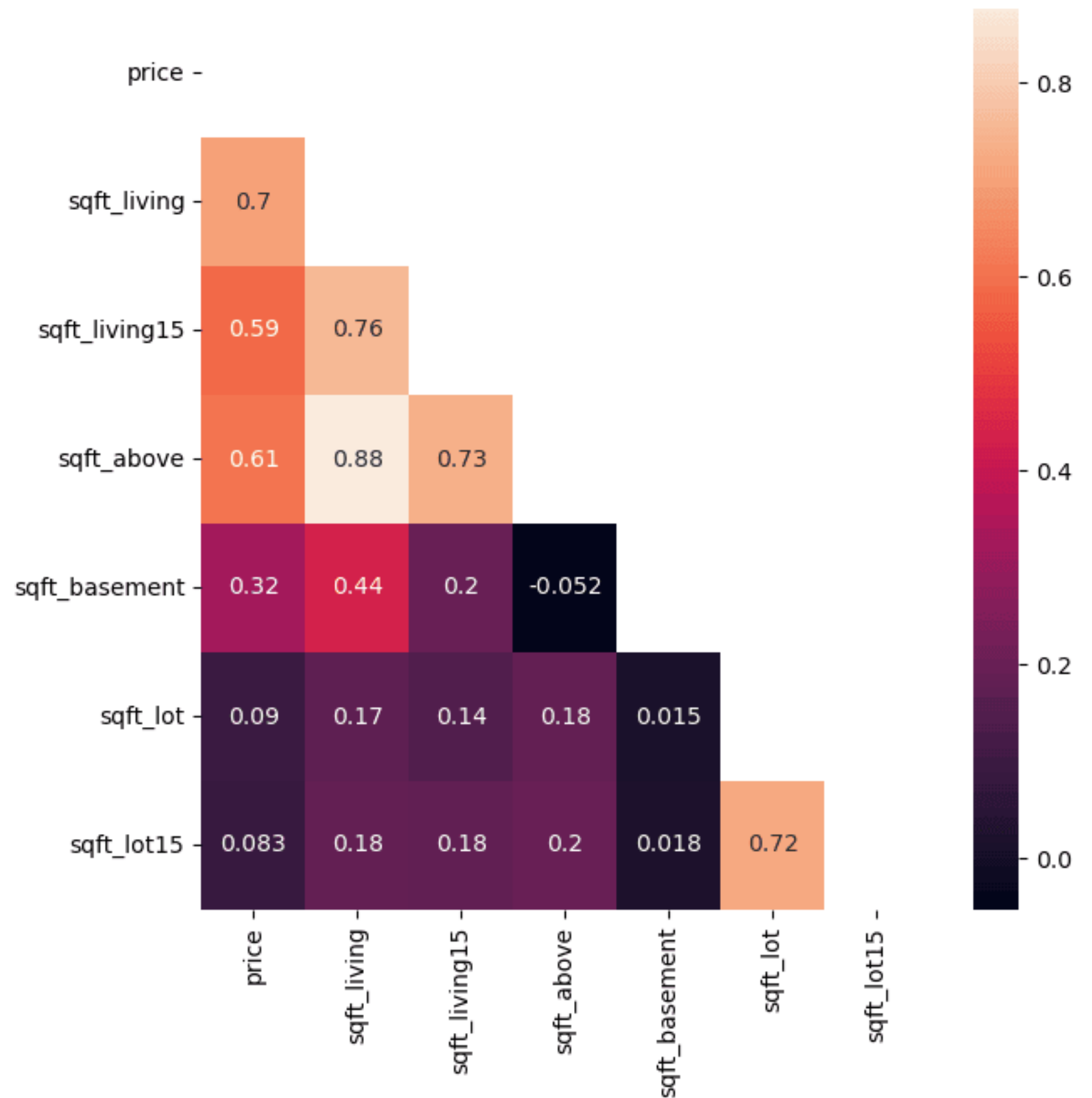


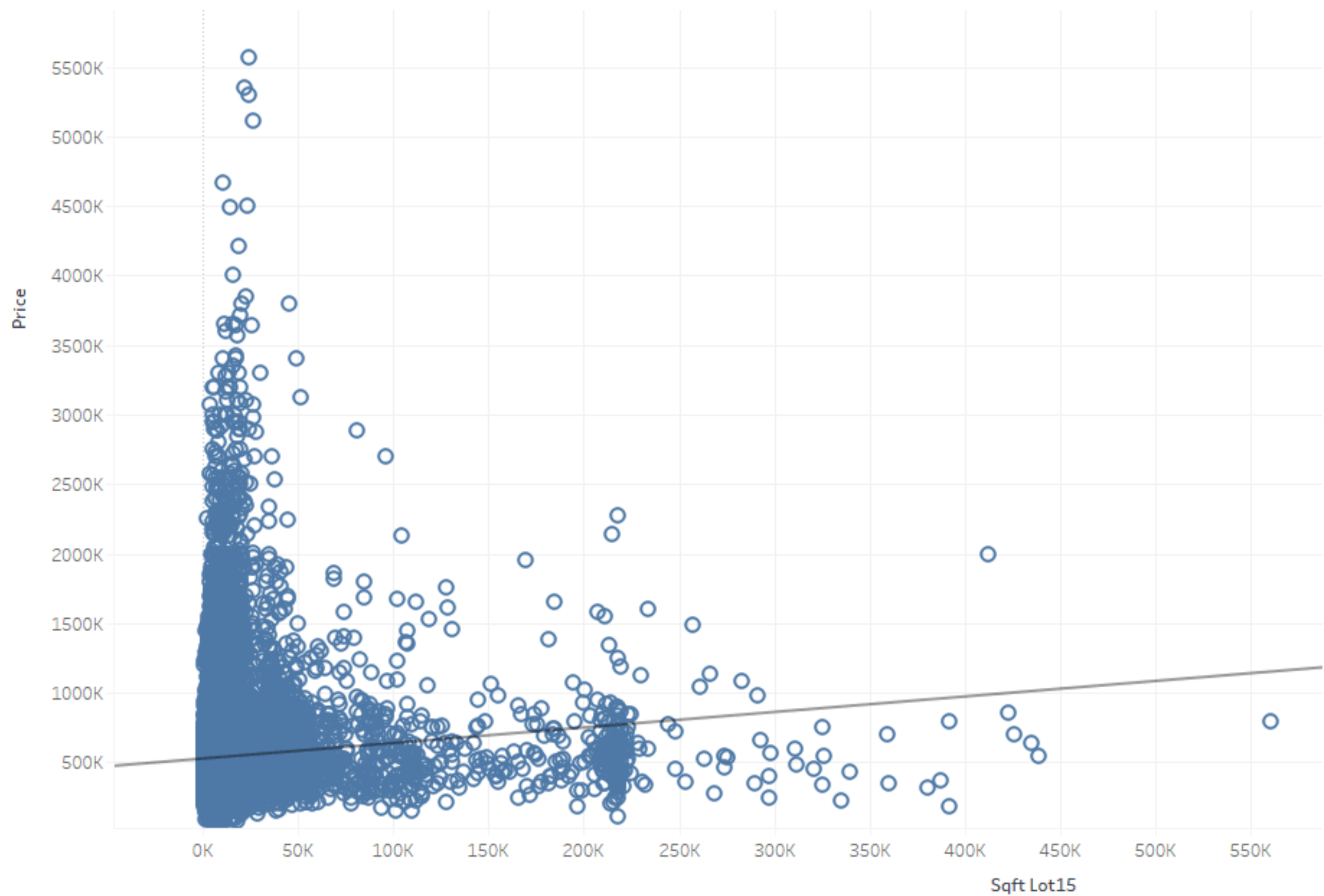
4.2 SIZE

```
corr_matrix["price"].sort_values(ascending=False)
```

