

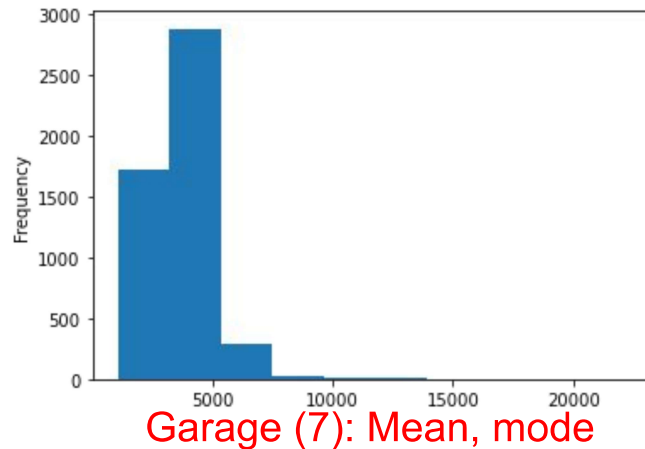
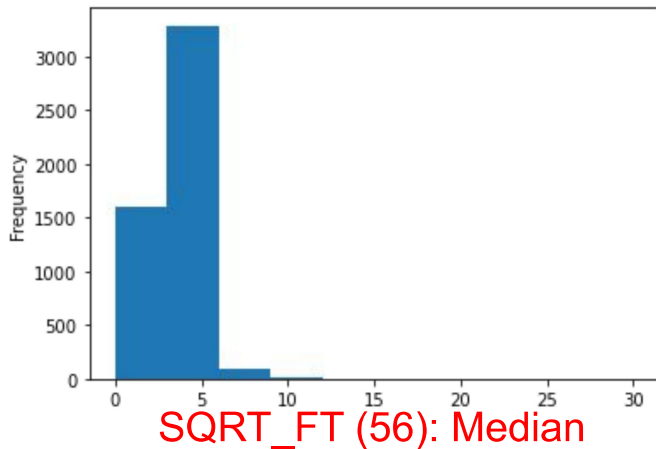
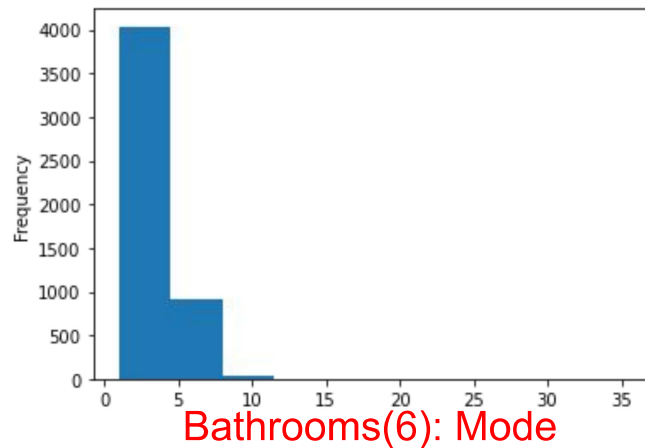
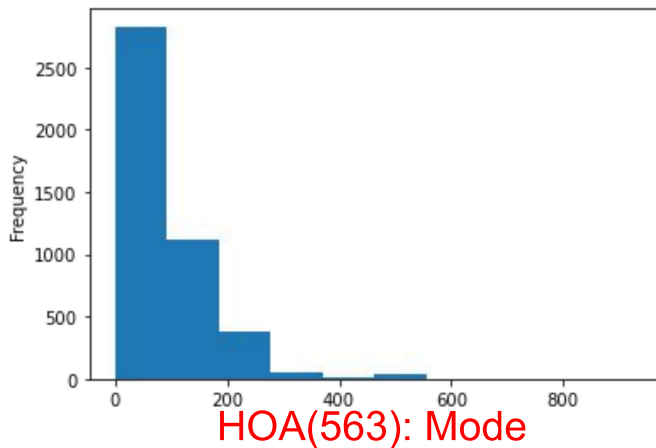
# Data Cleaning and EDA

David Lizama

# Introduction

- ❑ Database of 5000 observations and 16 variables related to houses.
- ❑ The first assumption is that our main objective is to predict house prices.
- ❑ Three main task: fill in missing values, separate features from strings, and handling bad records.
- ❑ No data normalization or correction of outliers in this assignment.
- ❑ The most correlated variables with house prices are lot\_acres, bathrooms, sqrt\_ft, and fireplaces.
- ❑ It is recommended to drop MLS and zipcode from the database.

