# PREDICTING HOUSE PRICES WITH LINEAR REGRESSION

by Ruba Alnashwan

# METHODOLOGY

- Proplem understanding
- Data collection
- Data cleaning
- Exploratory data analysis (EDA)
- Feature Engineering and selection
- Data modeling

# Proplem understanding

Overview:
In this project, we will use data for villas offered for sale in the city of Riyadh from the (Aqar) website, which provides villas for sale and their prices.

Problem statement:
The objective of this project is to predict the price of the villa by providing some features for each vila by using supervised linear regression model

Scope:
The scope of this project was the villas for sale in the areas of Riyadh, whether old or new villas. The dataset contains 2099 rows × 10 columns.
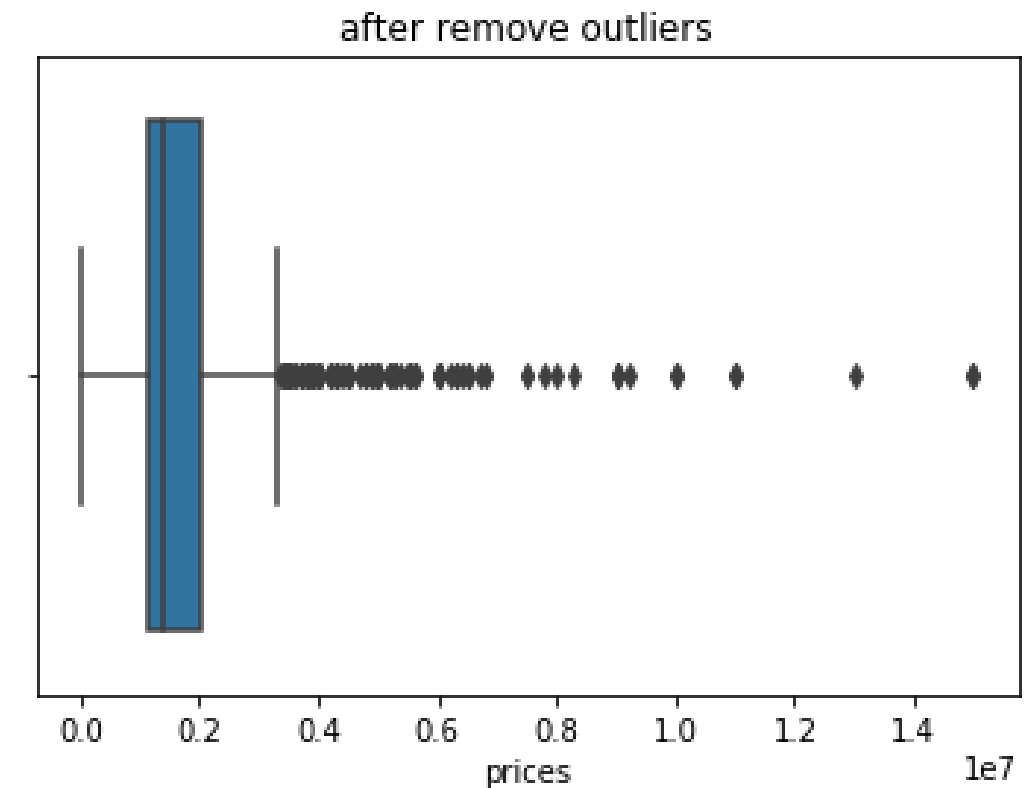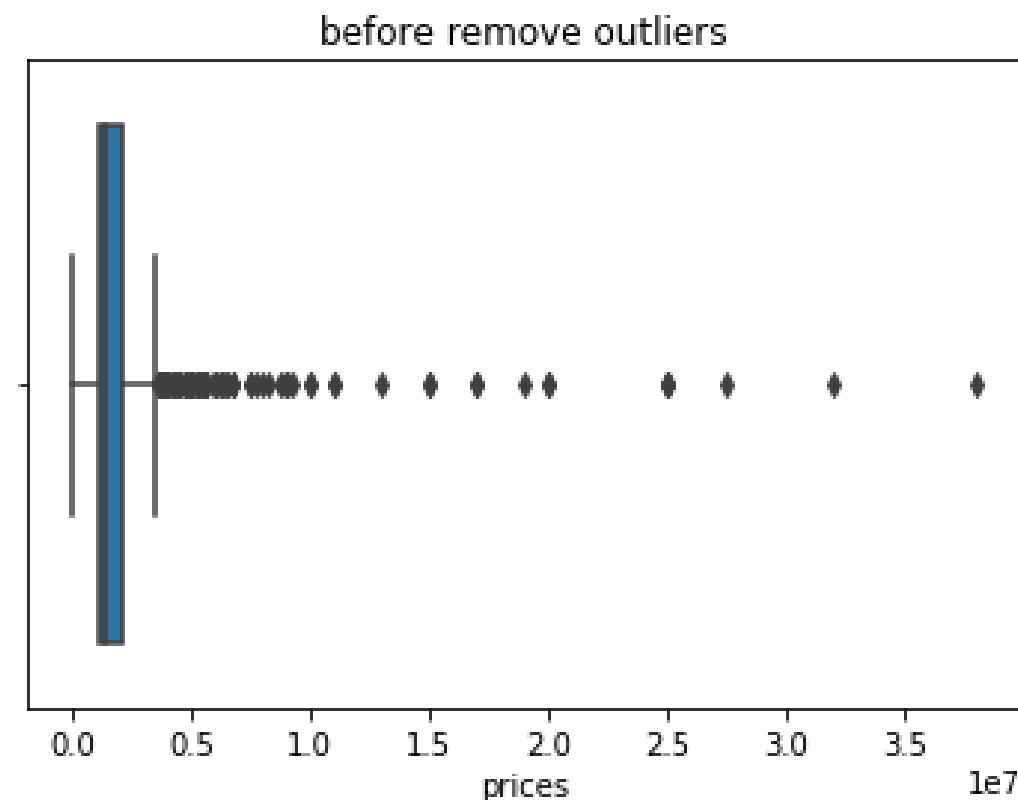
# Data cleaning

