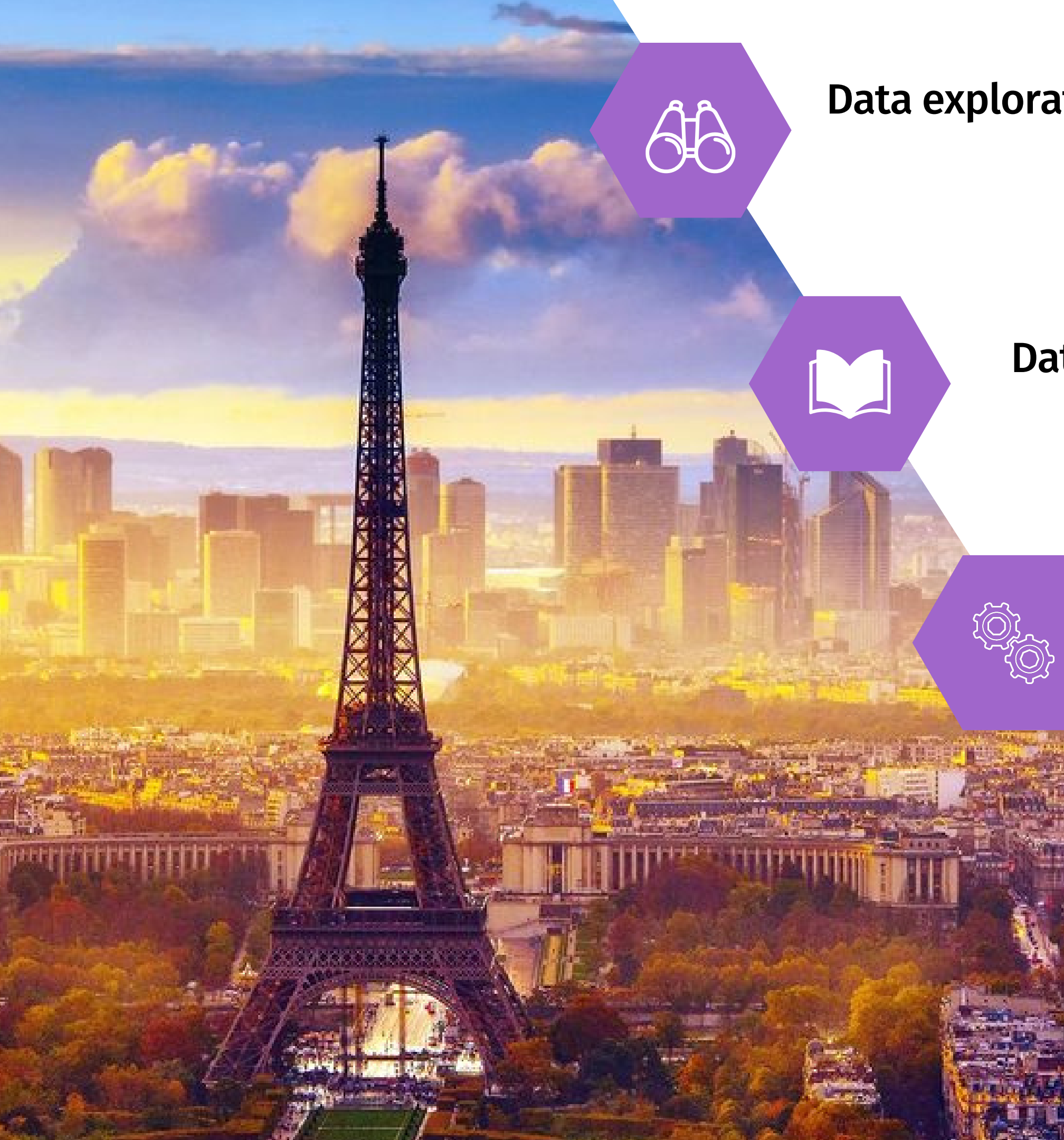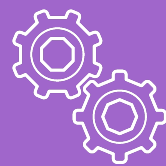**Elliot & Chi**

# Mid Project Presentation

Subject : Housing Price Prediction Model

**Data exploration**

**Data visualization**

**Training the model**

**Analyzing the results**

# Presentation of the Dataset

## Context

Our dataset is composed of :
- houses and their caracteristics
- Their selling prices

## Where ?

United States
Washington state
King's county

## Objective ?

To define a model that can predict the selling price.

# Some key elements

**The average price in King's county to buy a house.**

540 296 $
--> Not very expensive if we compare some other county that are near the ocean

