



Real Estate Industry of Belgium

Analysis and Prediction

Md Imran
Data Scientist
07/03/2022

Mission & Objectives

Mission:

- ▶ Cleaning and doing a complete analysis and interpretation of the dataset created during the previous challenge.
- ▶ In order to create a machine learning model to predict prices on Belgium's real estate's sales.

Objectives:

- ▶ Using Pandas for data manipulation.
- ▶ Using Matplotlib and/or Seaborn for plotting.
- ▶ Finding and understanding correlations between dataset's variables.

Data Collecting

- ▶ ▶ A dataset of **50k+** real estate's observations, Collect from Kaggle. Which is published as Belgium real estate industry
- ▶ ▶ It has a lot of entries : more than 50k ! By having the maximum amount of data to discover interesting correlations, and have a meaningful Analyse.

Data Cleaning



Identifying the needs:

To proceed to the analysis, we needed a clean dataset containing at least:

- ▶ Prices, postal code and per sqft price

Removing the outliers (error, incorrect or absurd).

- ▶ It's good to have a lot of columns, as it can create more correlations between them. However, it's bad to have columns with errors, incorrect, missing or absurd values.

Data Cleaning



Two phases of data cleaning:

1. Cleaning the raw:

- ▶ A very first clean to the raw data. We were focused on "**dropping the big lies**":
- ▶ **Dropping** the duplicated rows
- ▶ **Dropping** columns with unique value
- ▶ **Checking** each columns' properties

2. Refining the values

- ▶ Some tweaks were made on the dataset to **remove outliers and useless columns**, due to their high rate of *None* value. This step required deeper investigation in top the data.

Data Cleaning

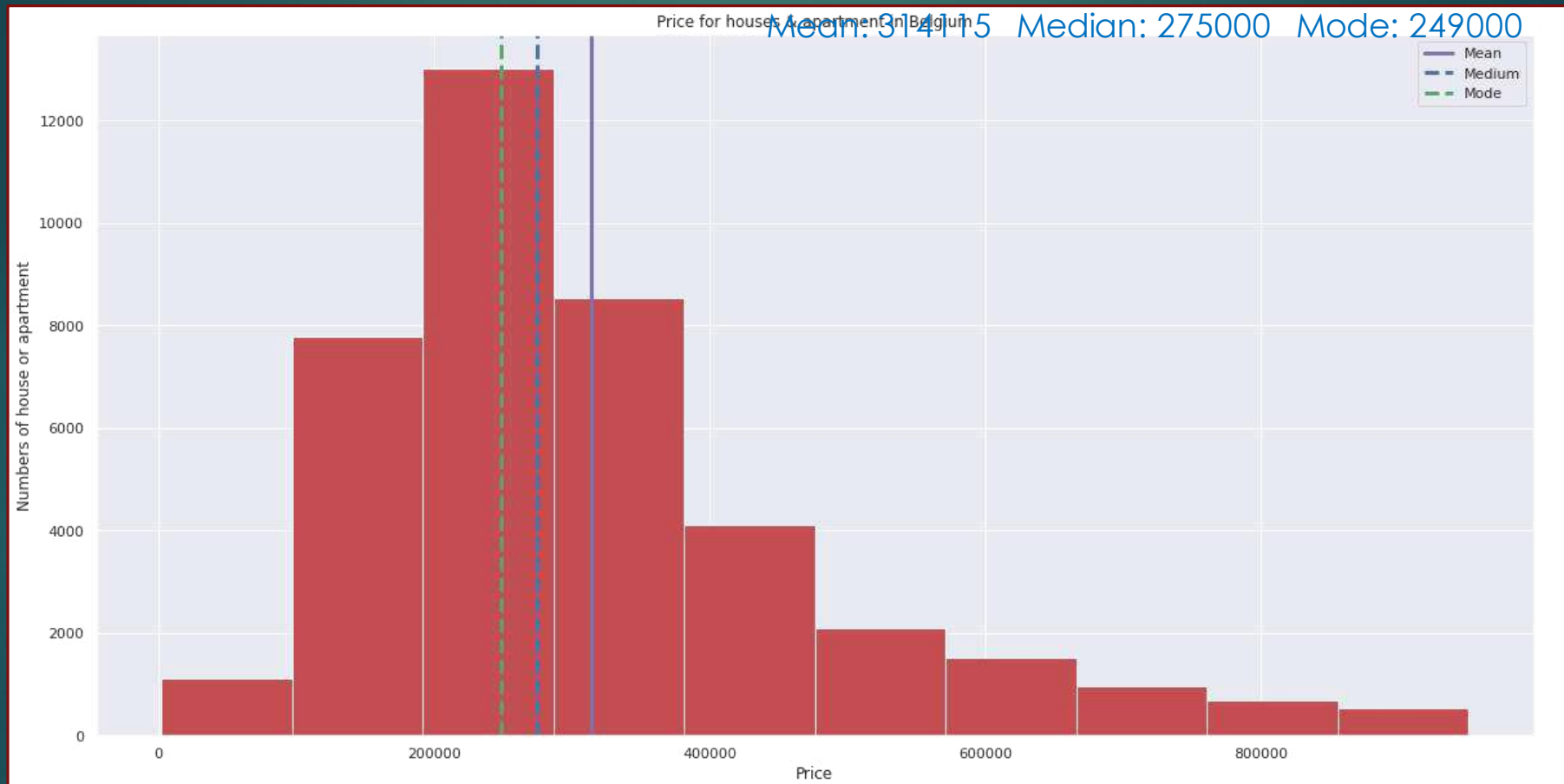
Details:

- ▶ Dropping "terrace_area" column
 - It has more than 30% of None.
 - ▶ Dropping "garden_area" column
 - It has more than 50% of None.
 - ▶ Dropping "subtype" column
 - Lots of property subtype. Some with less than 100 entries, in a dataset of 50.000.
 - This column was not relevant.
 - ▶ Removing the "Apartment blocks" entries
 - Apartment blocks are a whole building. It's not the kind of real estate sales we want here.
 - ▶ Changing None to "unknow"
- ▶ We also refactored all *float* to *int*. At the end of the cleaning, **we merged our dataframe with the two other ones created during the request study.**

40395 rows , 18 columns

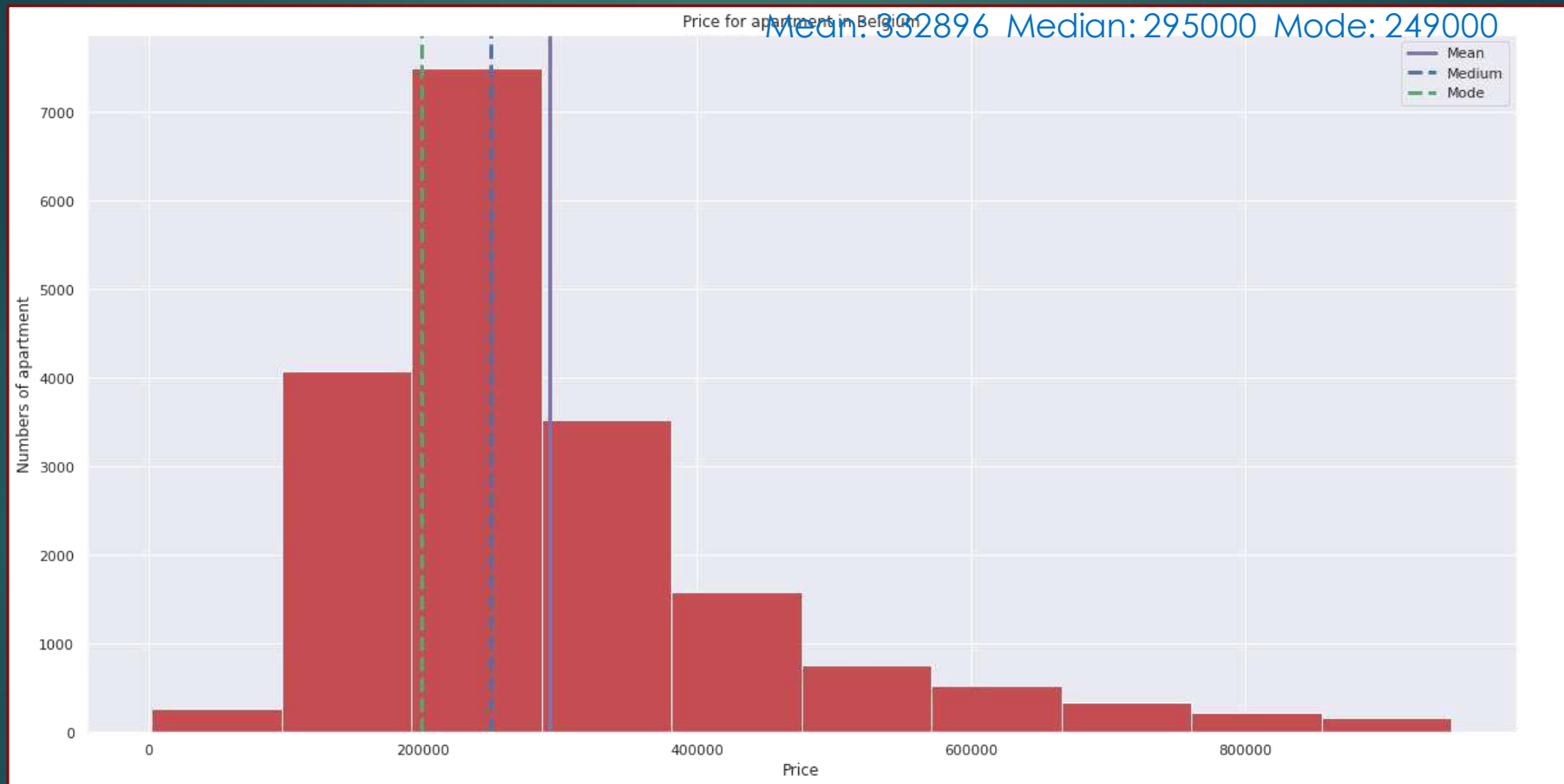
Data Visualisation

Our target: The Price



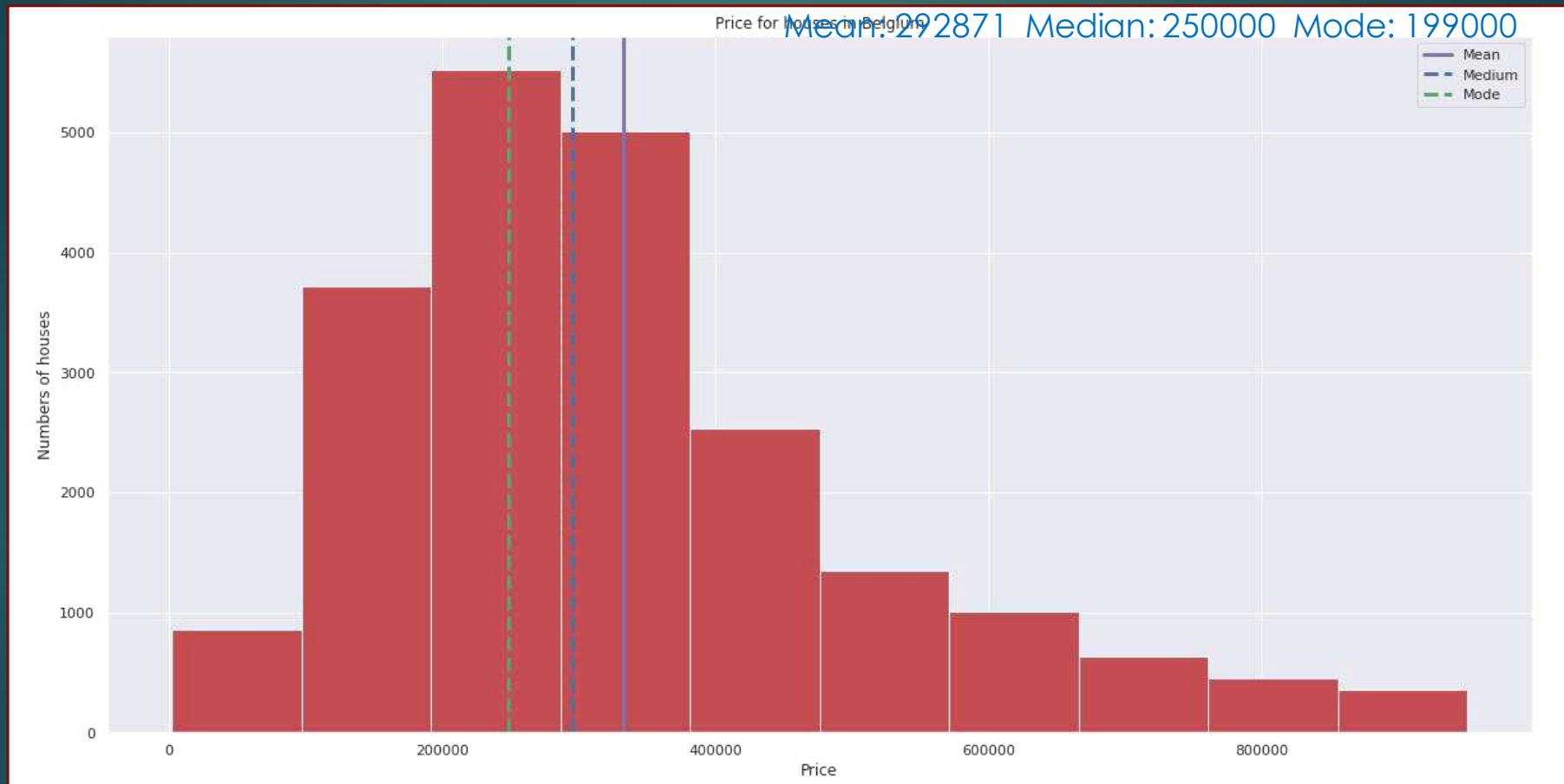
Data Visualisation

Our target: The Price

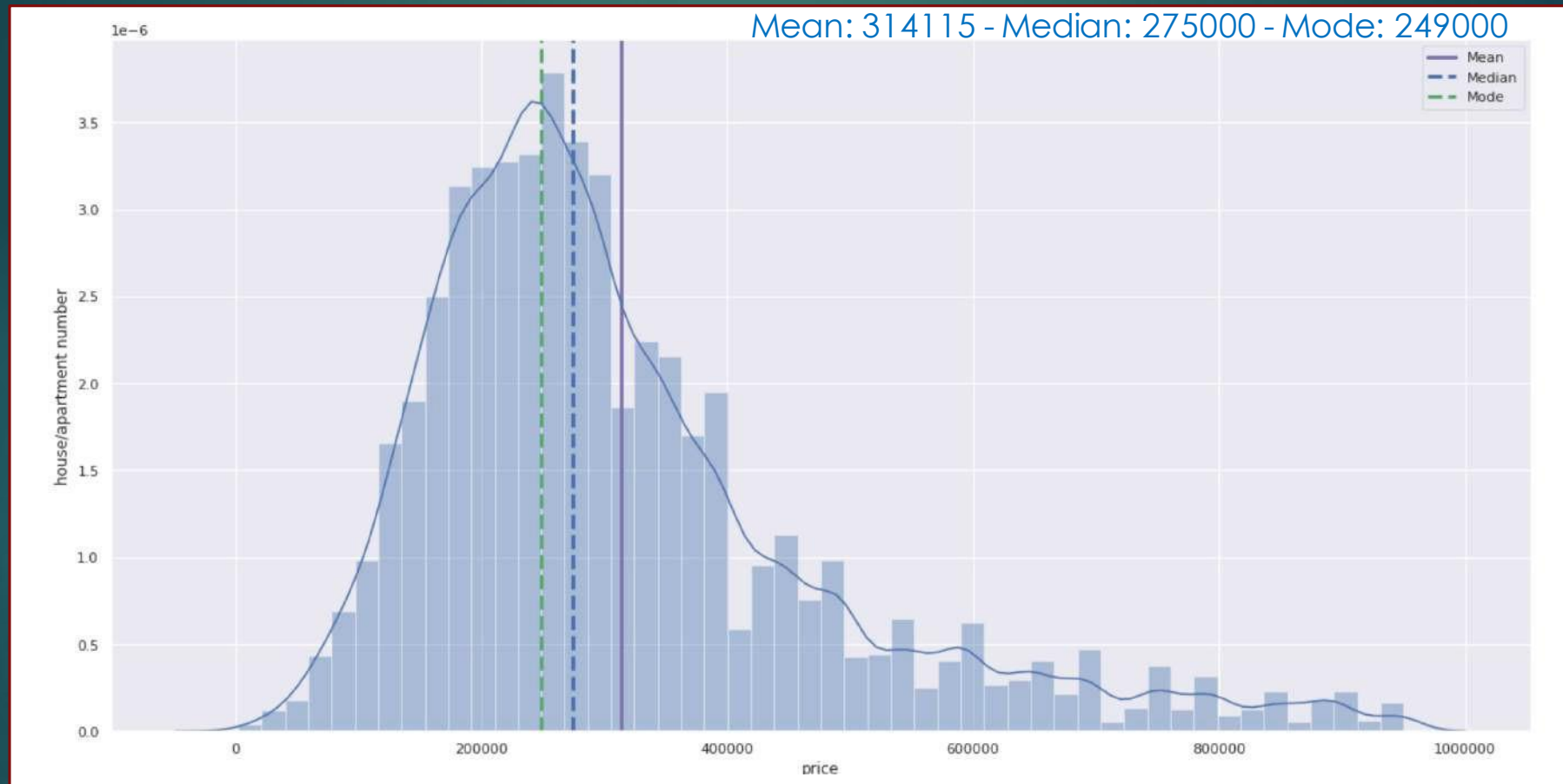


Data Visualisation

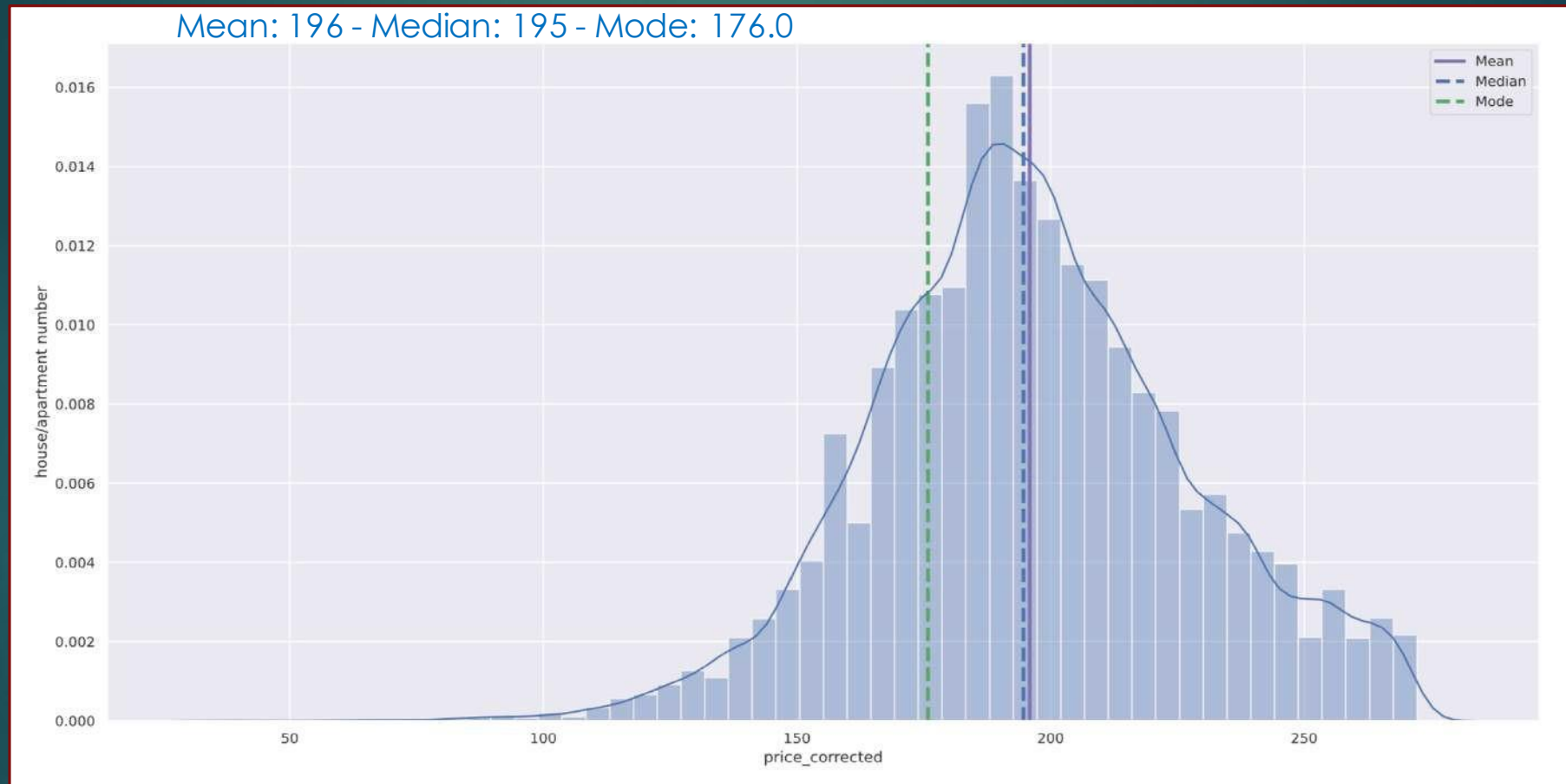
Our target: The Price