# Cleaning data: Fill in missing values



There are 563 observations with 'None' and 35 missing values. It is not possible to drop them (>5%), so we need to impute them.

# Cleaning data: Bad records

## ❑ Numbers separated by commas.

- ❑ `df['HOA'] = df.HOA.str.replace(',', '.')`
- ❑ `df['HOA'] = df.HOA.str.replace('None', '0')`
- ❑ `df['HOA'] = df.HOA.astype(float)`

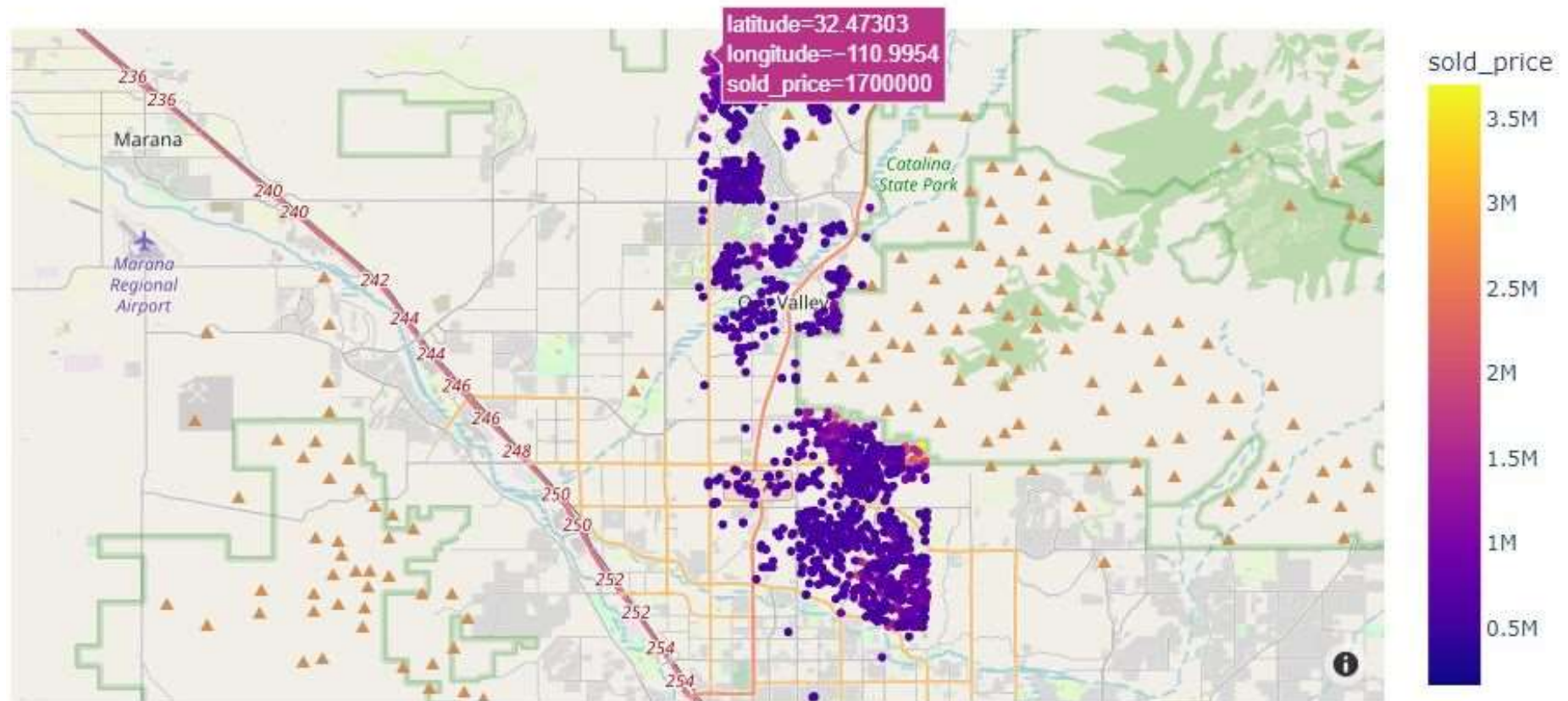
```
File ~\anaconda3\lib\site-packages\pandas\core\dtypes\c
1177     raise ValueError(msg)
1179 if copy or is_object_dtype(arr.dtype) or is_obj
1180     # Explicit copy, or required since NumPy ca
-> 1181     return arr.astype(dtype, copy=True)
1183 return arr.astype(dtype, copy=copy)
```