check nulls

I filled in the null values in the oldness feature with zero. My notes. If there is no oldness, the villa will be new and filled the street width, number living rooms with mode and median.
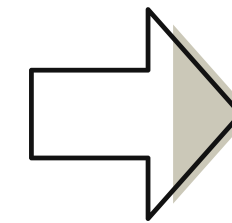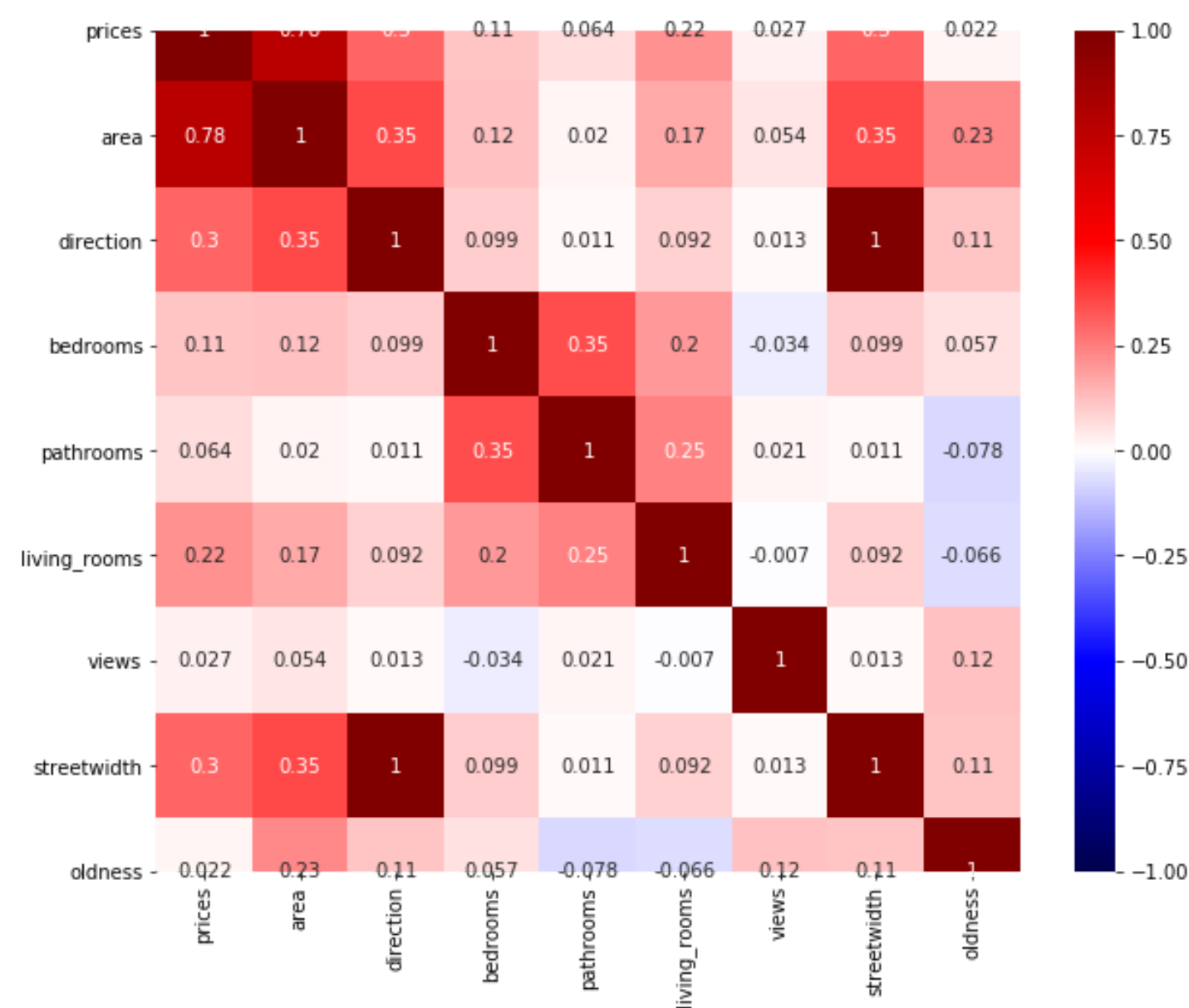
check duplicates