**Luxury houses in the county**

We have 1490 houses with a price that is above 1 million $ which represent only 6% of our data
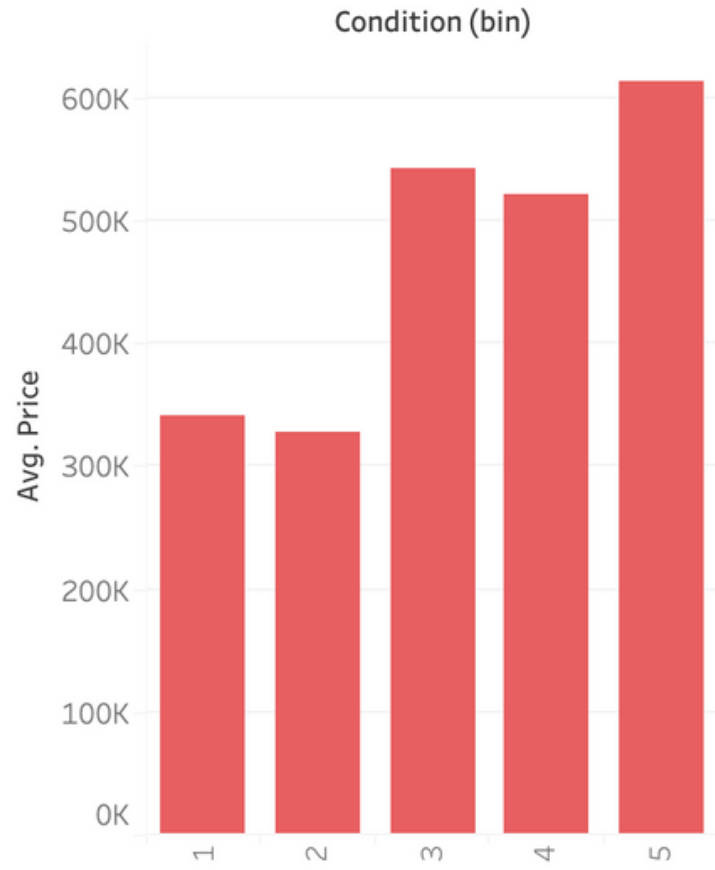
**Recent houses ?**

Most of our houses were built around the 70s. Something to take into consideration is that we only have approximatively 4700 houses that were built in the second millennium.
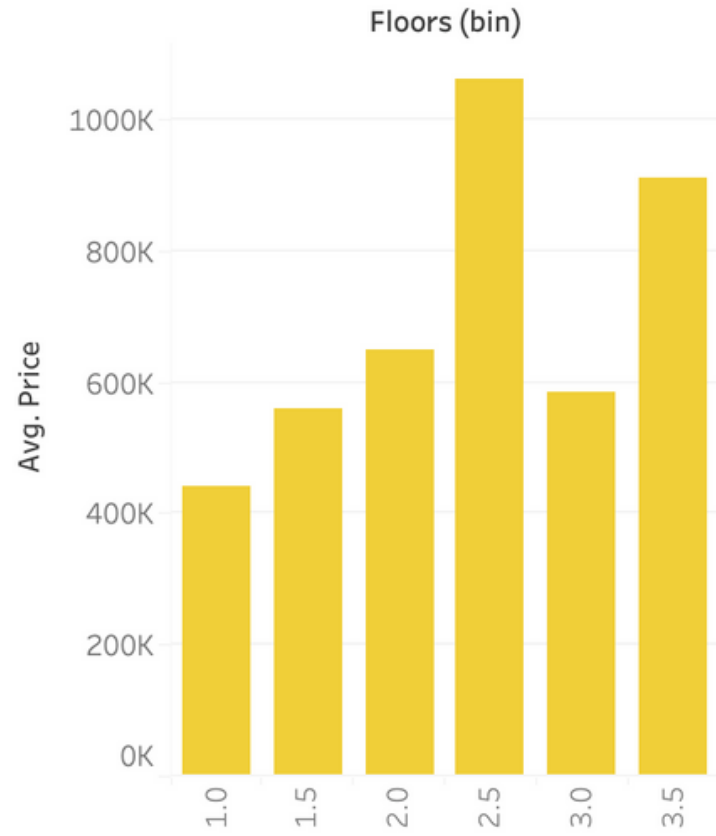--> Not a lot of construction, does the county still attractive ?
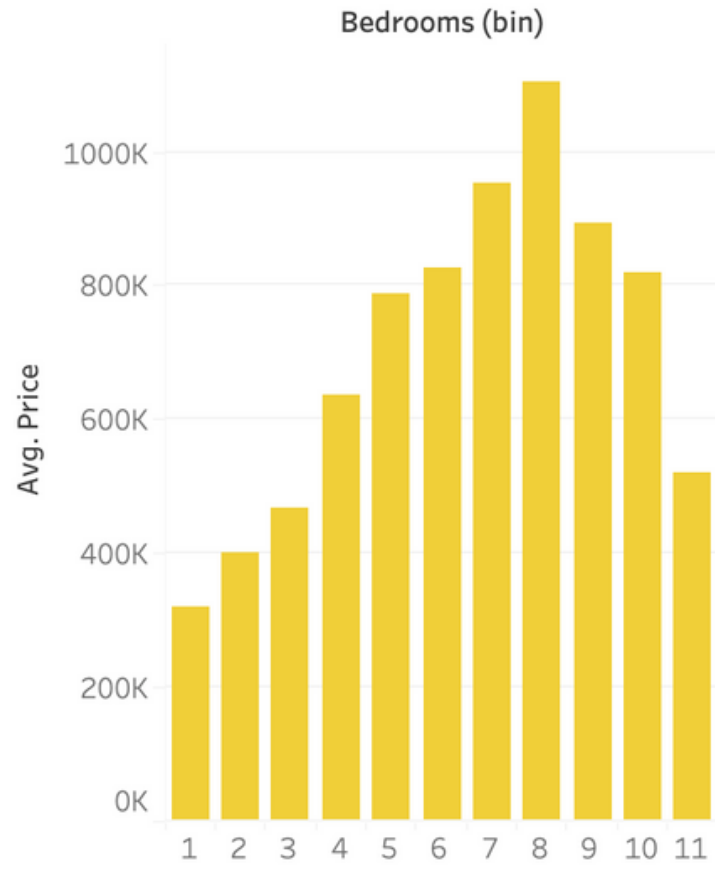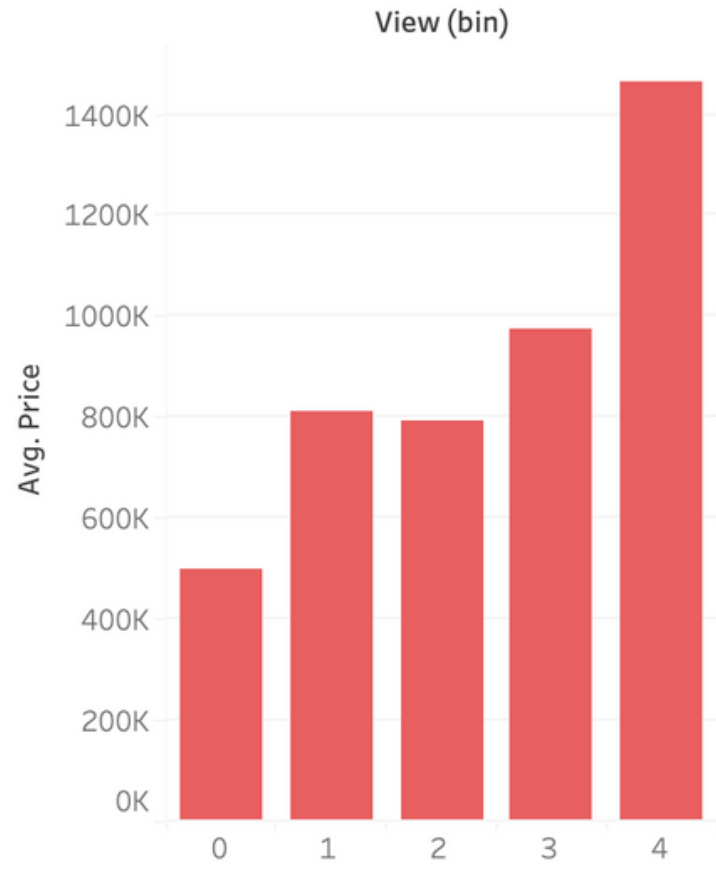
# Different figures...