**ValueError:** could not convert string to float: '1,717'

## ❑ Latitude and longitude with numbers separated by dots.

- ❑ Remove all the dots from columns.
- ❑ Divided by  $10^6$  all observations.
- ❑ Correct some observations which were out of range.
- ❑ All houses should be located in Houston, Texas. Some of them were located in the middle of the ocean.

## Cleaning data: Bad records



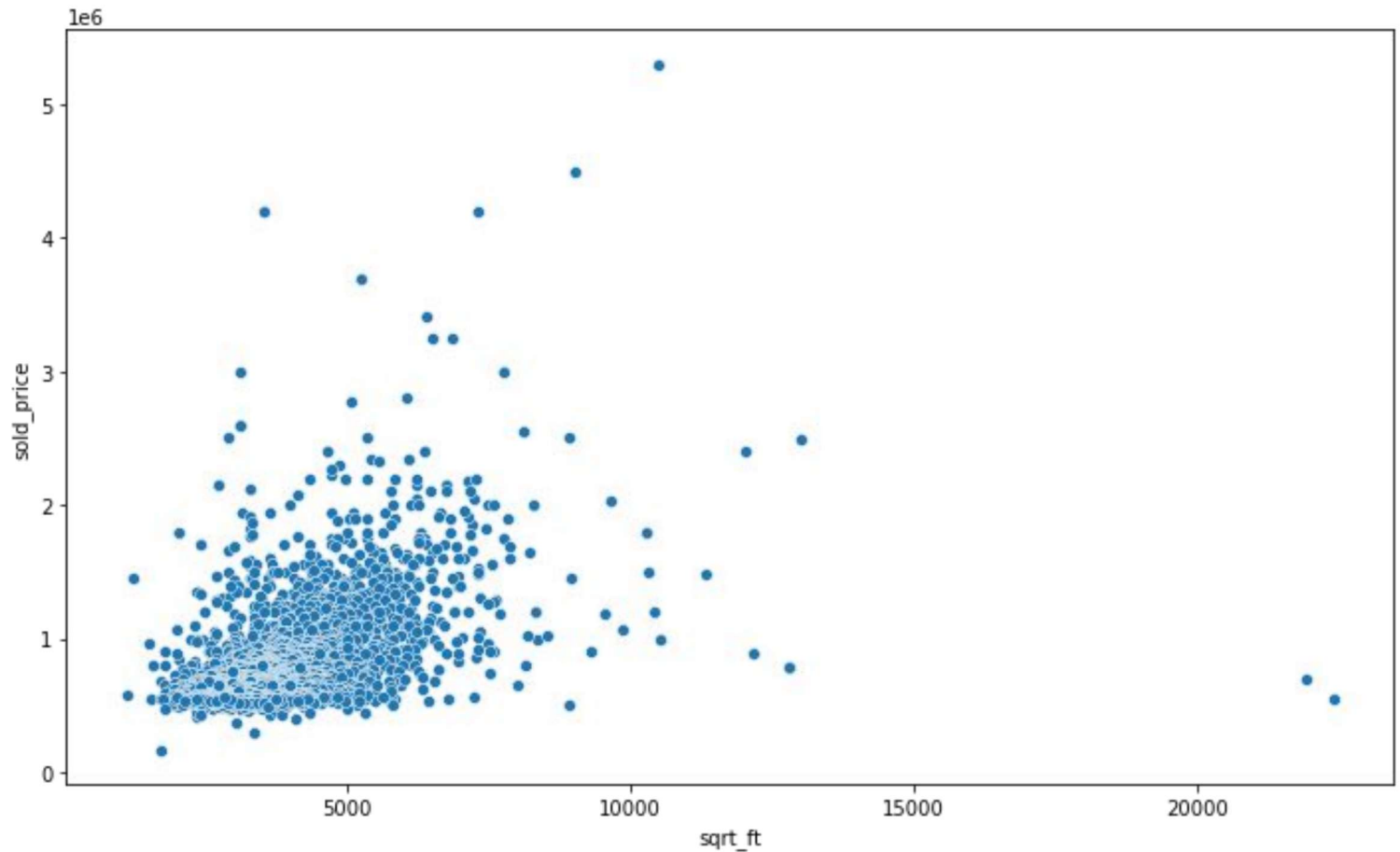
# Cleaning data: String variables

- ❑ There are two columns with strings: floor covering and kitchen features.
- ❑ In this case, I detected some keywords in this columns and I created new columns with zeros or ones to count if houses have these keywords.

Floor covering	Kitchen features
<ul style="list-style-type: none"><li>• FC_Stone 1499</li><li>• FC_Ceramic 2527</li><li>• FC_Laminate 86</li><li>• FC_Wood 1248</li><li>• FC_Carpet 3509</li><li>• FC_Concrete 756</li><li>• FC_MexicanTile 660</li></ul>	<ul style="list-style-type: none"><li>• KF_Dishwasher 4857</li><li>• KF_GarbageDisposal 4520</li><li>• KF_Refrigerator 4234</li><li>• KF_DoubleSink 1164</li><li>• KF_Microwave 3625</li><li>• KF_Oven 3977</li><li>• KF_Compactor 432</li><li>• KF_Freezer 395</li><li>• KF_ElectricRange 401</li><li>• KF_Island 1252</li><li>• KF_GasRange 1307</li><li>• KF_Countertops 1482</li><li>• KF_Desk 327</li></ul>