# Data cleaning

▸ replace some character and arabic words like (شمال)
(م٢، سنة) and ( الرياض ، غرب الرياض)

▸ check outliers using EDA and remove



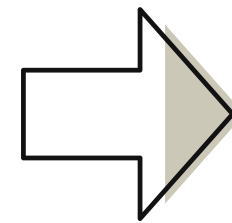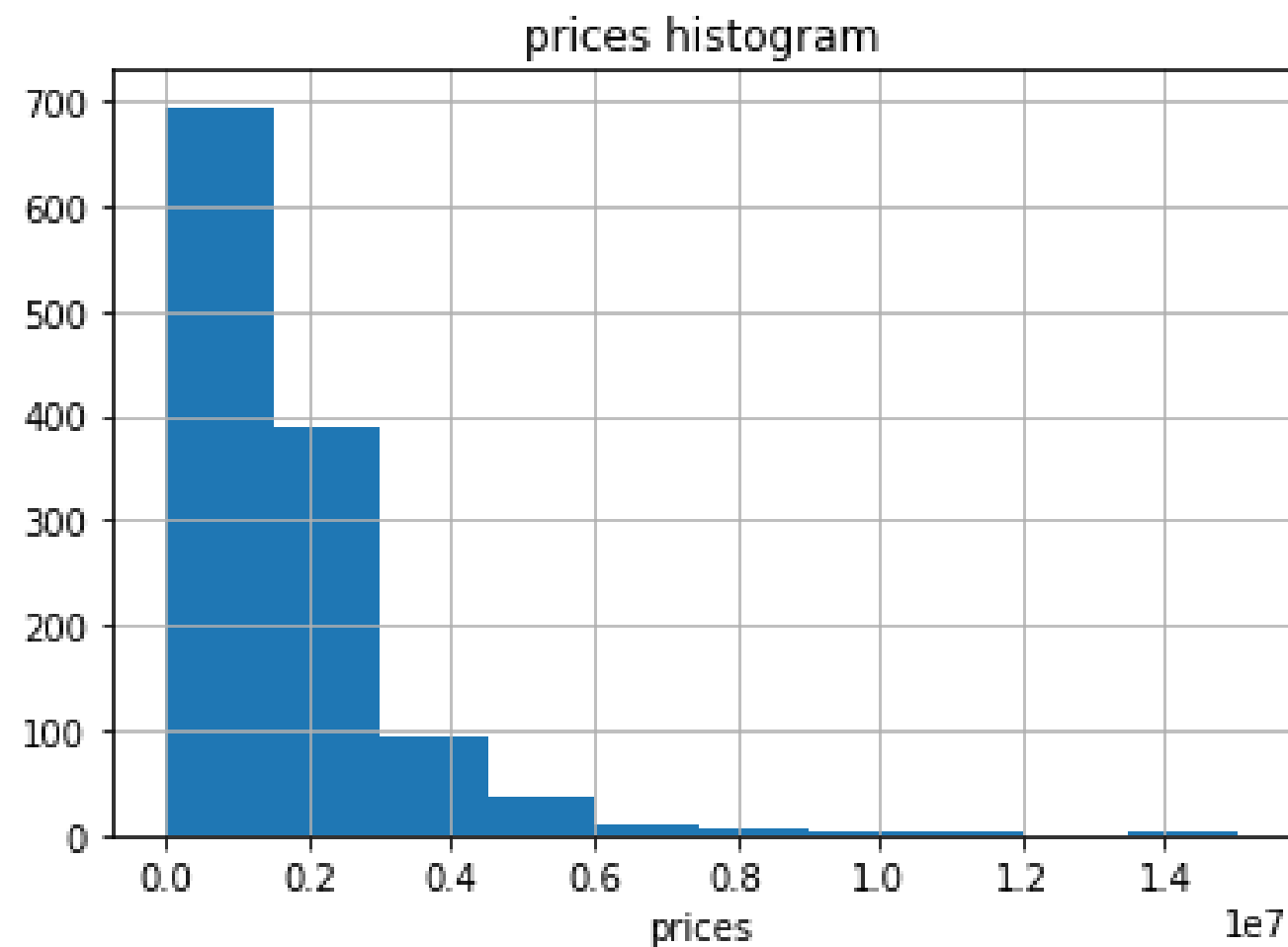before remove outliers