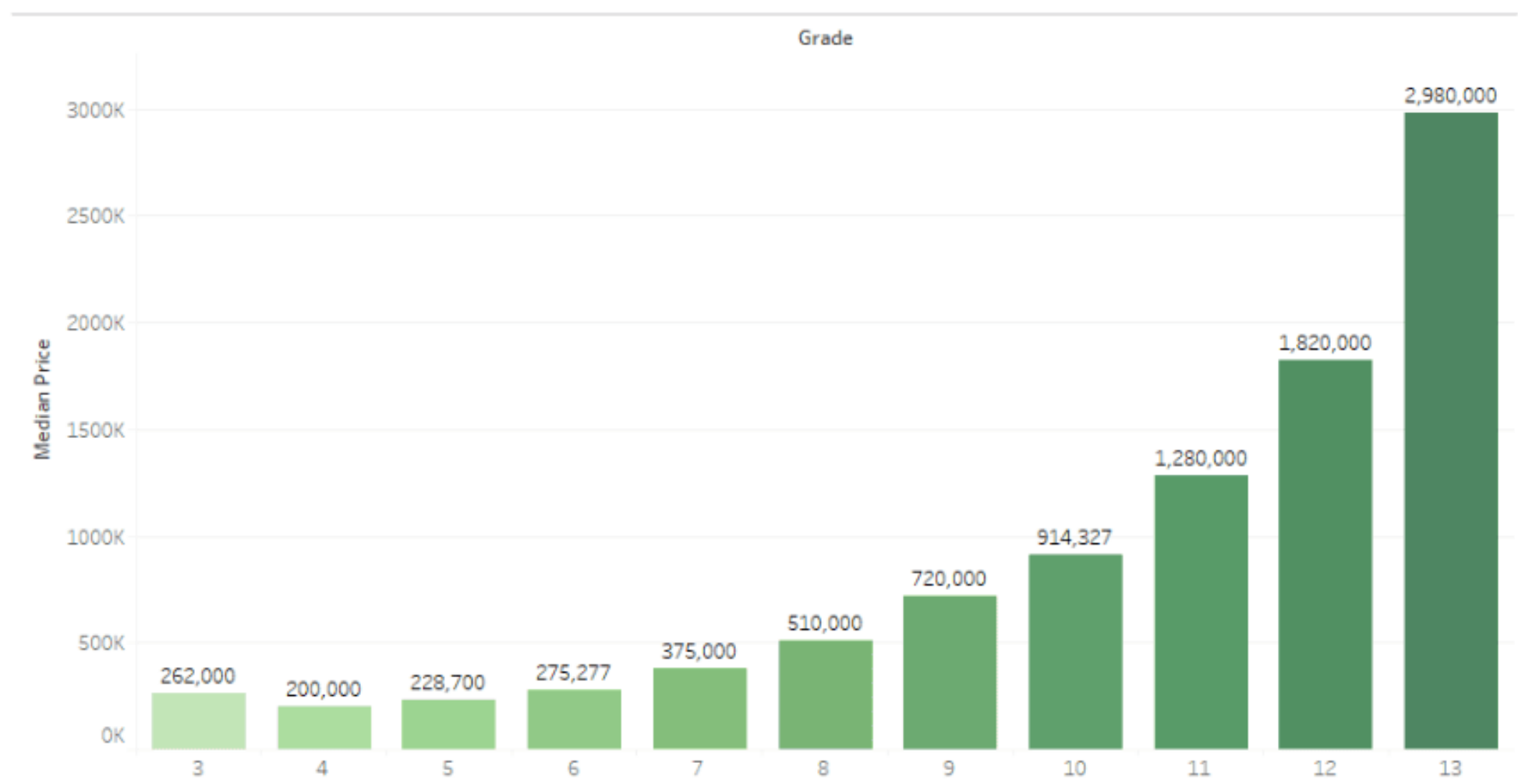
price	1.000000
sqft_living	0.701917
grade	0.667951
sqft_above	0.605368
sqft_living15	0.585241
bathrooms	0.525906
view	0.397370
sqft_basement	0.323799
bedrooms	0.308787
lat	0.306692
waterfront	0.266398
floors	0.256804
yr_renovated	0.126424
sqft_lot	0.089876
sqft_lot15	0.082845
yr_built	0.053953
condition	0.036056
long	0.022036
zipcode	-0.053402