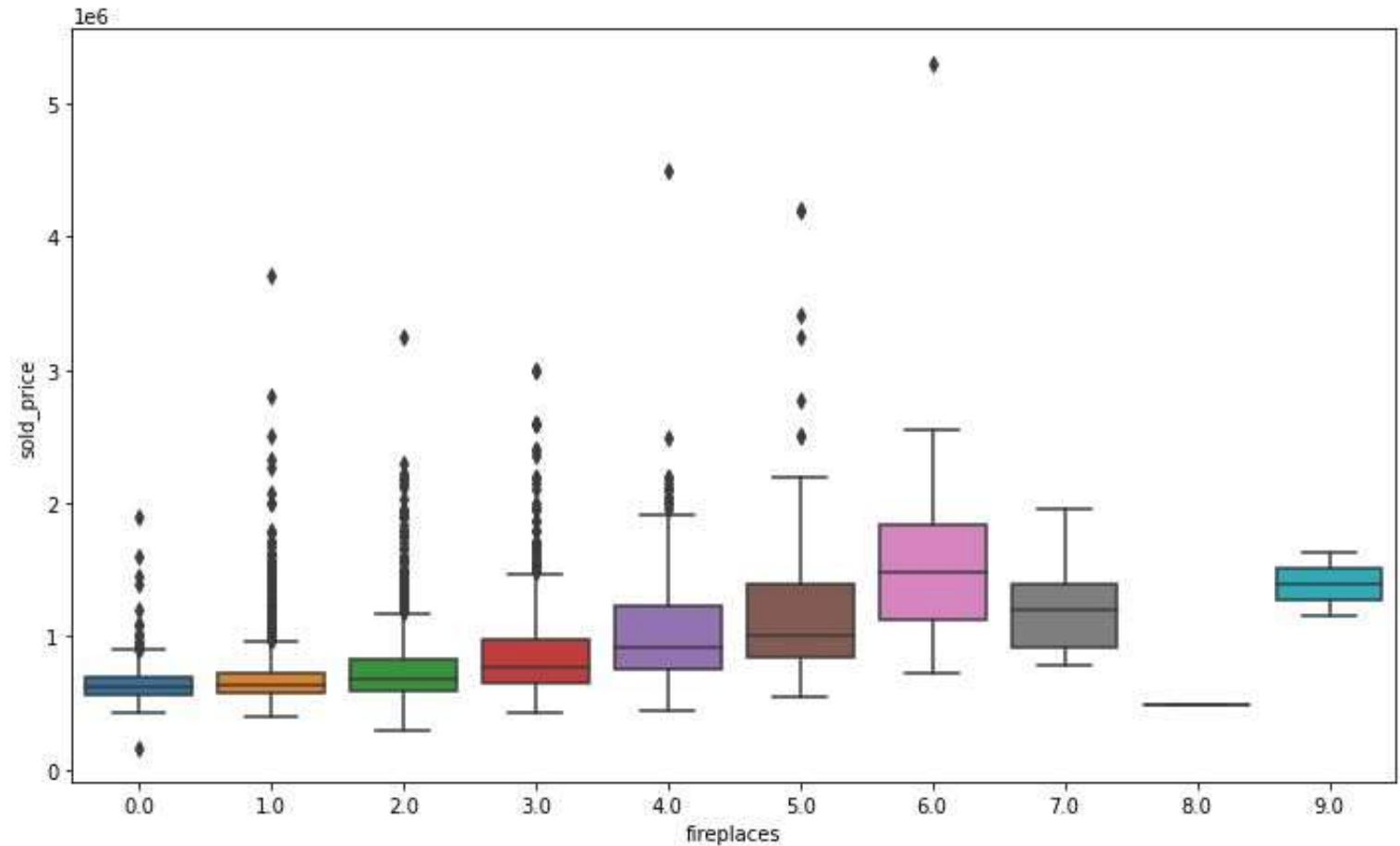
# EDA: Price houses and square footage

- High correlation between the house prices and sqrt\_ft (0.52).
- When sqrt\_ft increases, the price of the house increases.



