

EDA Project

Client: Jennifer Montgomery

Requirements: High budget, wants to show off, timing within a month, waterfront, renovated, high grades, resell within 1 year (buyer)

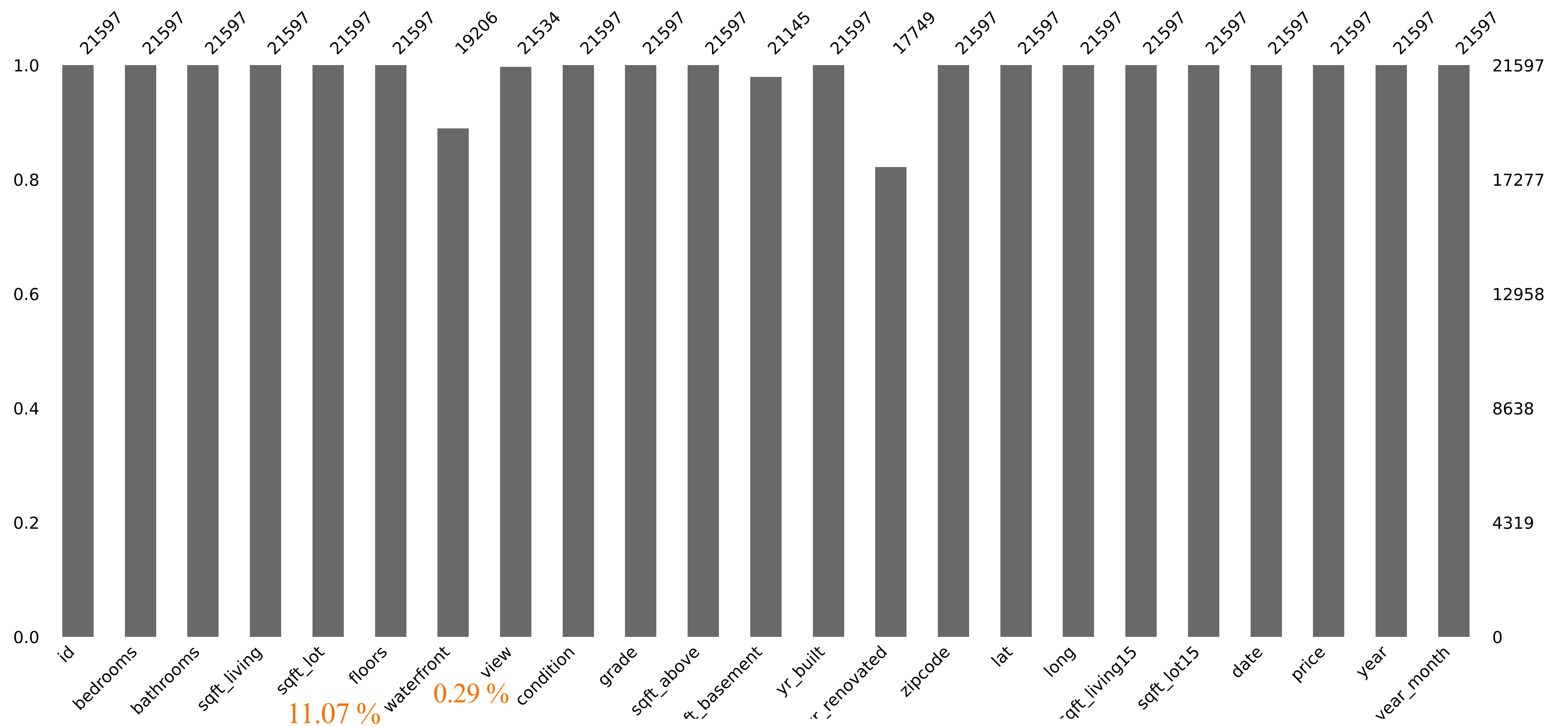
Workflow

- Understanding data
- Data cleaning
- Feature engineering
- Data visualization
- Insights regarding over all data
- Data preparation for client
- Recommendation to client

Total no of rows

Missing data

missingno



- waterfront

- yr_renovated

- view

- sqft_basement

- Over all missing values in data frame : 1.36 %

- No duplicate data sets

Missing data and imputation

Column 'year of renovation'

Fill Method

Column yr_renovation weird?

NaN

Zero

19900 → 1990

```
df_n0[["waterfront", "yr_renovated", "view", "sqft_basement"]] =
```

```
df_n0[["waterfront", "yr_renovated", "view", "sqft_basement"]].fillna(0)
```

```
df_n0[["waterfront", "yr_renovated", "view", "sqft_basement"]].reset_index()
```

```
df_n0['yr_renovated'] = df_n0['yr_renovated'].apply(lambda x: str(x)[:4])
```

```
df_n0['yr_renovated'] = df_n0.yr_renovated.replace('nan', np.NaN)
```

```
df_n0['yr_renovated'] = df_n0['yr_renovated'].astype(float).astype('Int64')
```

```
df_n0 = df_n0.dropna().reset_index()
```

Hypothesis

Client requirements: High budget, wants to show off, timing within a month, waterfront, renovated, high grades, resell within 1 year (buyer)

- Additional data: **age of the house**
- Additional data: **age of the house renovation**
- Additional data: **price per sqft**

Feature engineering

Additional columns:

- change "date" dtype to datetime with format %Y/%m/%d.

add two new columns: (1) year-month (2) year

- Additional column age of the house

```
df_n0['house_age'] = df_n0['date'].dt.year - df_n0['yr_built']
```