# The map

- **Higher prices on the coast in general**

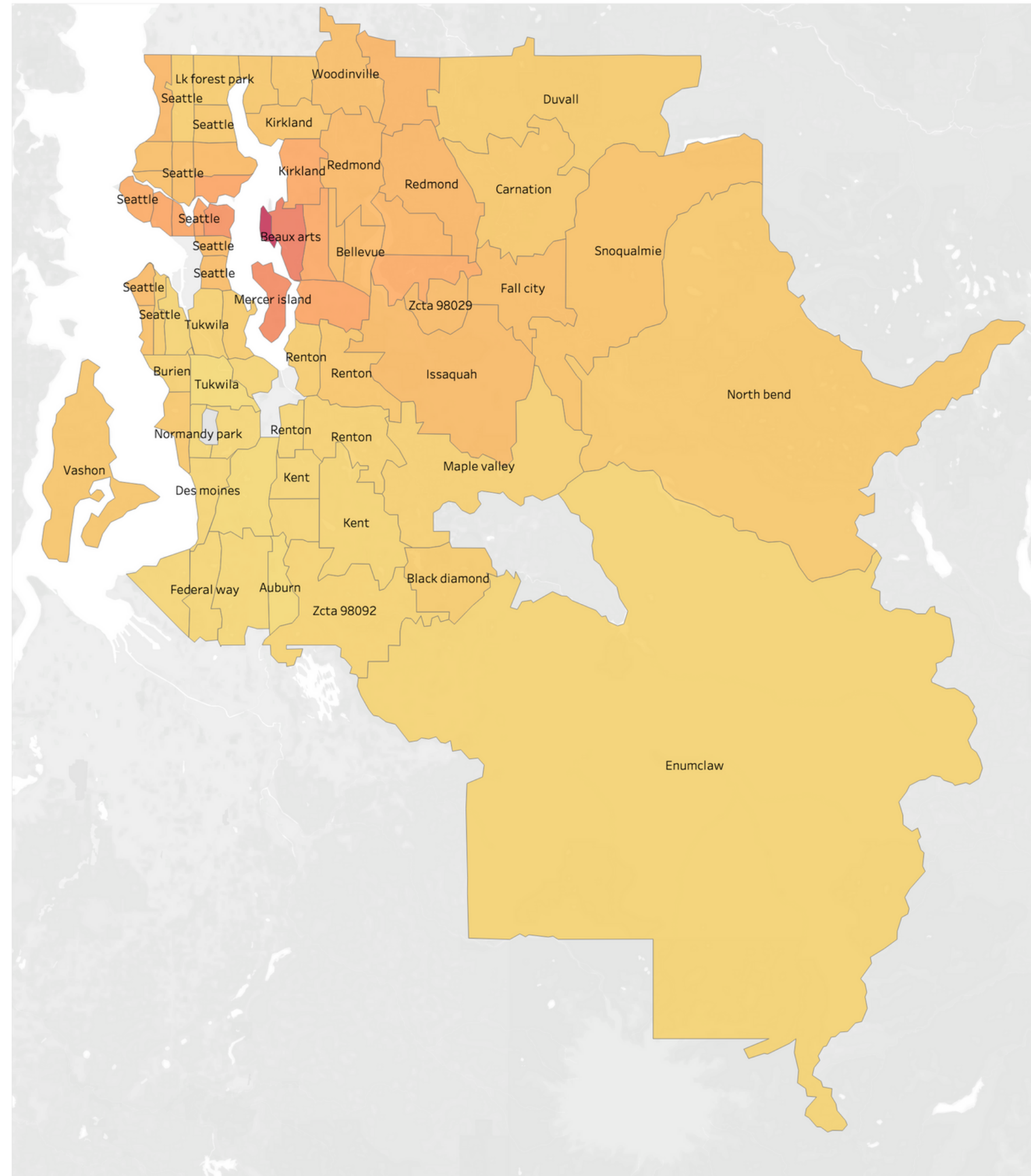- **No significant differences between the regions**