Name: price, dtype: float64





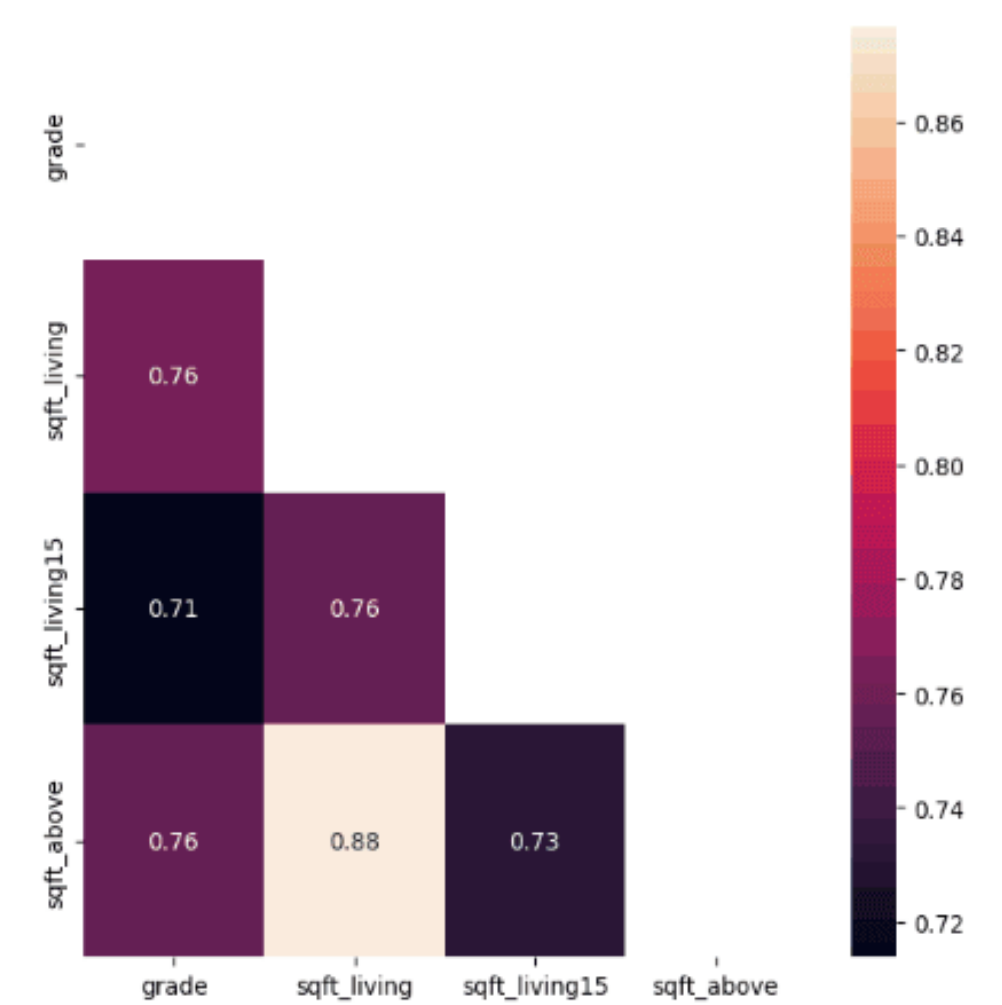
4.3 GRADE