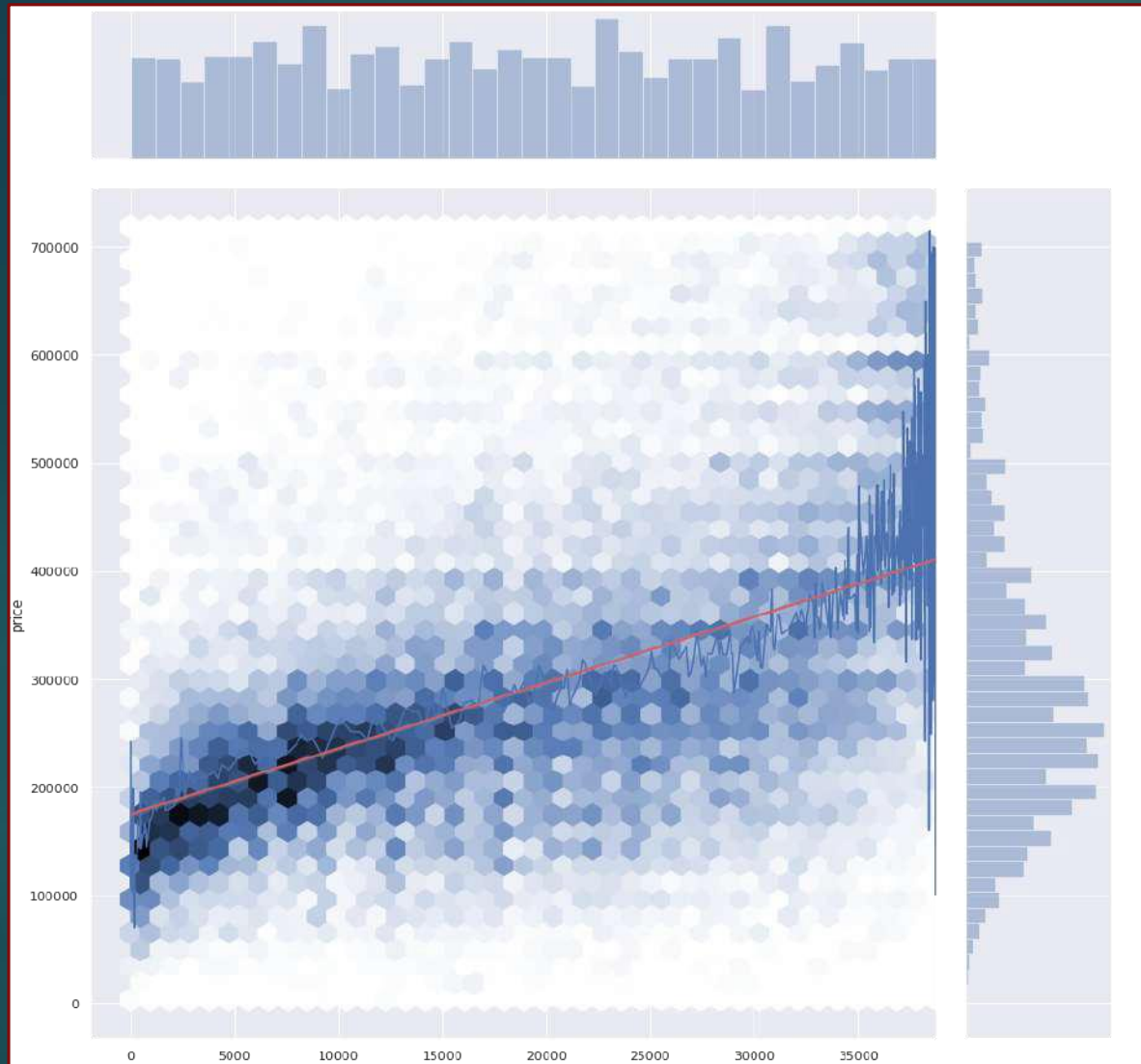


Data Visualisation



Data Visualisation

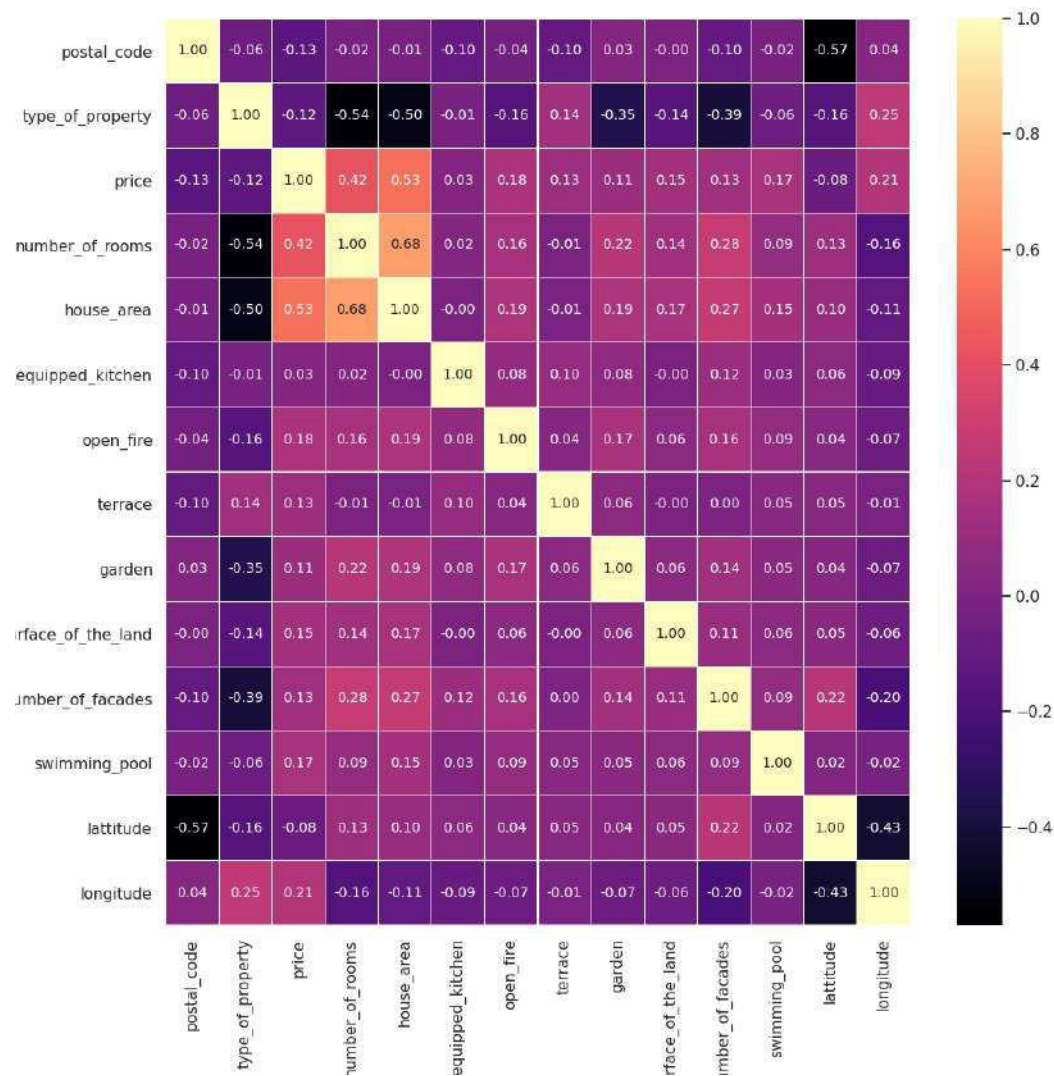




Data Visualisation

1.The Price of a house is correlated with its area: **The higher is the area, higher is the price.**

2.However, this correlation is not very strong, especially for big houses (houses with a area bigger than 35000 m2): The Price may vary a lot ! It may have other factor that influence the price of "big" houses.



Data Interpretation

Correlation Heatmap

Observations:

- 1.The **Price** is mainly correlated with the **Number of rooms** and the **House area**.
- 2.The **Number of rooms** and **House area** seems mainly correlated with each other.
- 1.The **Type of property** is the variables which has the most correlation with other variables.

Correlation does not imply causation

Links



Github Repository

- ▶ <https://github.com/mdimran1/Real-Estate>

Linkedin :

- ▶ <https://www.linkedin.com/in/imran-pro/>