- **North part seems richer than the south Part**



Price vs. Location
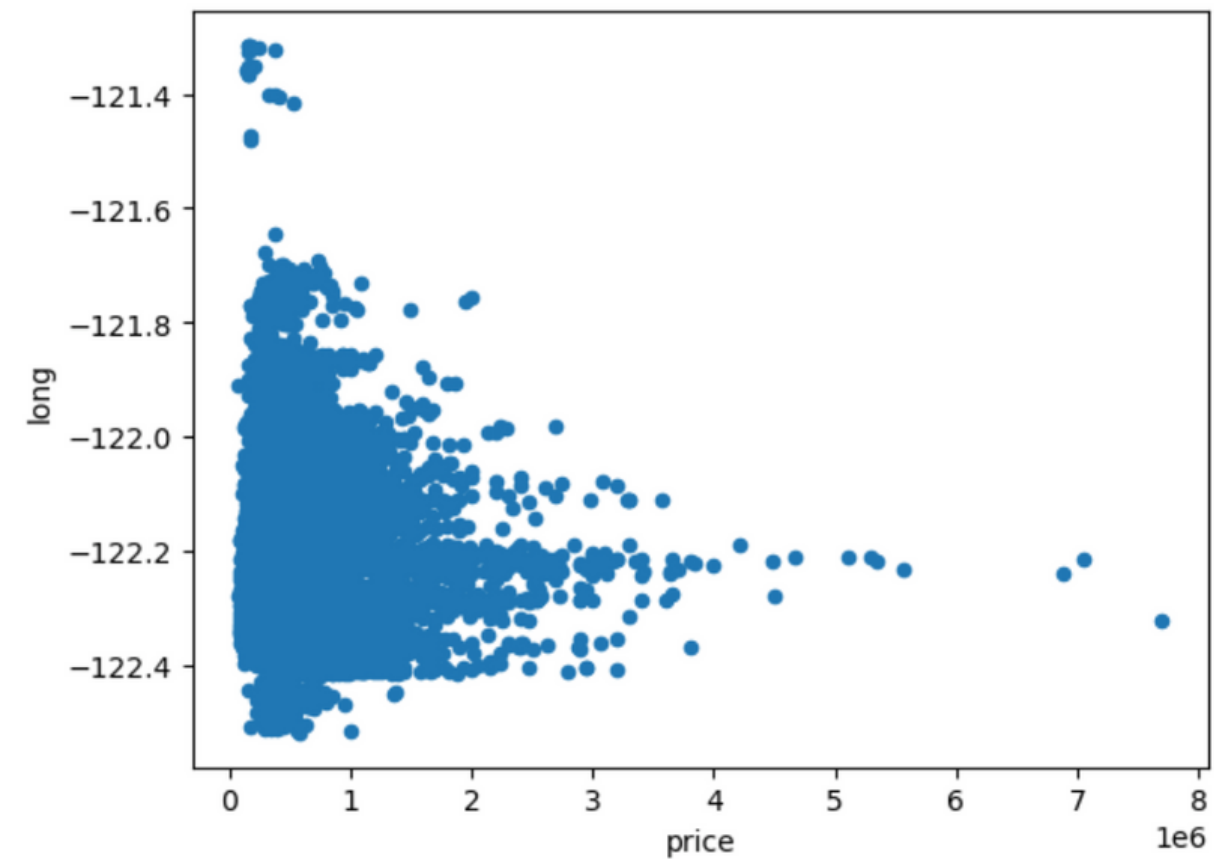