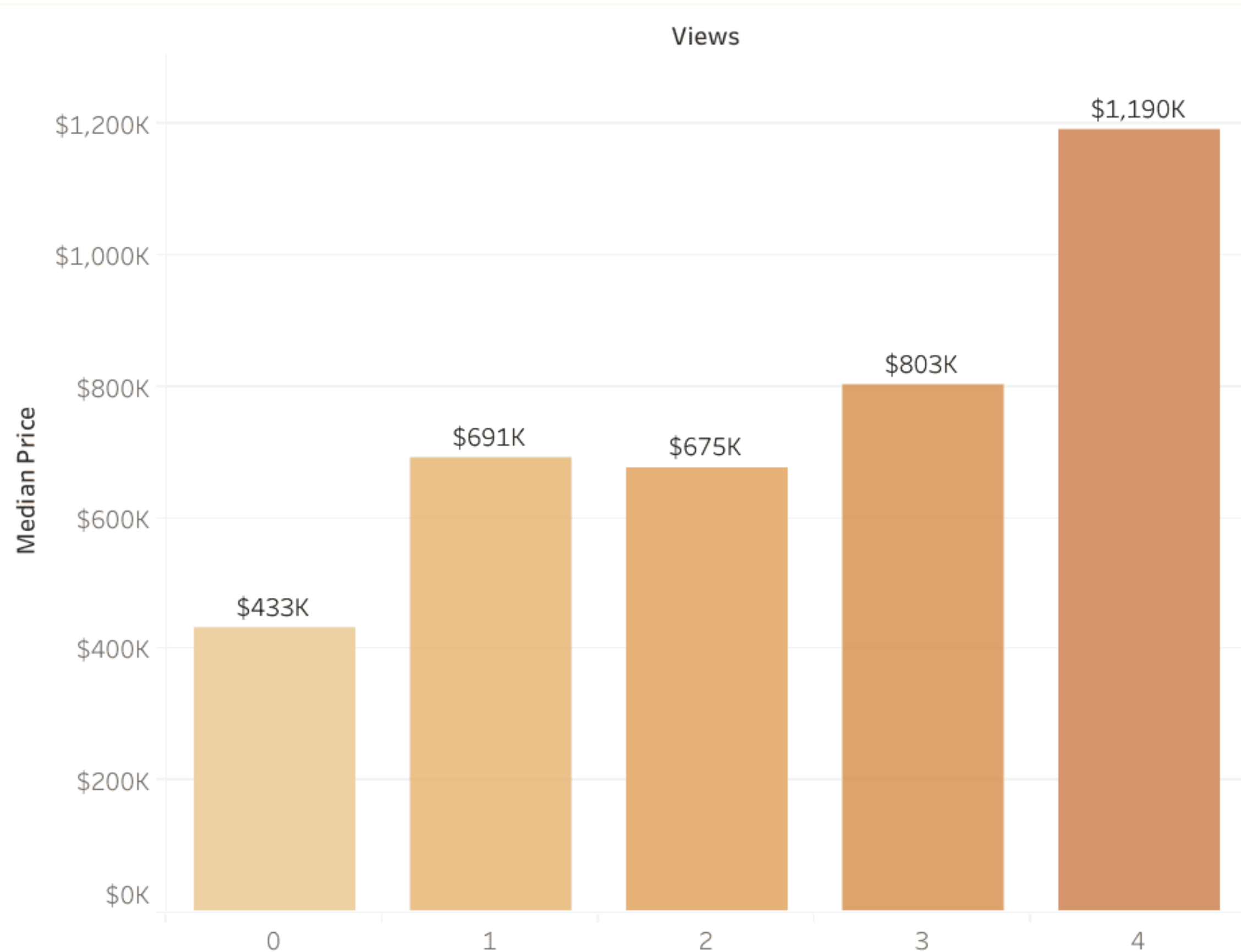
```
corr_matrix["price"].sort_values(ascending=False)
```