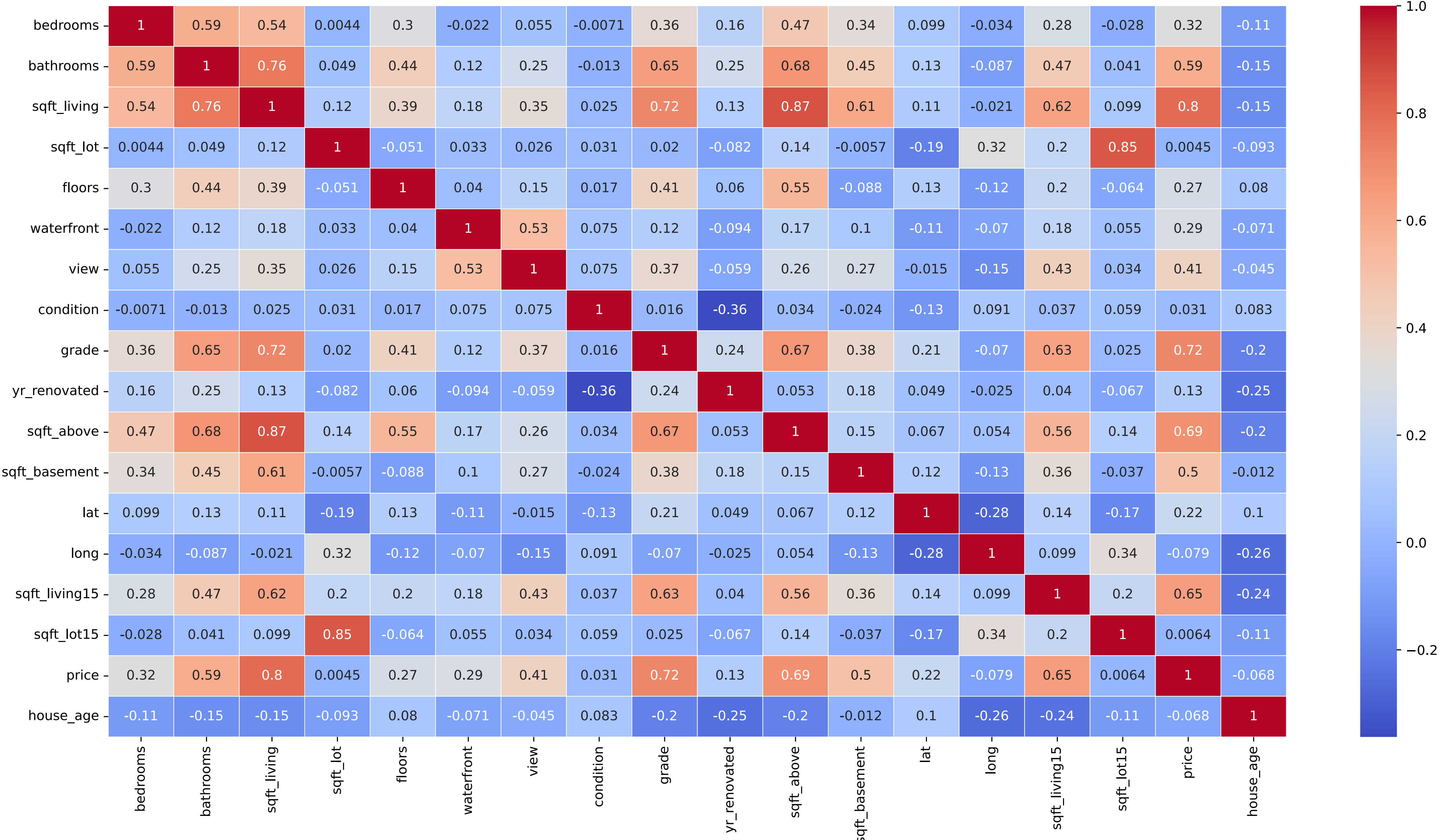
- Additional column age of the house renovation

```
df_n0['house_age_yr'] = df_n0['date'].dt.year - df_n0['yr_renovated']
```