## EDA: Price houses and fireplaces

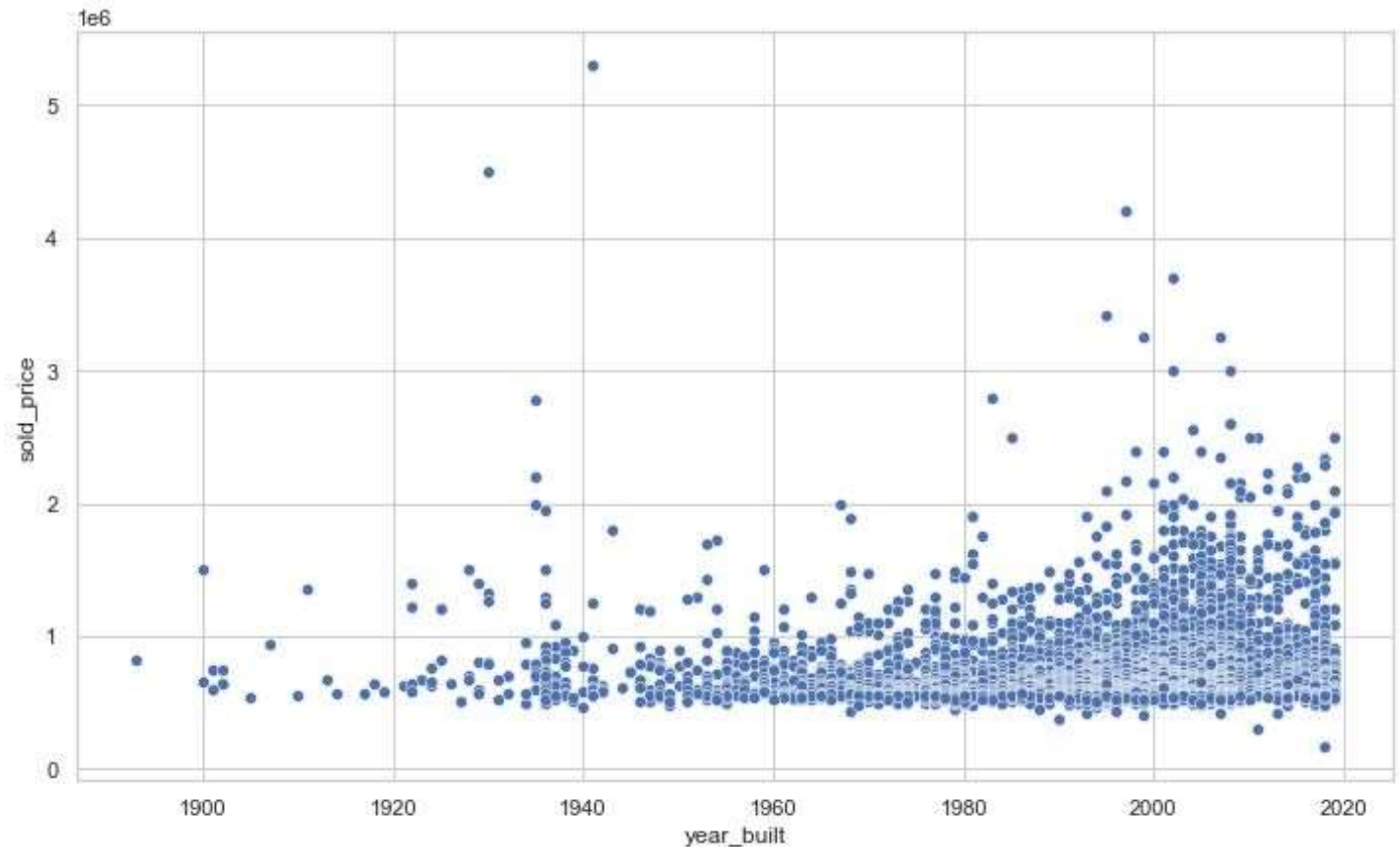
- There is a positive association between the amount