price	1.000000
sqft_living	0.701917
grade	0.667951
sqft_above	0.605368
sqft_living15	0.585241
bathrooms	0.525906
view	0.397370
sqft_basement	0.323799
bedrooms	0.308787
lat	0.306692
waterfront	0.266398
floors	0.256804
yr_renovated	0.126424
sqft_lot	0.089876
sqft_lot15	0.082845
yr_built	0.053953
condition	0.036056
long	0.022036
zipcode	-0.053402

Name: price, dtype: float64



4.4 VIEWS

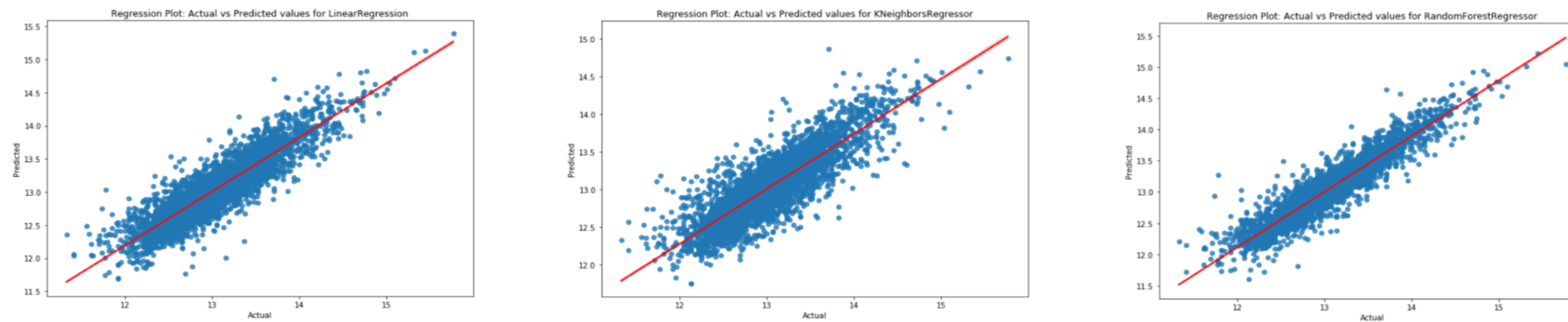


5. HOUSE PRICE PREDICTION MODEL

- Employed Prediction Models: Linear Regressor, KNN Regressor & Random Forest Regressor.
- Model Validation: R2 Score, MAE, MSE
- Model Improvement: Scaling (Log Transform - to reduce outliers), Feature Selection (Avoid Multicollinearity using correlation matrix)