- Additional column price per sqft

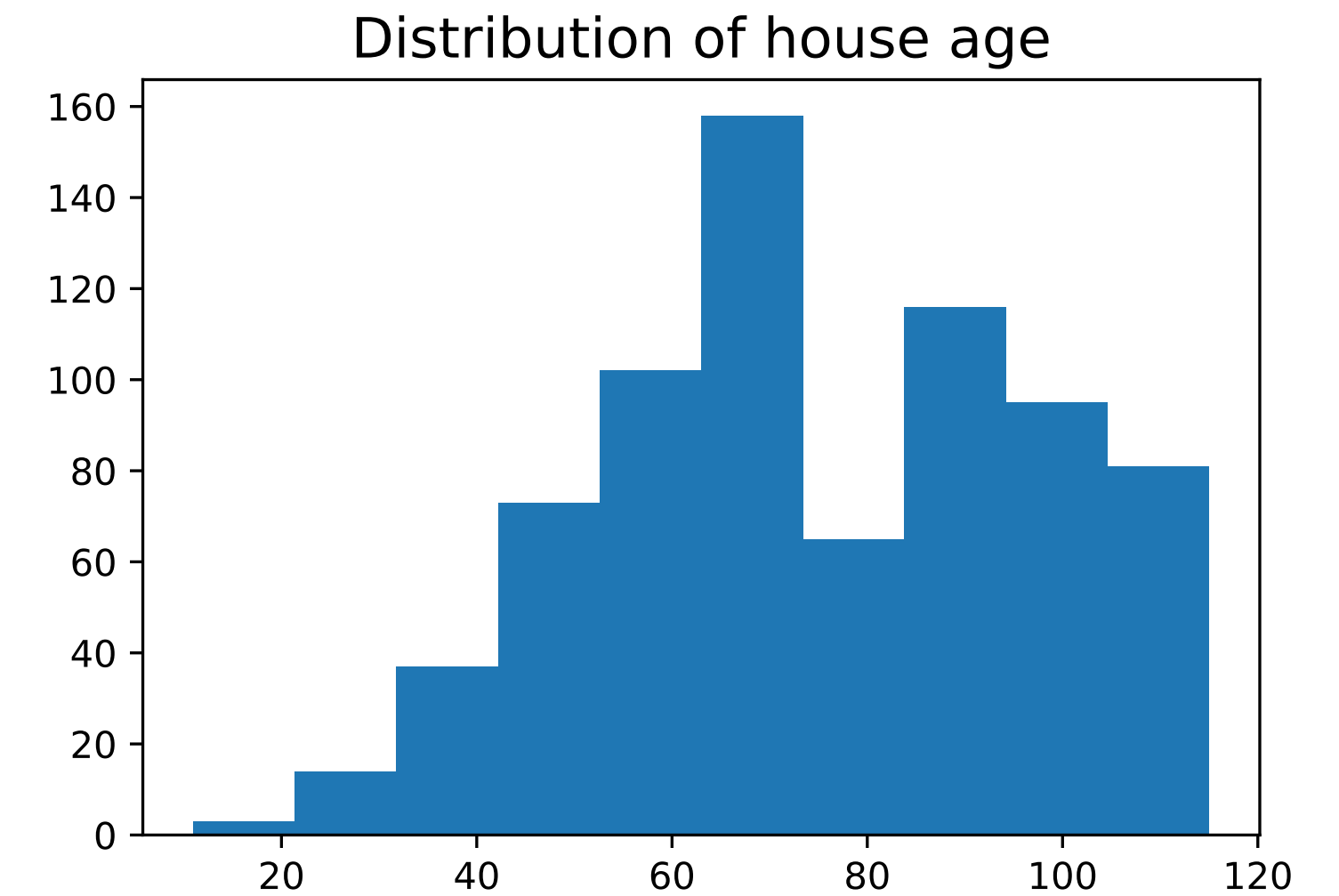
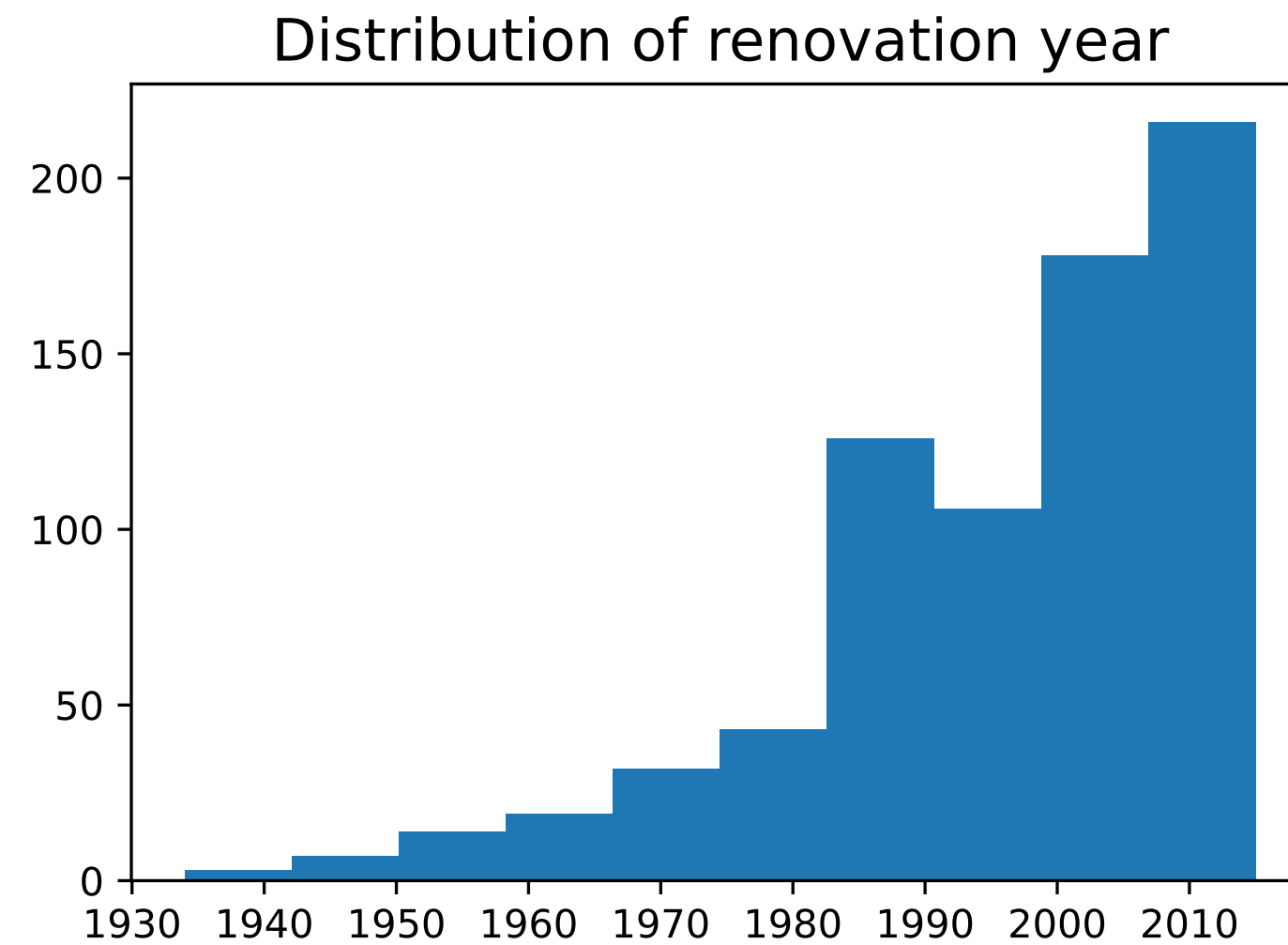