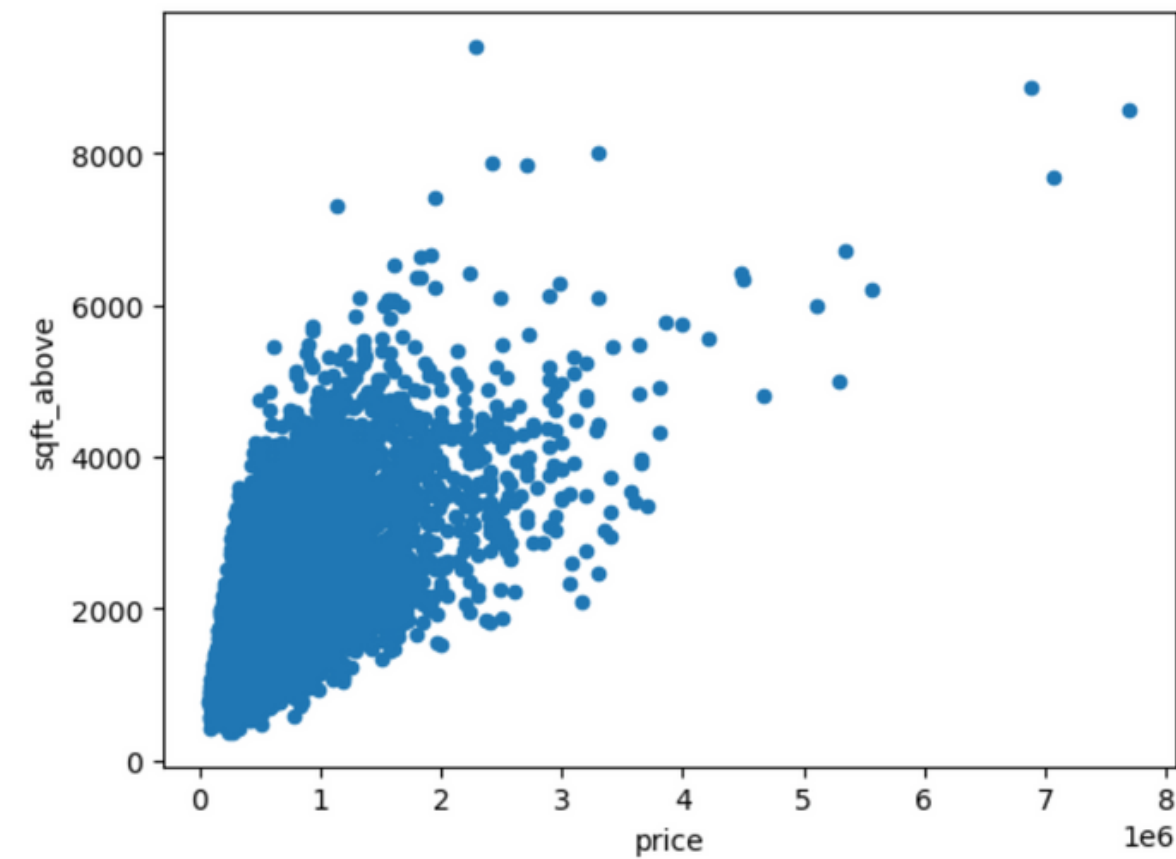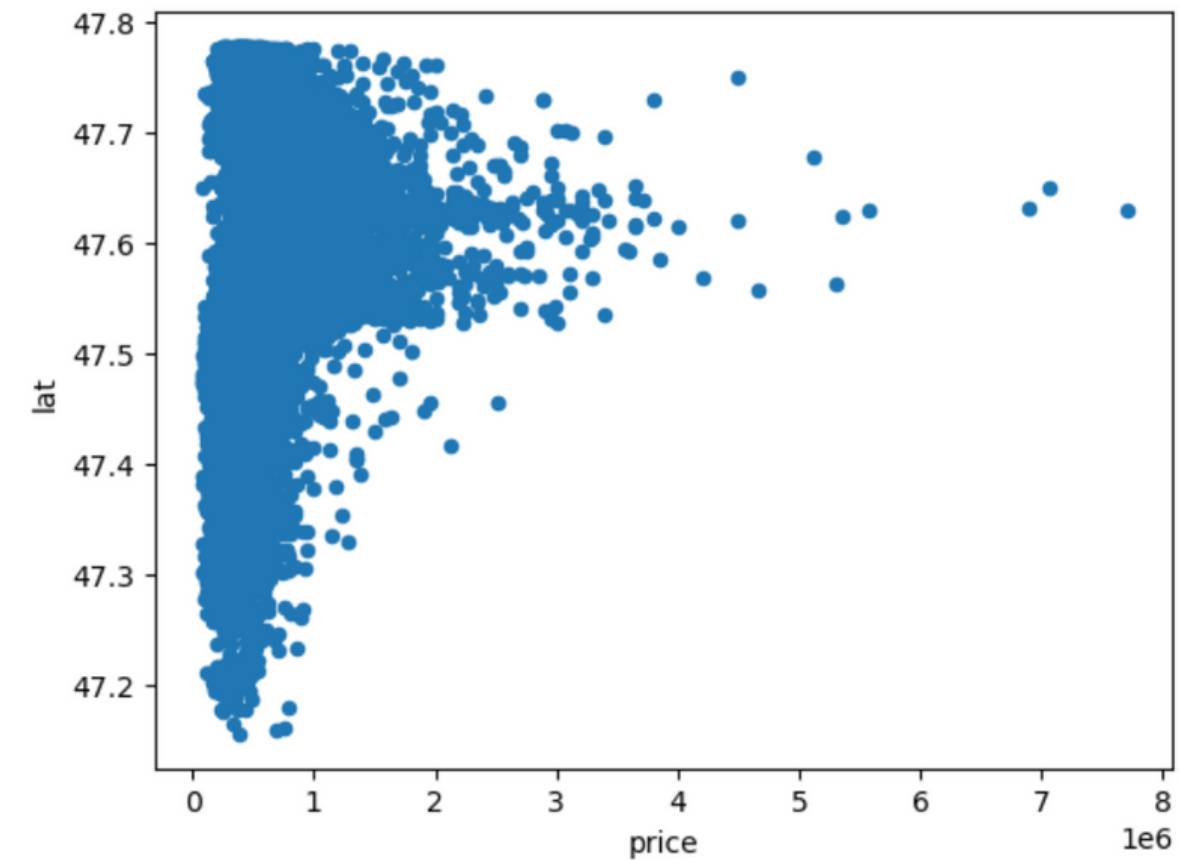
# Training the model with ?
