

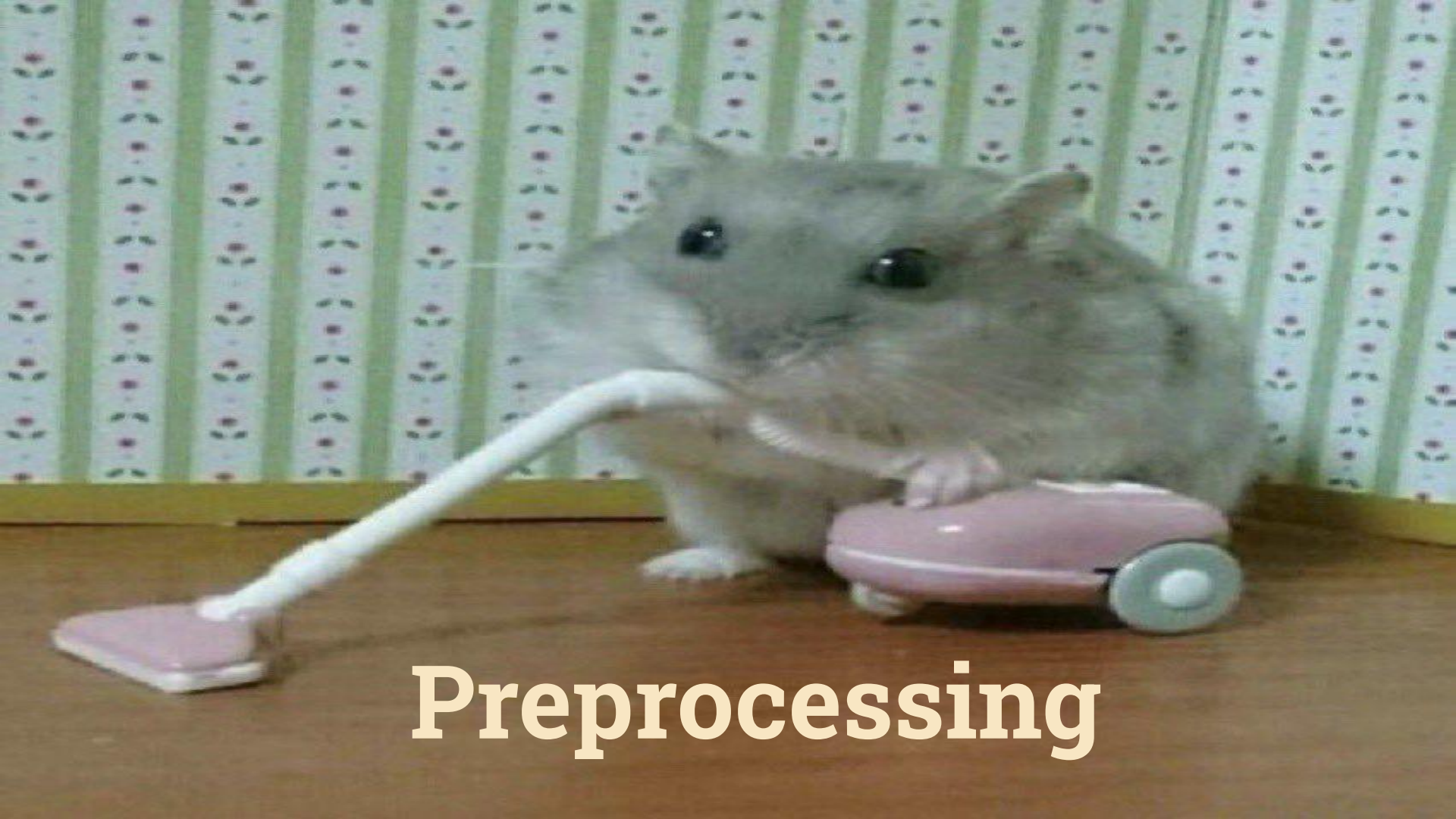


# Machine Learning



Final presentation





**Preprocessing**

# Data Cleaning steps

---

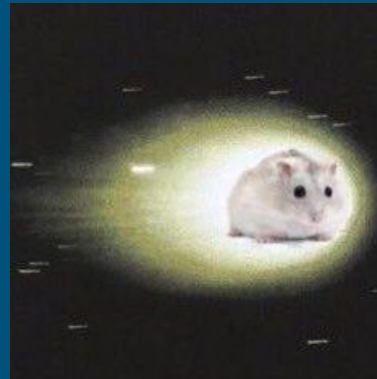
- remove columns we assume are unnecessary and have no effect on the price
- drop columns with high amounts of NaN
- reduce text columns to only numerical or boolean values
- use LabelEncoder() for numeric values