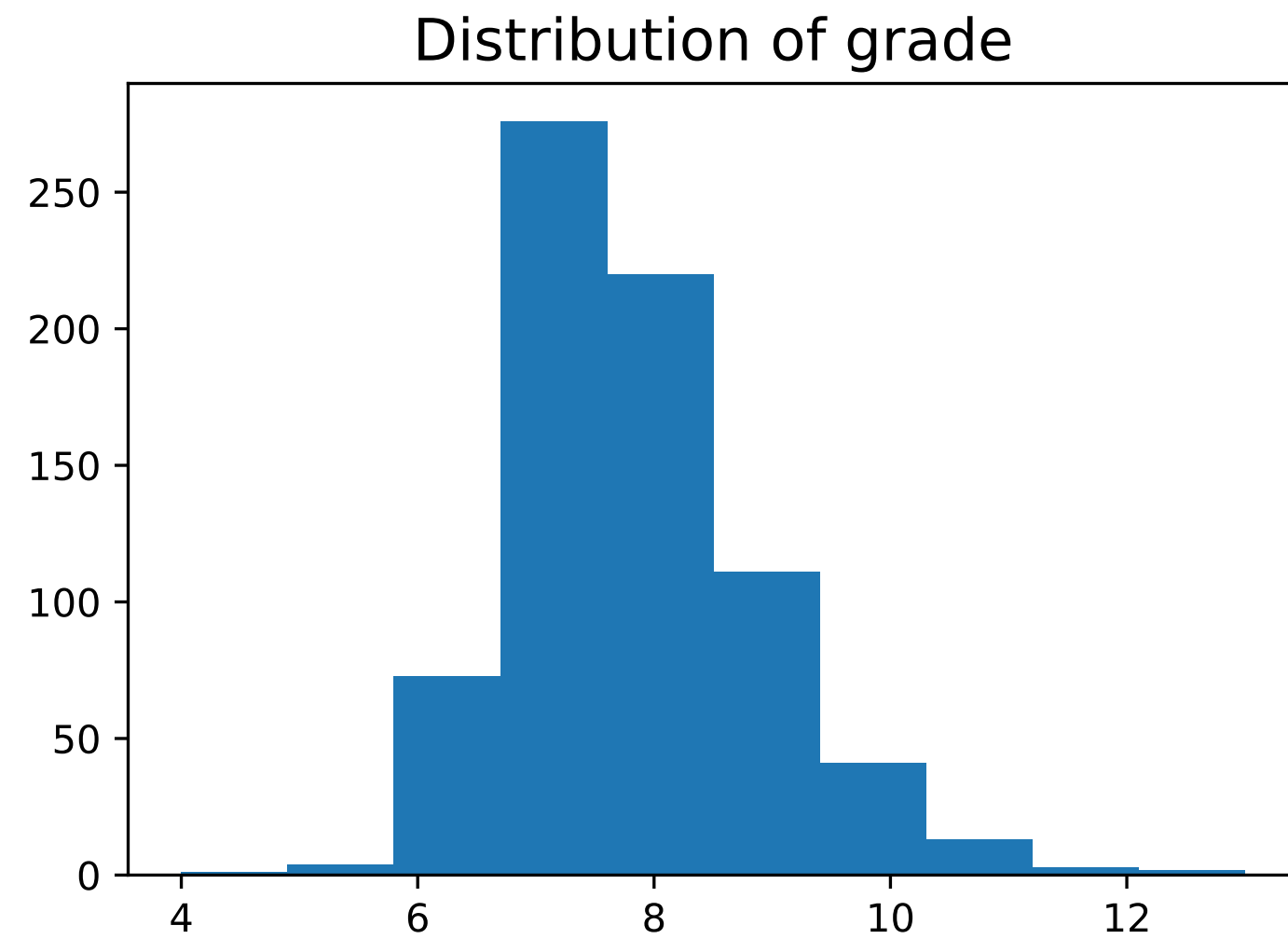
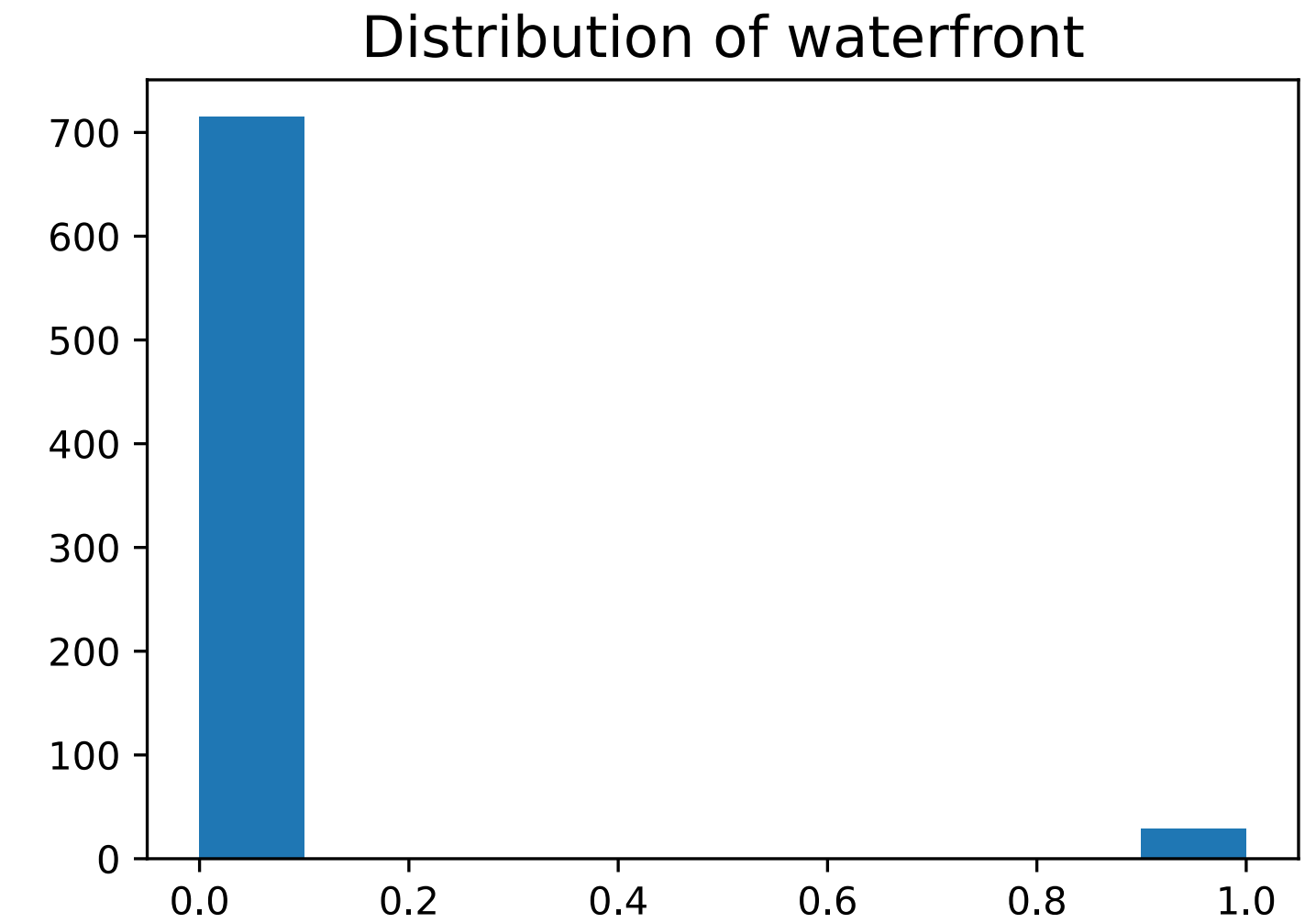
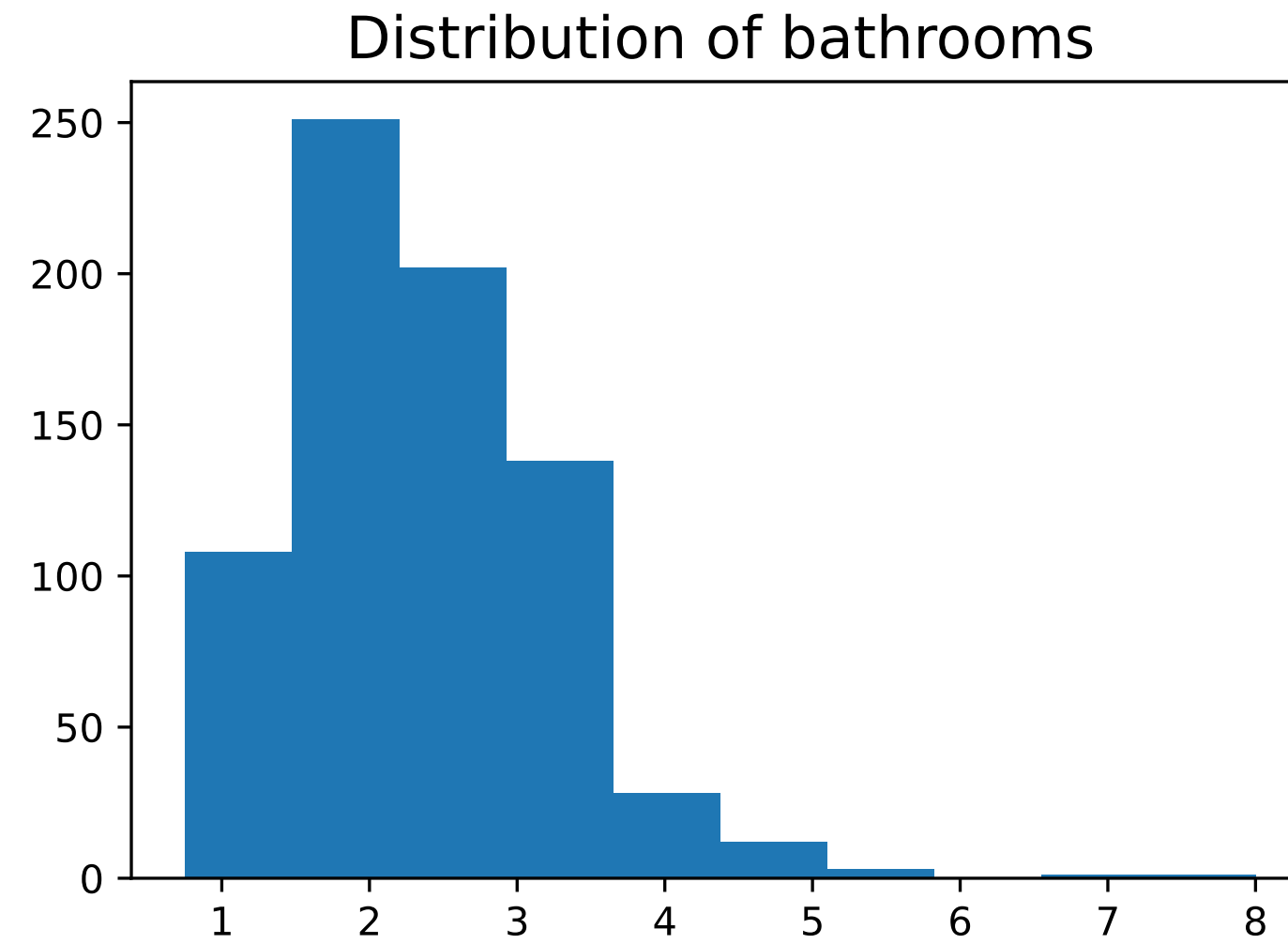