after remove outliers

# Exploratory data analysis (EDA)



Observations:

This heat map shows that there is no strong relationship between the dependent variable and features except for the area.

# Exploratory data analysis (EDA)



prices histogram
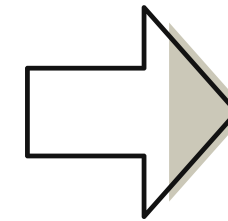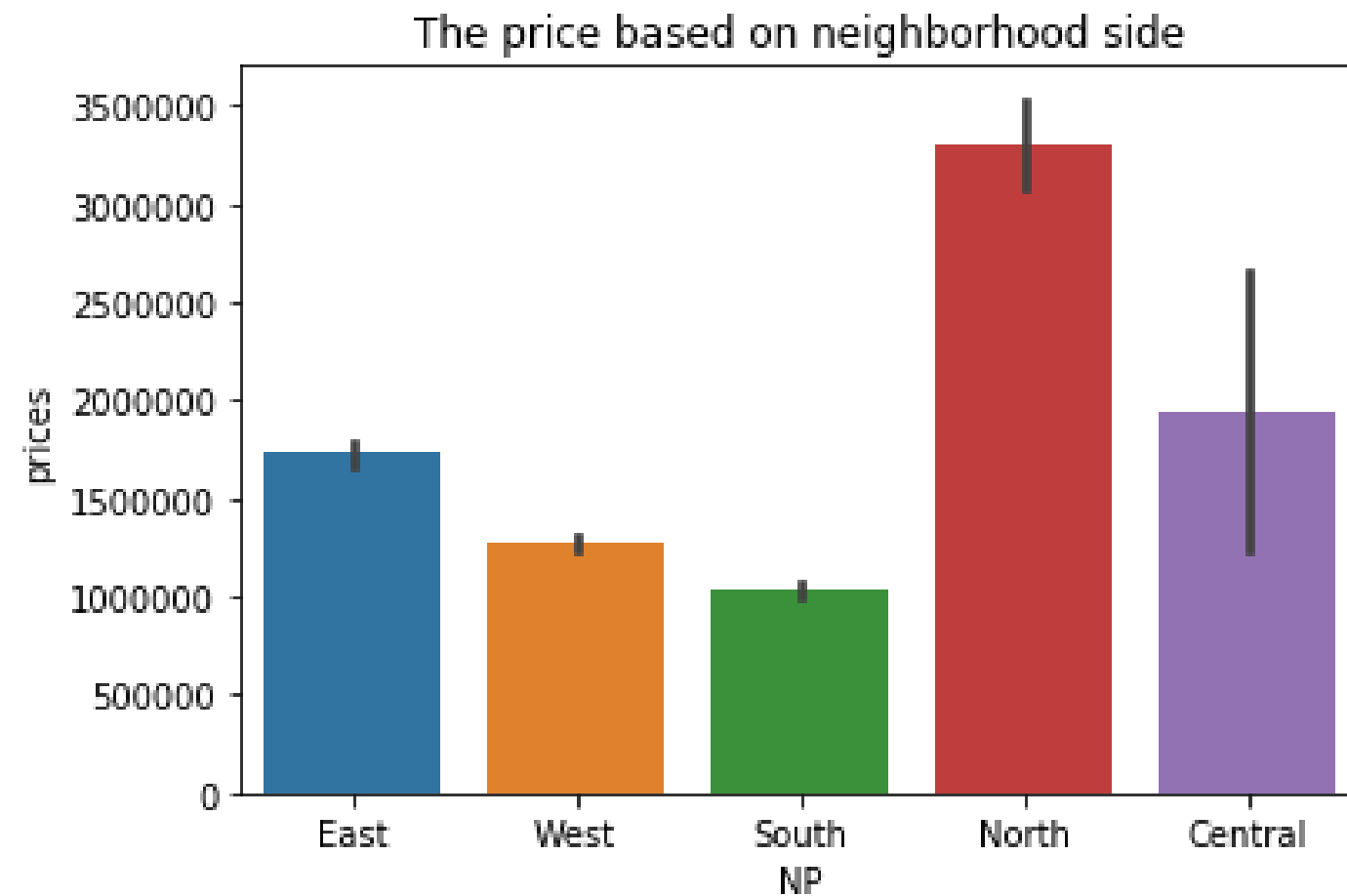
Observations:

In this histogram shows the most of the density lies between 1M and 2M, but there appears to be a lot of outliers on the pricier side.

# Exploratory data analysis (EDA)



The price based on neighborhood side

Observations:

The par plot shows the highest prices are located in the north

# Feature Engineering and selection

Principal component analysis (PCA) to reduce the data

Polynomial Features But the results were not good

Dummies for Neighborhood Feature ,the model has been greatly improved

**Scaling:** MinMaxScaler, StandardScaler

**Feature selection:** lasso, ridge

# Data modeling

MAIN METRIC USED FOR EVALUATION : R2, RMSE

MODEL USED: LINEAR REGRESSION

## basline

TRAINING SCORE:
 0.37746

VALIDATION SCORE:
0.5955

## one feature power of 2

TRAINING SCORE:
 0.37963

VALIDATION SCORE:
0.59218

## dummies

TRAINING SCORE:
 0.60967

VALIDATION SCORE:
0.56537

## log target

TRAINING SCORE:
 0.60571

VALIDATION SCORE:
0.59849

## Polynomial Features

TRAINING SCORE:
 0.56041

VALIDATION SCORE:
0.63017

## PCA

TRAINING SCORE:
 0.6097

VALIDATION SCORE:
-0.10154

# Data modeling

MAIN METRIC USED FOR EVALUATION : R2, RMSE

MODEL USED: LINEAR REGRESSION

## basline

TRAINING SCORE:
0.37746

VALIDATION SCORE:
0.5955

## one feature power of 2

TRAINING SCORE:
0.37963

VALIDATION SCORE:
0.59218

## dummies

TRAINING SCORE:
0.60967

VALIDATION SCORE:
0.56537

## log target

TRAINING SCORE:
0.60571

VALIDATION SCORE:
0.59849

## Polynomial Features

TRAINING SCORE:
0.56041

VALIDATION SCORE:
0.63017

## PCA

TRAINING SCORE:
0.6097

VALIDATION SCORE:
-0.10154

# CONCLUSION

**result the best model :**

Training score:  0.60555

Validation score:  0.59771

**Difficulties:**

Real estate is always changing

Prices in each neighborhood are different

# Thank you

Any questions?