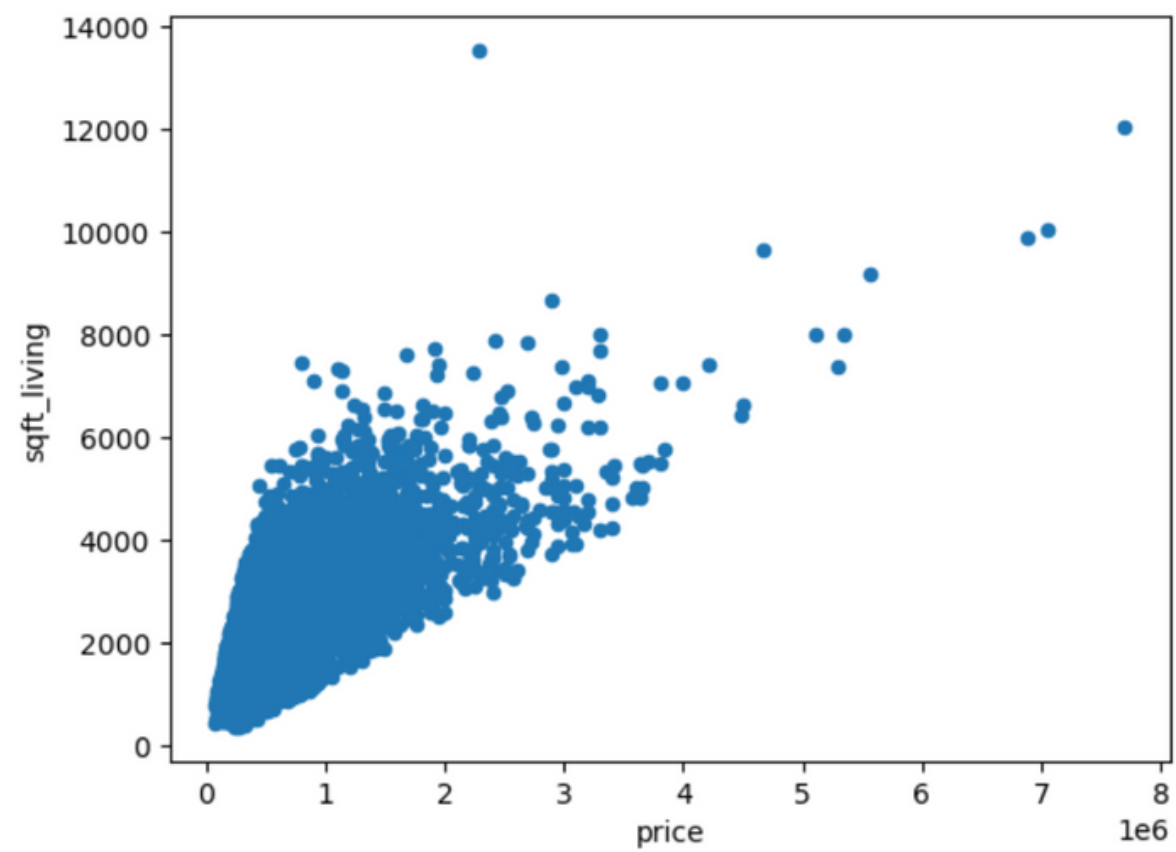
→ Exploratory Data Analysis
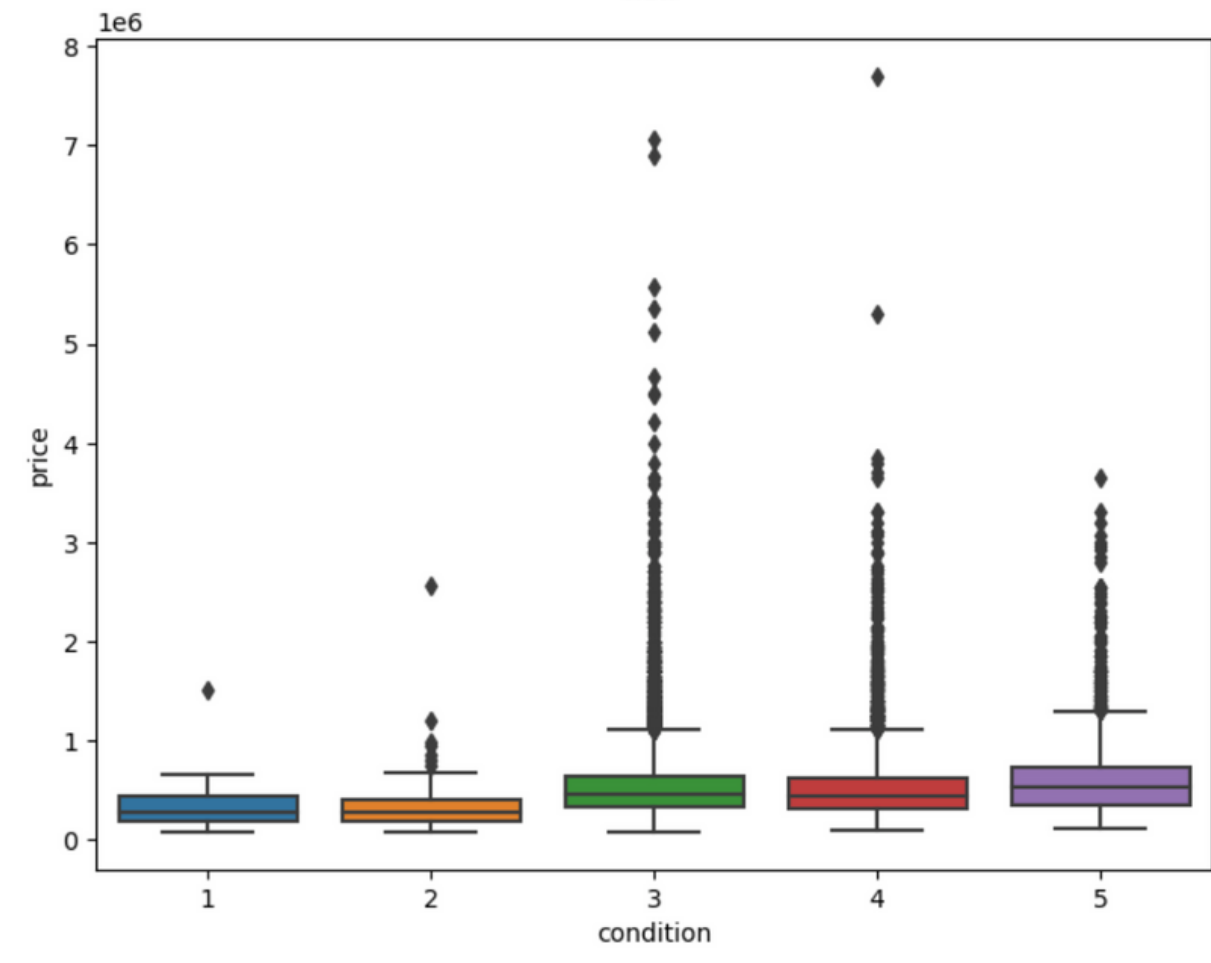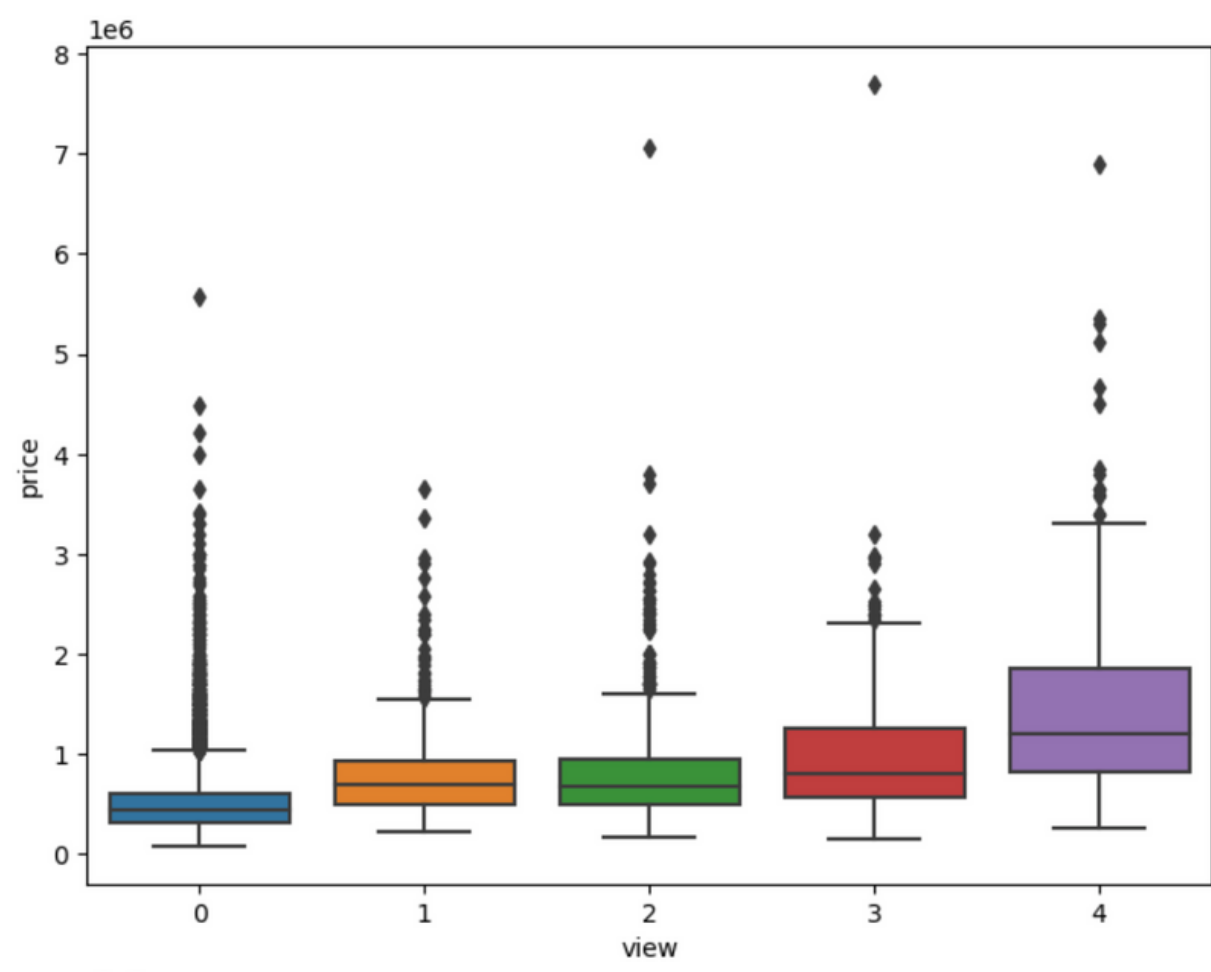