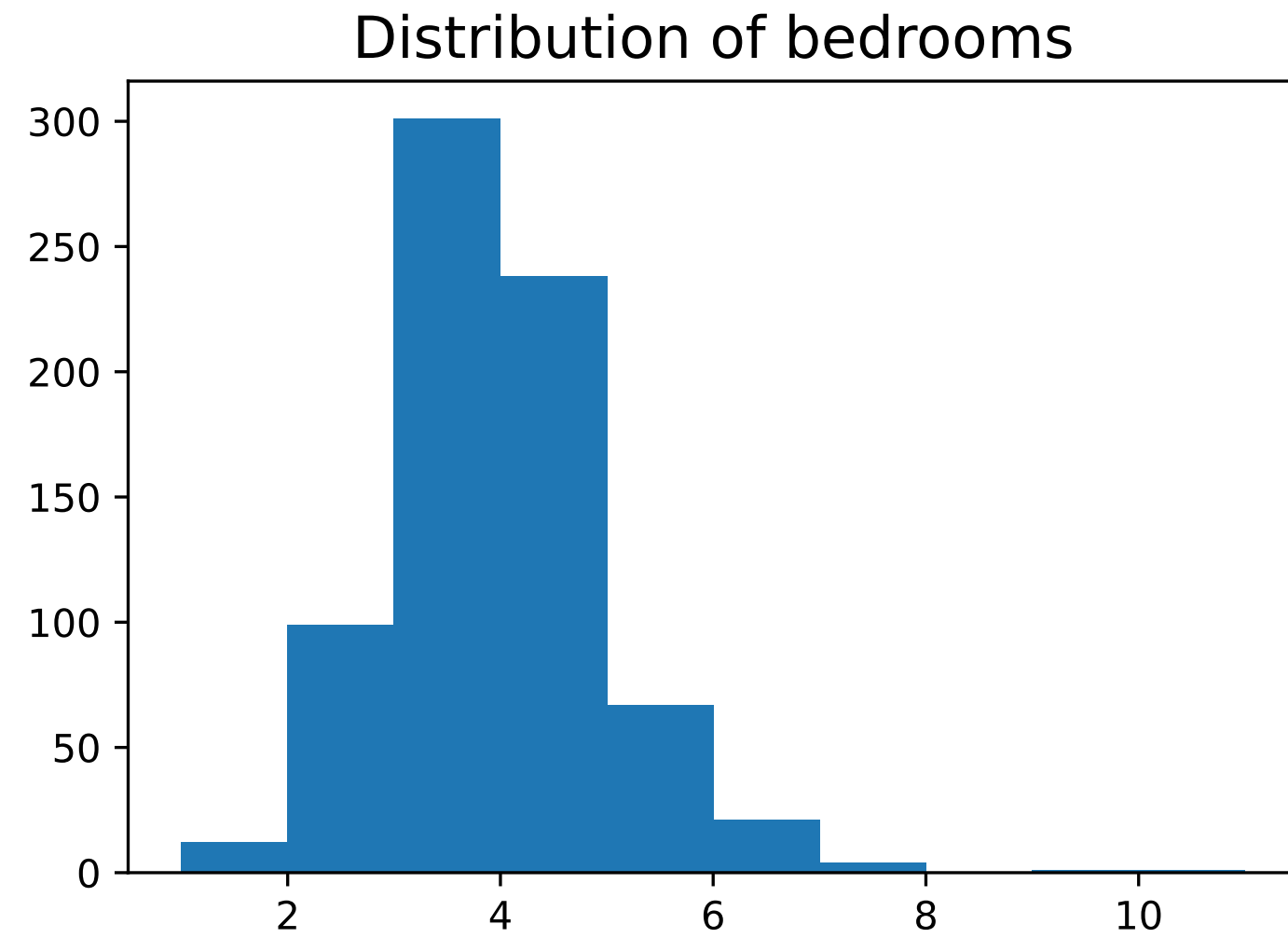
```
df_n0['price_by_sqft'] = df_n0['price'] / df_n0['sqft_lot']
```