- remove outliers with a price above 2500€, empty URL links or unusual values
- turn amenities into a binary matrix and filter according to luxury keywords
- calculate the median price (calendar listings) and replace missing price values

# Data Engineering

---



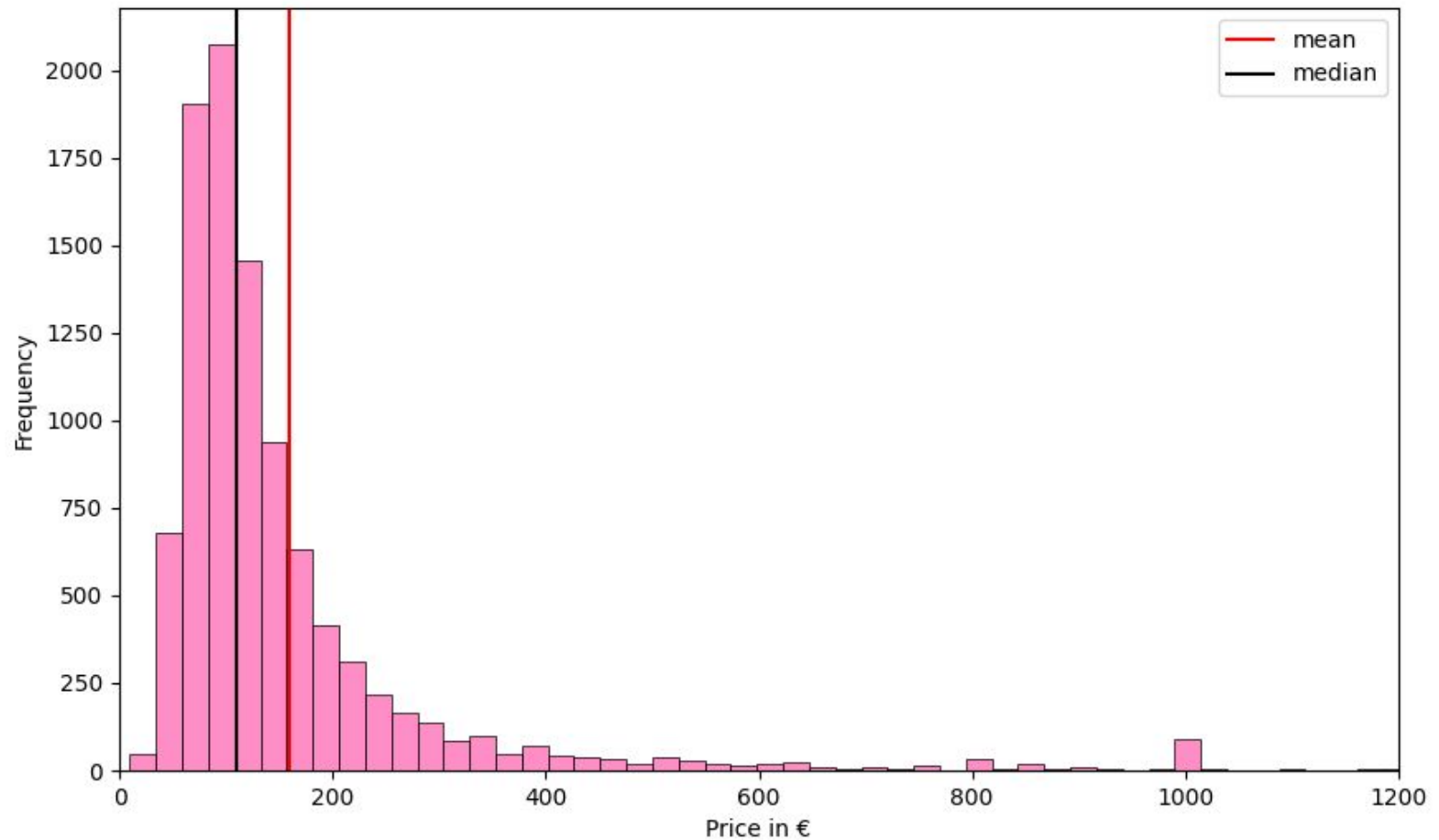
- Sentiment score based on listings reviews
- calculate distance to city center