of fireplaces and the price of the house.





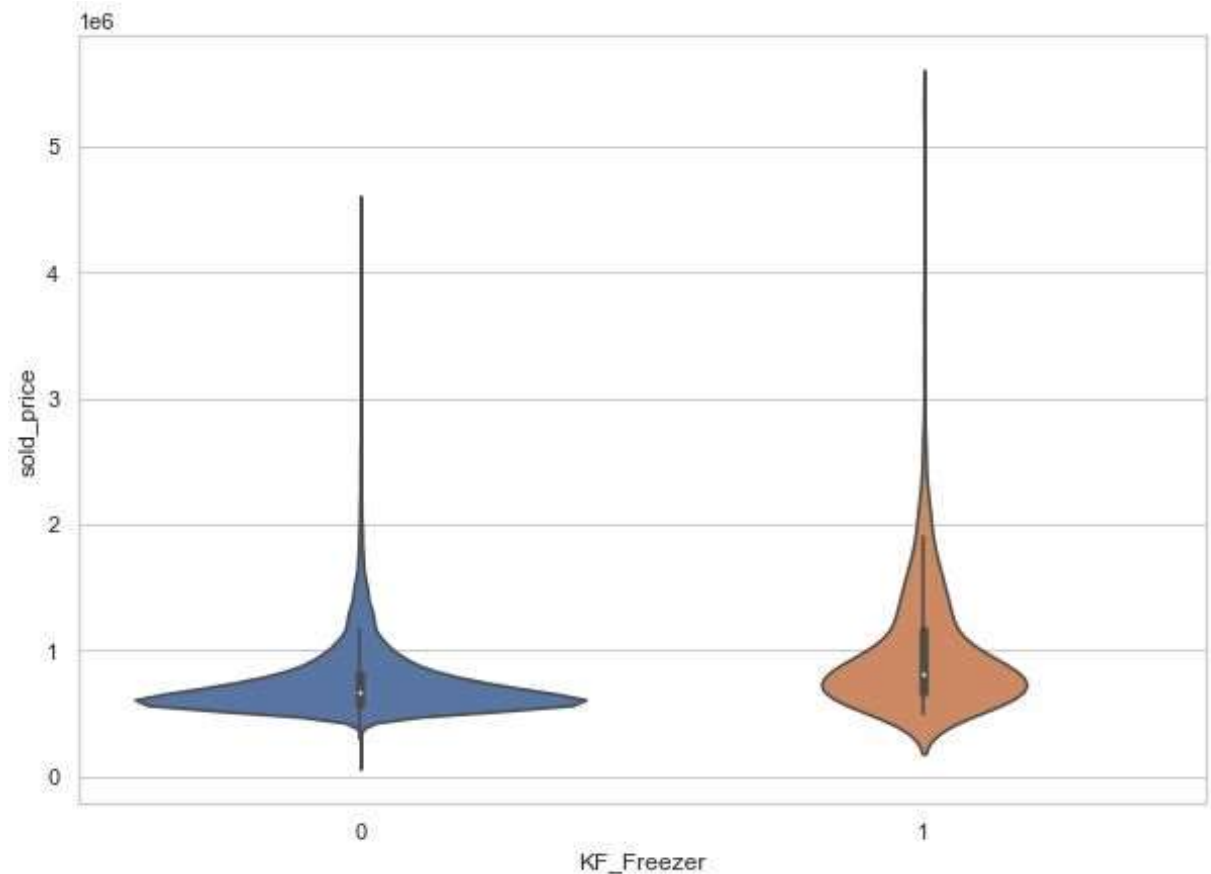
# EDA: Price houses and year of construction