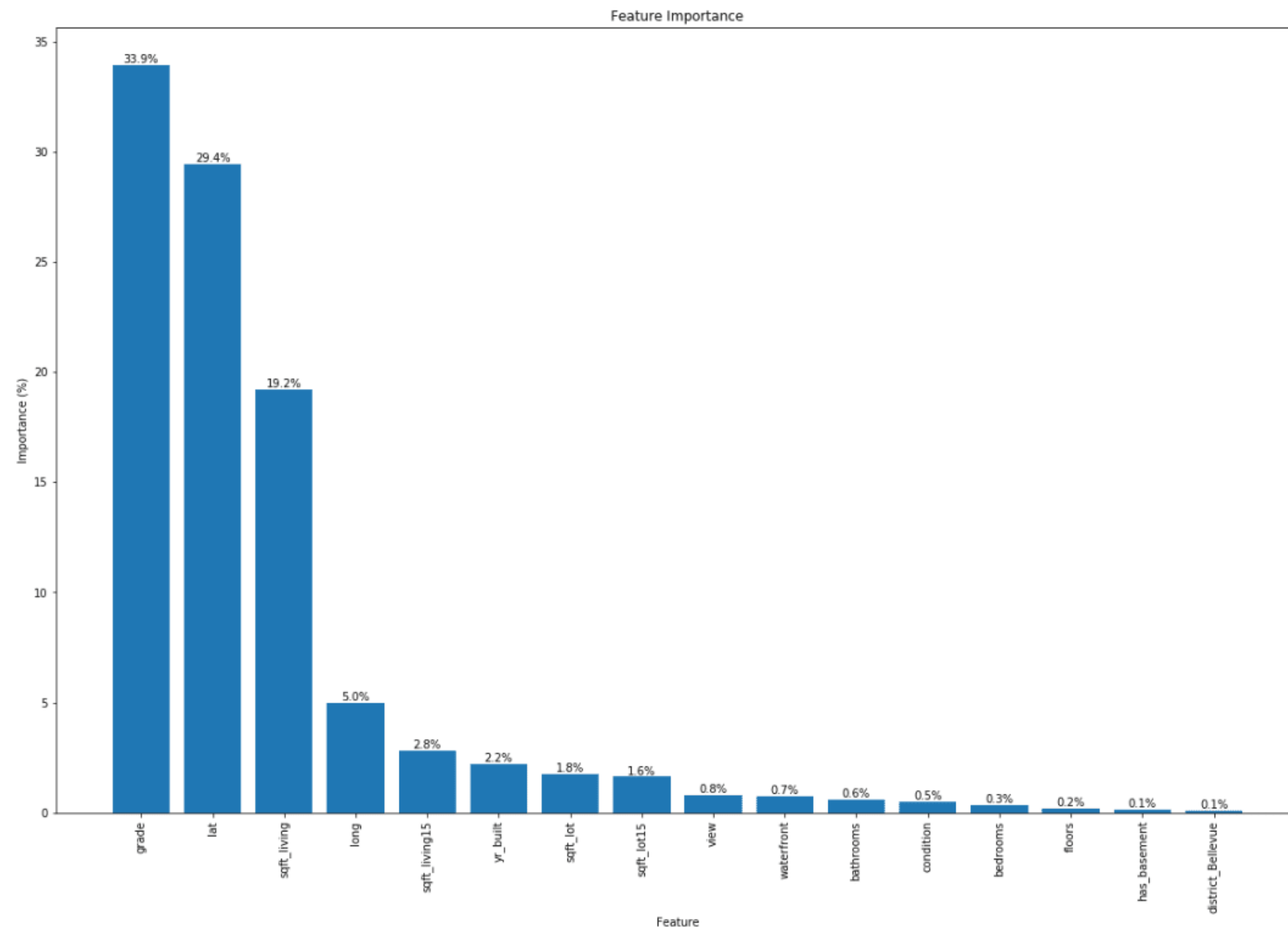
Cross- Examination of Different Algorithms



	Model	r2	mse	mae
0	LinearRegression	0.742048	3.358960e+10	112183.571208
1	KNeighborsRegressor	0.498109	6.535432e+10	156184.477917
2	RandomForestRegressor	0.882405	1.531284e+10	68196.365438

	Model	r2	mse	mae
0	LinearRegression	0.809026	0.052012	0.171247
1	KNeighborsRegressor	0.734856	0.072212	0.194961
2	RandomForestRegressor	0.887906	0.030529	0.123204