# Visualization

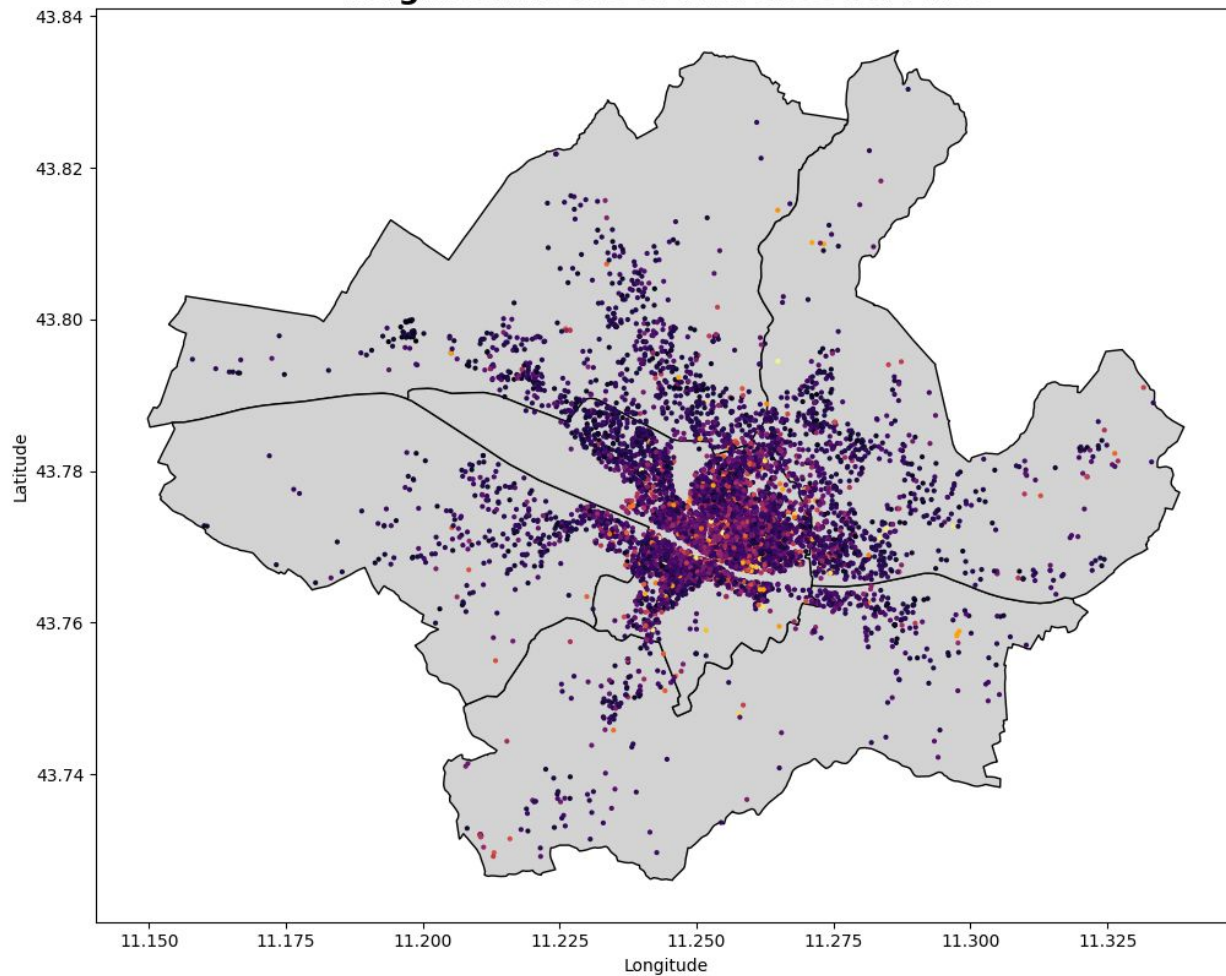


# Price Distribution

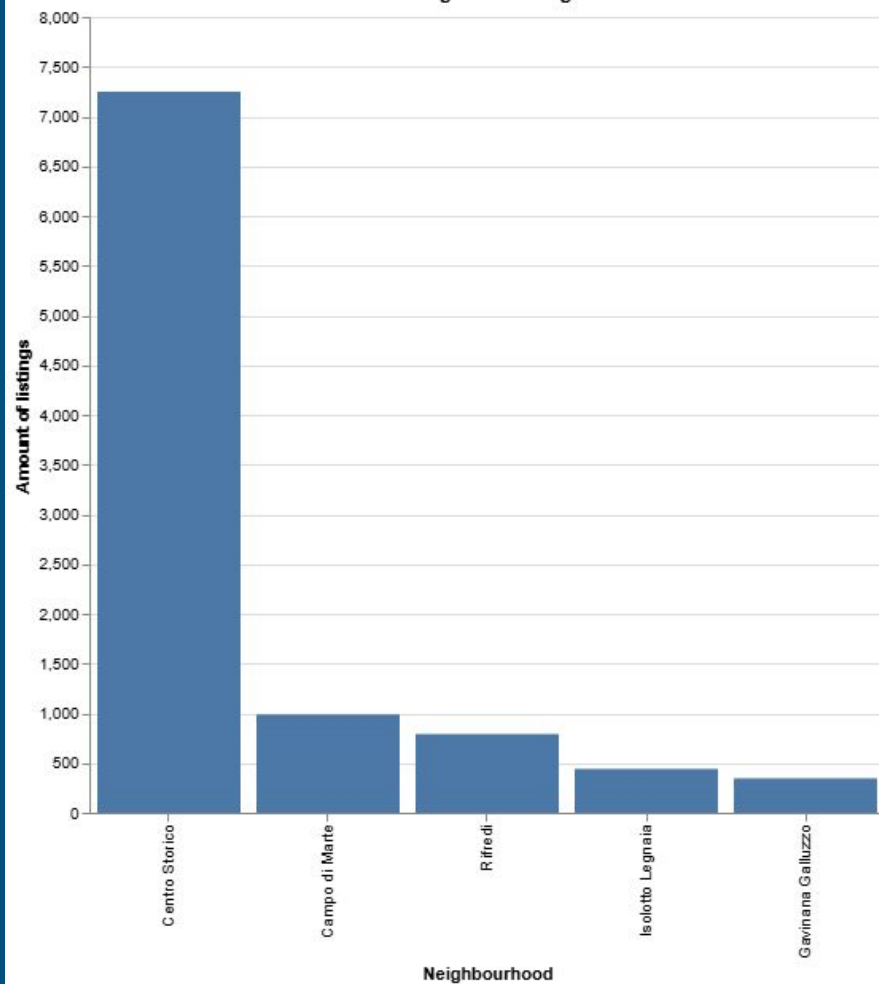