- There is a positive association between the year of construction of the house and its price.



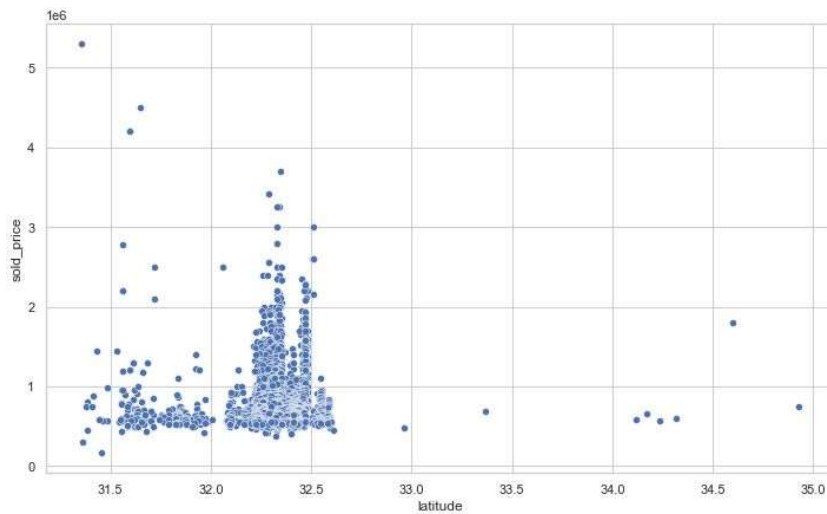
## EDA: Price houses and kitchen features

- Most of the features in the kitchen were not significantly at explaining the price of the houses.
- Although the medium price of the houses are similar, houses with a freezer in the kitchen have a higher price than houses without freezer.