Feature Importance

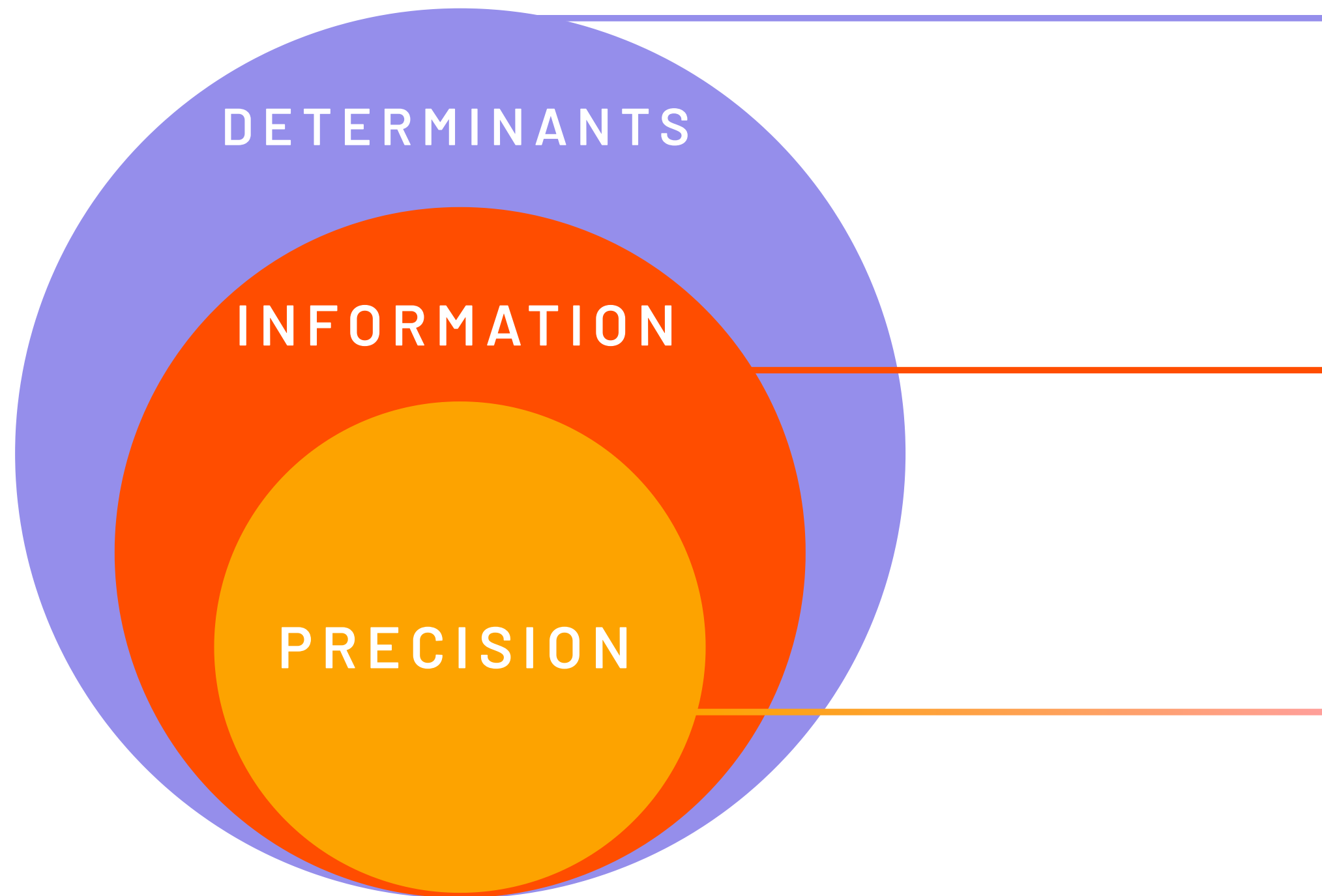


Important house price factors:

- Grade
- Location (lat & long)
- Size (sqft_living & sqft_living15, sqft_lot)
- View
- Waterfront
- Bathrooms, Bedrooms,
- Condition, Floors, Basement



6. CONCLUSIONS



Our Forest Random Regression analysis indicates that property grade, location, and size are key price determinants.



Yet, we lack specific information about the grading system and our dataset omits influential external factors such as proximity to amenities and crime rates.



To increase the precision of our housing price estimation model, we recommend incorporating these social-economic indicators into further analyses.