Correlation matrix

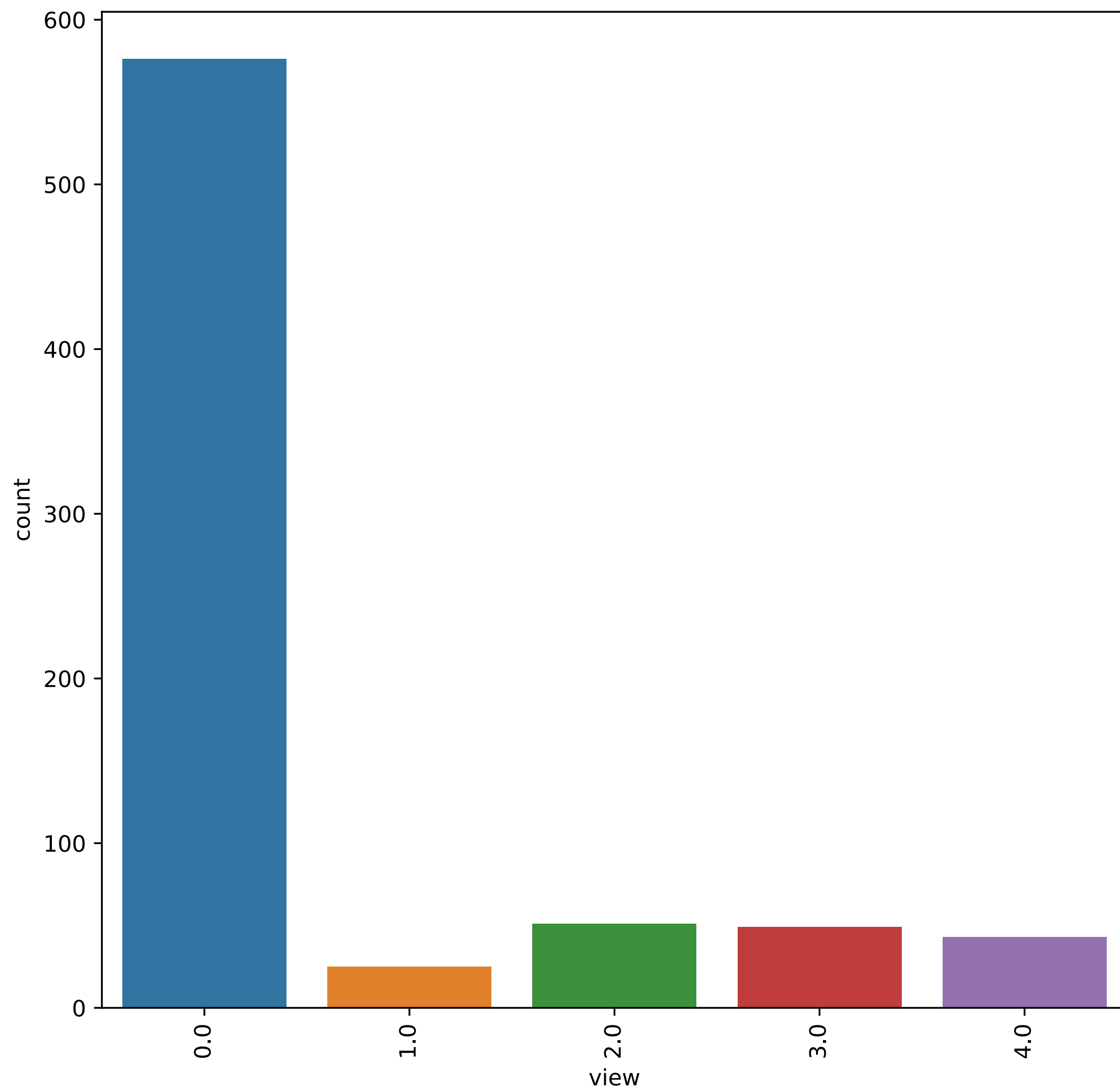


Histograms

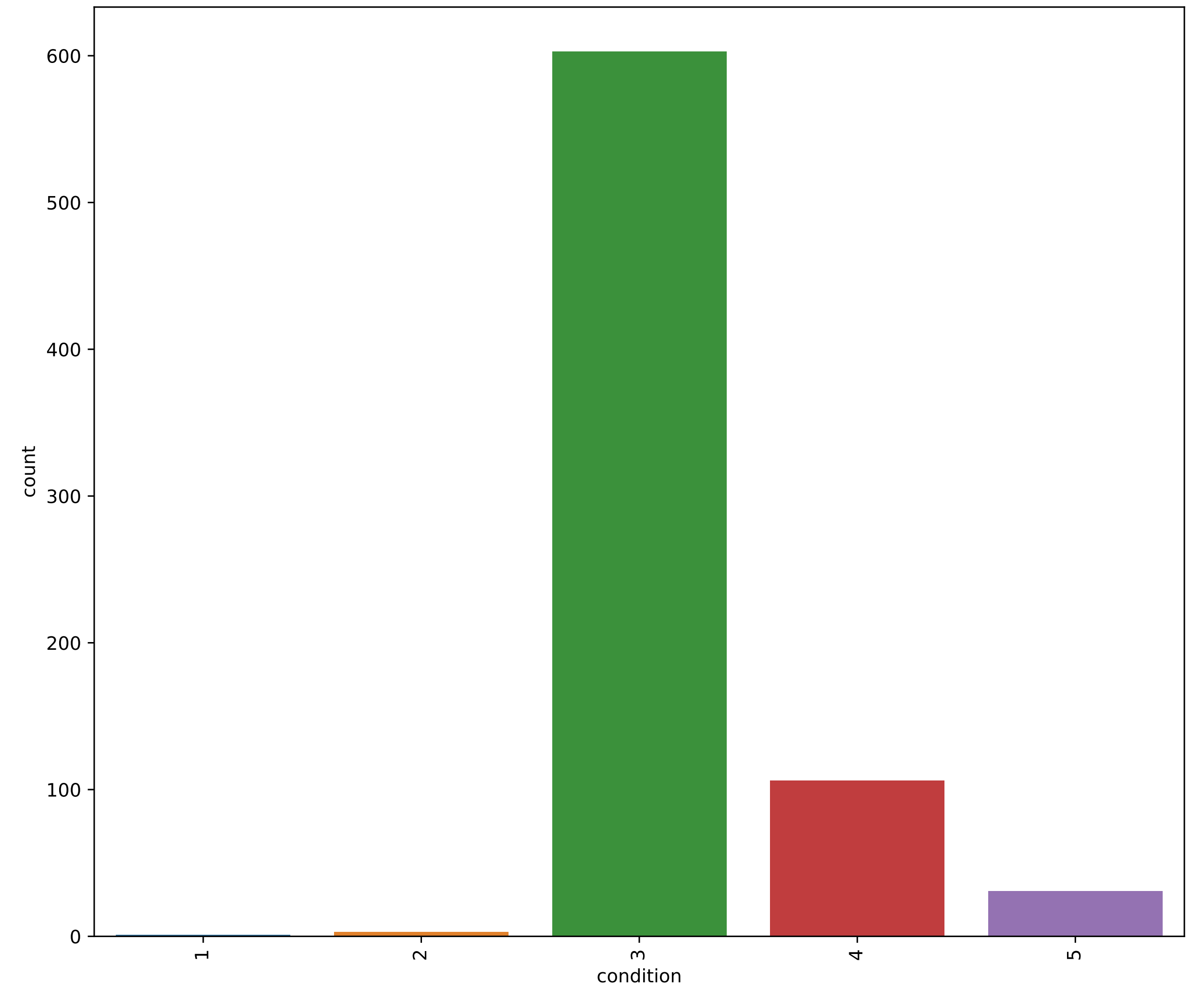
Distribution of numeric columns

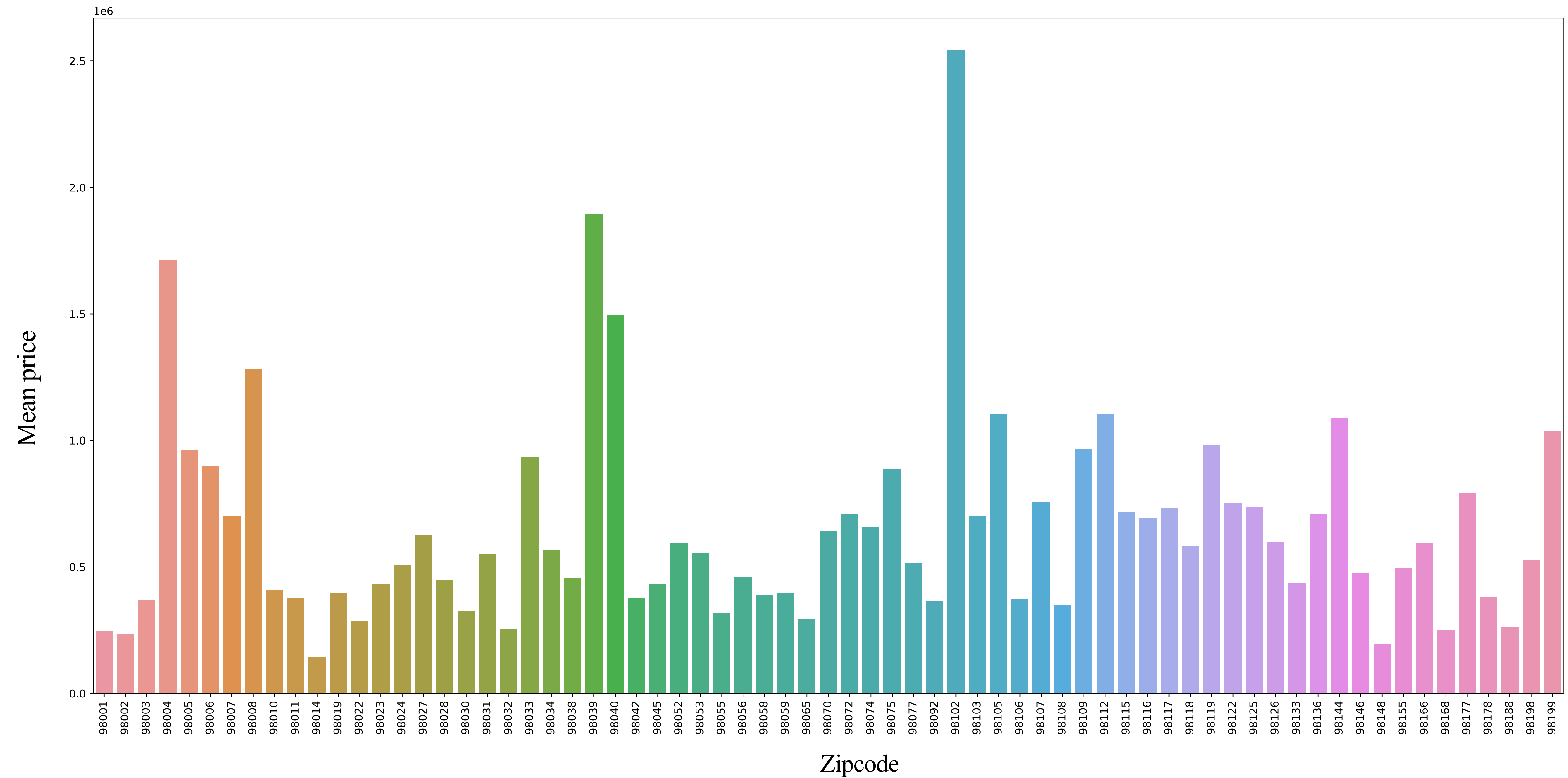


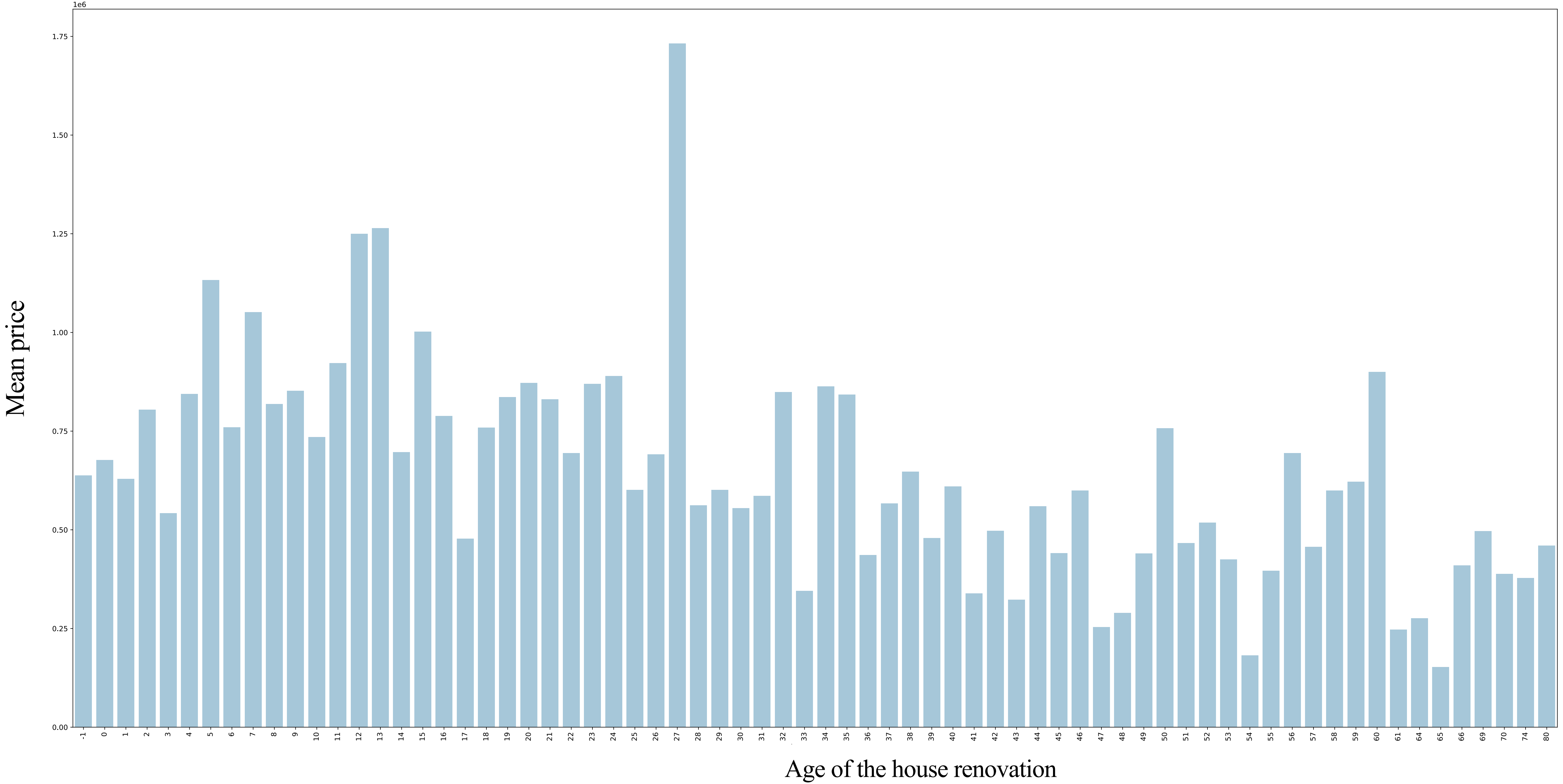
Sales count versus view



Sales count versus conditions







Summary

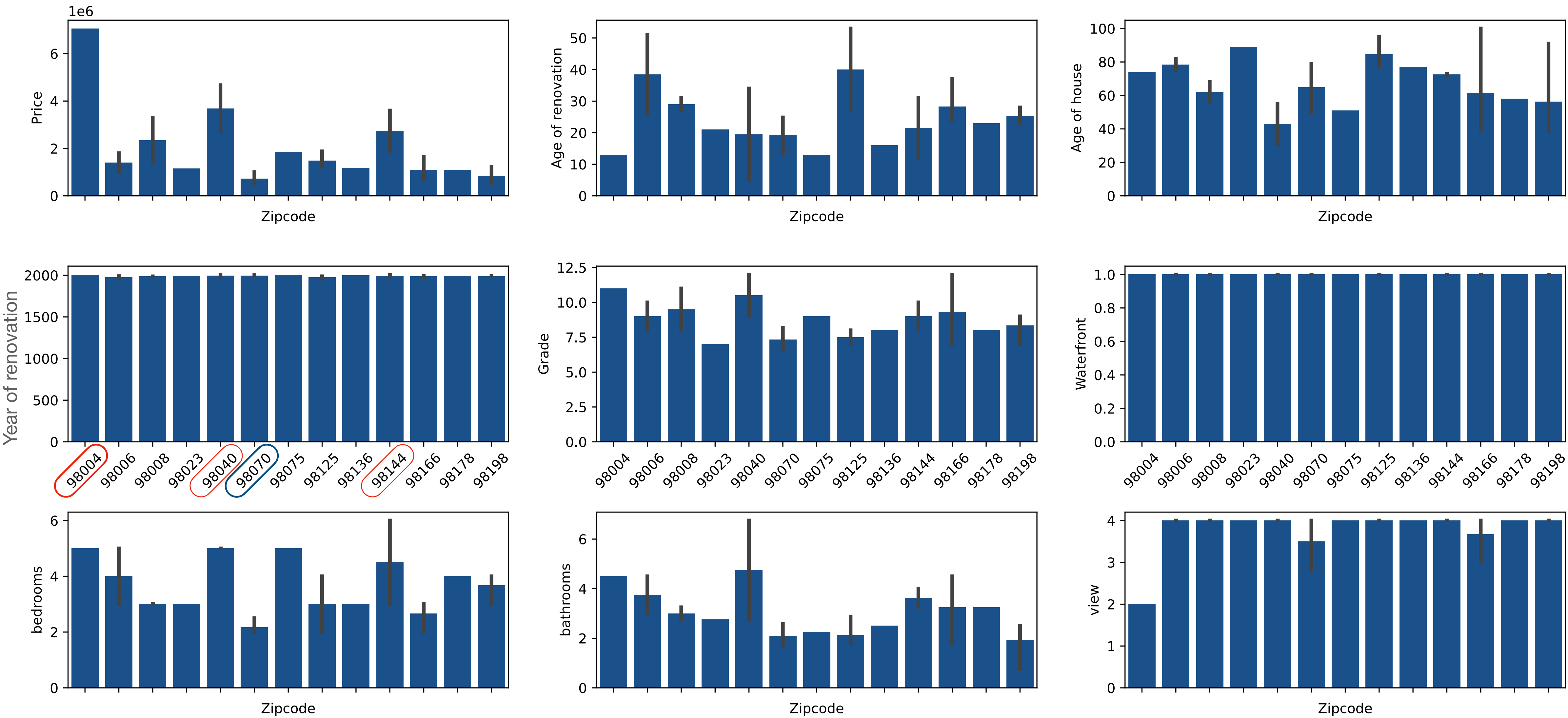
- **Zipcode is an identifier of houses**
- **Zipcodes 98039, 98004, 98040 - highest price . Zipcodes 98002, 98168, 98032 - lowest price**
- **Majority of houses sold in grade between 6 and 9**
- **Most of the house sold has 0 view. Buyers do not care about view**
- **Majority of houses doesn't have a waterfront**
- **Buyer prioritize 3 to 4 bedrooms and 1 to 2.5 bathrooms**
- **Living size is highly correlated with space of nearest 15 neighbors**
- **House with a waterfront has more price**
- **Buyer prioritize house with 2 to 3 floors. House with 2.5 floors has maximum price**
- **Most of the houses are in average condition**
- **Grade of a house is an important factor to decide the price of house**
- **Houses with less than one year of age has the highest median price**
- **The price of house is highly correlated with the living space**

Final data preparation

Pseudo code:

```
df_nump = []
for ind in df_n0.index:
    if df_n0['waterfront'][ind]==1:
        df_nump.append(
            [
                ind,
                df_n0['zipcode'][ind],
                df_n0['waterfront'][ind],
                df_n0['grade'][ind],
                df_n0['yr_renovated'][ind],
                df_n0['house_age'][ind],
                df_n0['house_age_yr'][ind],
                df_n0['price'][ind],
                df_n0['bedrooms'][ind],
                df_n0['bathrooms'][ind],
                df_n0['view'][ind]]
        )
df_nump = pd.DataFrame(df_nump)
```

Recommendations to Jennifer Montgomery: Houses details versus zipcode



Thank You