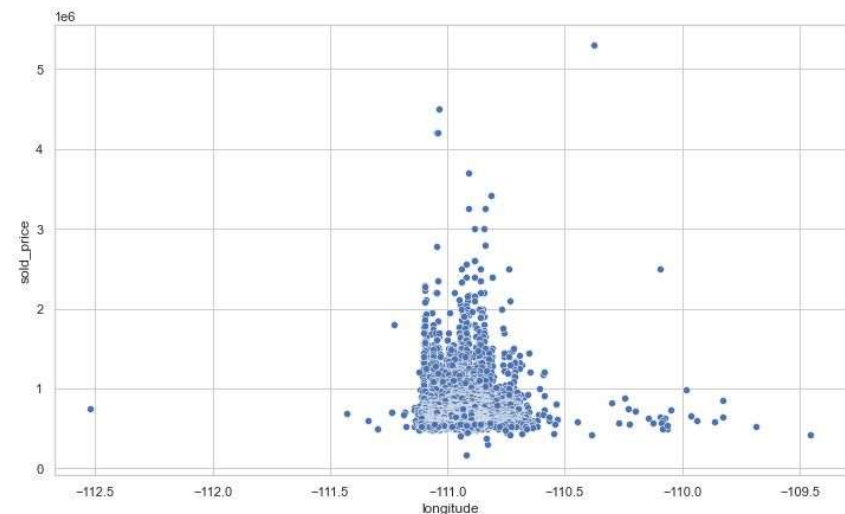


# EDA: Price houses, latitude and longitude

**Latitude**



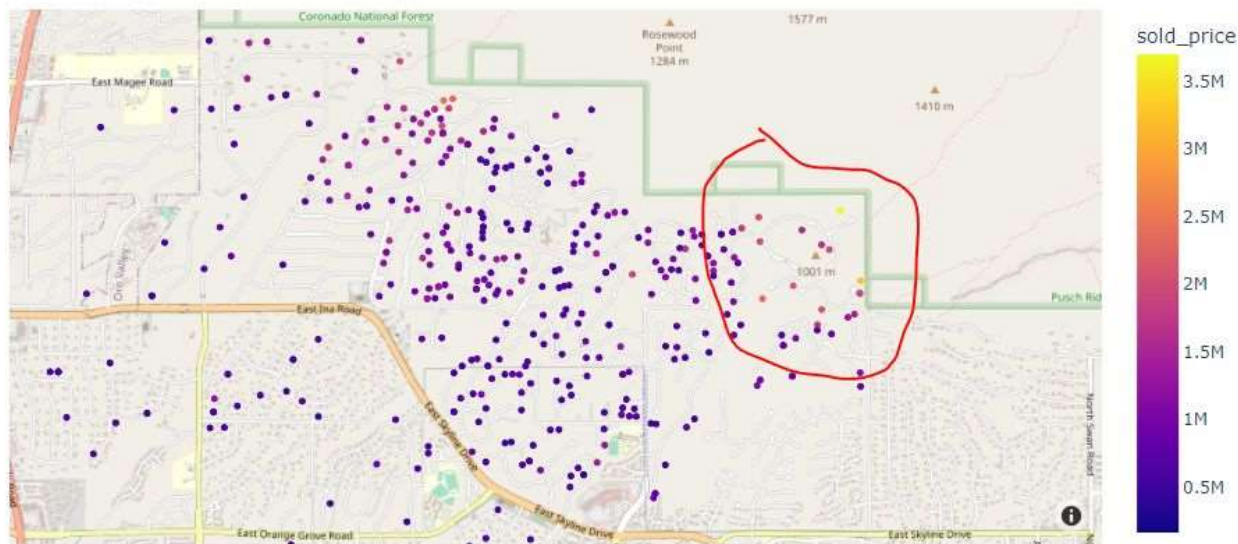
**Longitude**



- Most of the features in the kitchen were not significantly at explaining the price of the houses.
- Although the medium price of the houses are similar, houses with a freezer in the kitchen have a higher price than houses without freezer.

# EDA: Price houses, latitude and longitude

- The most expensive houses are located in a particular area of Houston which is difficult to find with a simple linear regression. A model base tree could obtain a better performance in this case.



# Conclusions

- The database was cleaning consist on filling missing information, handling bad records, and separating key features of the kitchen and floor covering.
- The most important information to predict house prices are square footage, year of constructions, bathrooms, lot\_acres, and fireplaces.
- There is not a positive association between longitude and latitude, and price house, but we can split specific areas in which we can find high value houses.
- Low association of house prices and the rest of variables. For example, features of kitchen and floor covering do not seem to be relevant.