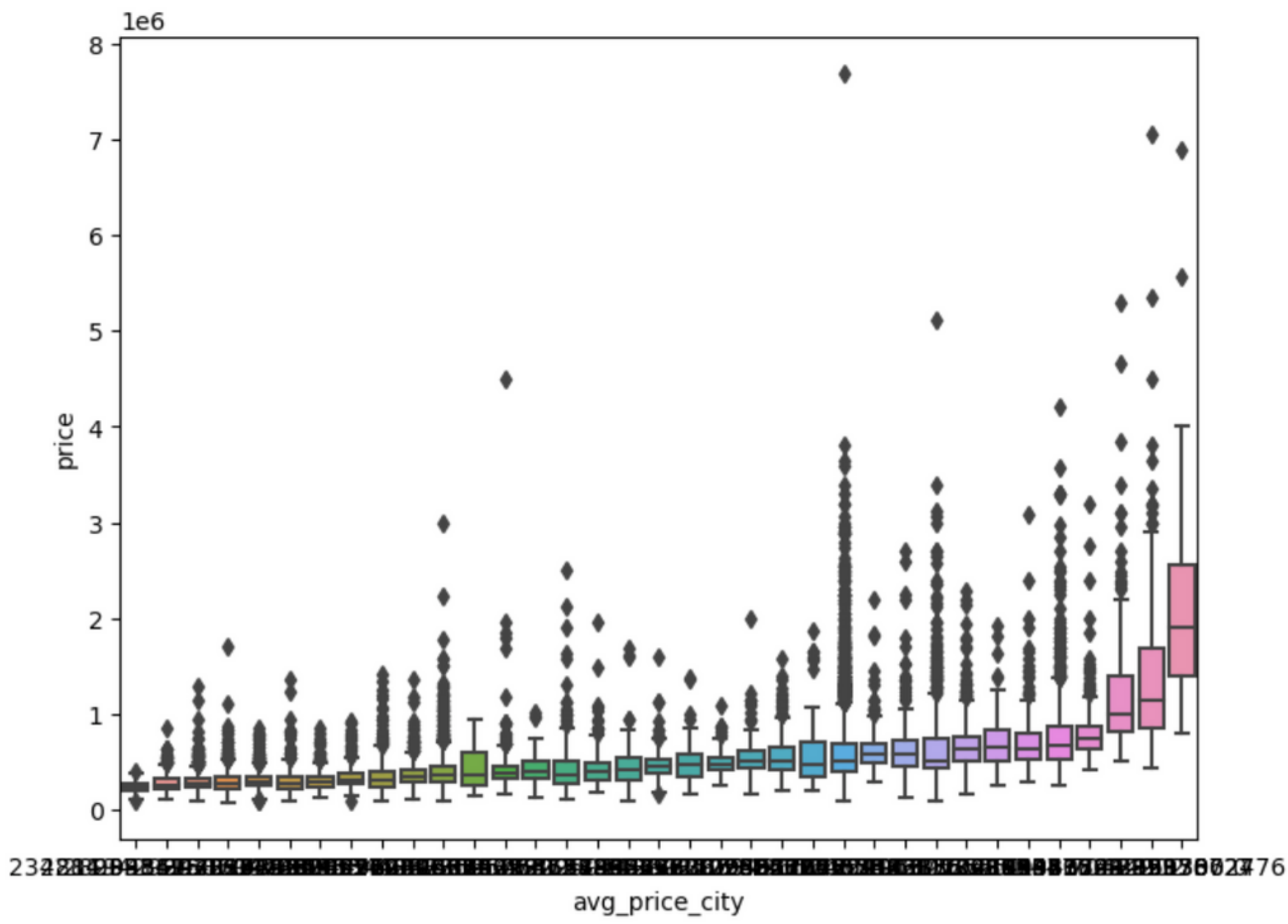
→ Introduce a new variable to the model: average_price_by_city

→ Linear Regression, KNN Model, Random Forest Regression

# Different results

## Linear Regression

```python
model = LinearRegression()
```

```python
model.fit(X_train_num, y_train)
```

LinearRegression()

```python
model.score(X_test_num, y_test)
```

0.7584504417397475

```python
predictions = model.predict(X_test_num)
```

```python
r2_score(y_test, predictions), mean_absolute_error(y_test, pr
```

(0.7584504417397475, 113489.91225463424, 182341.38546659282)

## KNN - K Nearest Neighbors

```python
from sklearn.neighbors import KNeighborsRegressor
```

```python
knn = KNeighborsRegressor(n_neighbors=10)
knn.fit(X_train_num, y_train)
```

KNeighborsRegressor(n_neighbors=10)

```python
knn.score(X_test_num, y_test)
```

0.7430534654723944

```python
predictions = knn.predict(X_test_num)
```

```python
r2_score(y_test, predictions), mean_absolute_error(y_test, pre
```

(0.7430534654723944, 97195.82953703705, 182916.98093456234)

## Random Forest Regression

```python
forest = RandomForestRegressor()
```

```python
forest.fit(X_train_num, y_train)
```

RandomForestRegressor()

```python
forest.score(X_test_num, y_test)
```

0.8776952968240401

```python
predictions = forest.predict(X_test_num)
```

```python
r2_score(y_test, predictions), mean_absolute_error(y_test, pre
```

(0.8776952968240401, 70319.40950684248, 126198.5082635587)

# Analyzing results

## In term of pure analytic, some numbers

### Our best prediction model result : Random forest regression

## In term of business

Due to our good result, 88% of chances that we predict the good price.
We are confident that we can use our model in order to predict the future prices.
Indeed, it is now much more easier to evaluate the selling price of our houses to align with the market and increase our profits now. It will also save us a lot of time so we can focus more on different regions.