## Neighbourhoods of Florence vs. Price

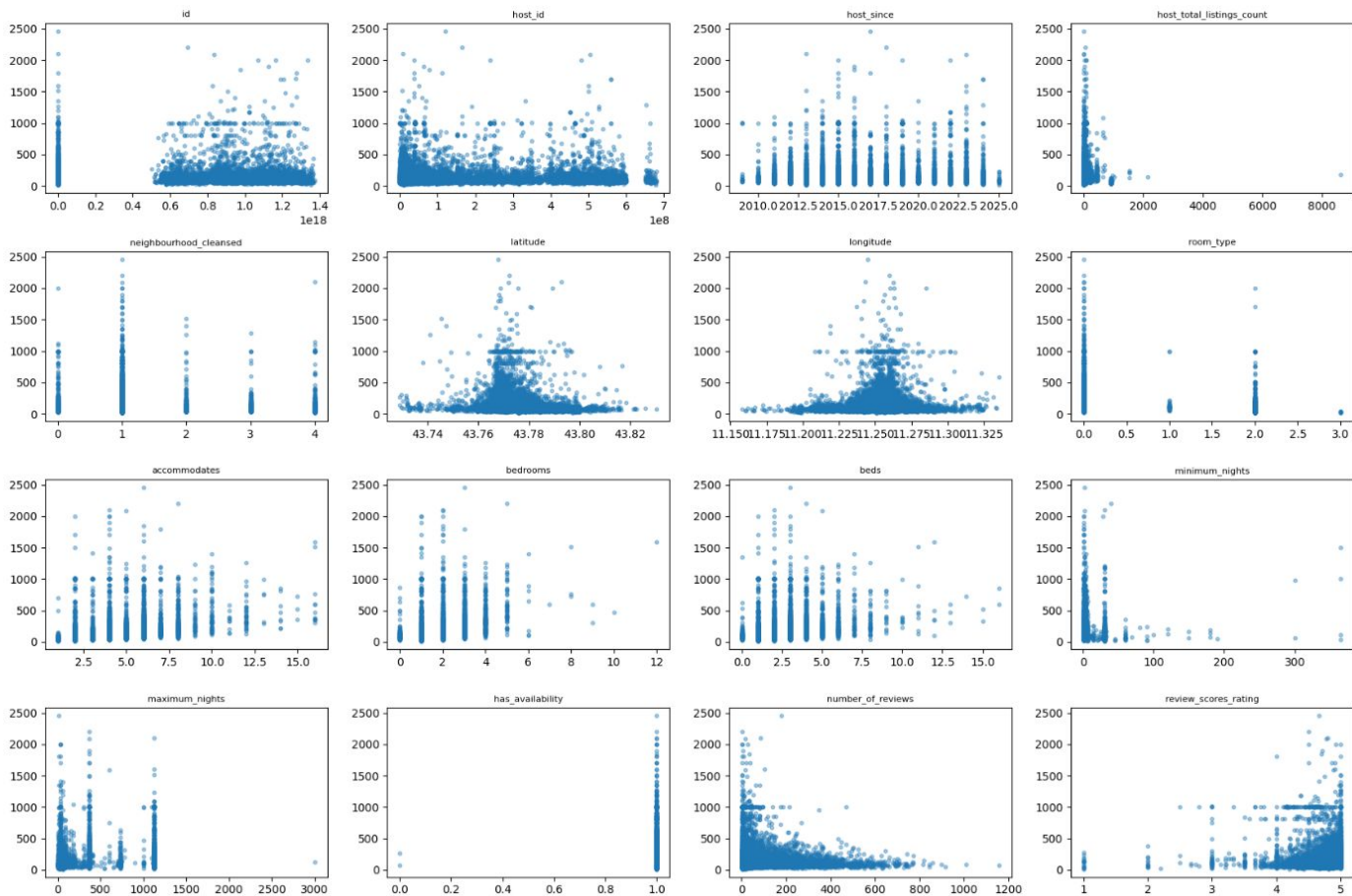


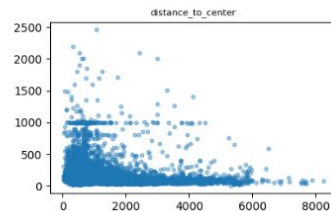
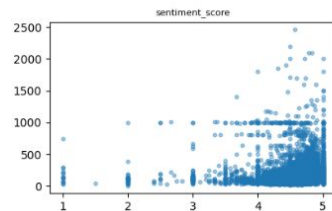
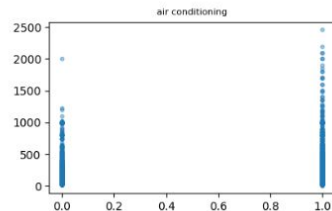
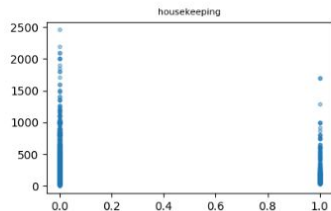
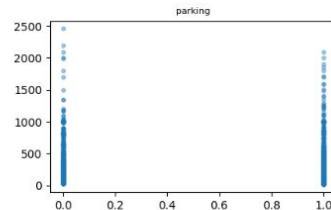
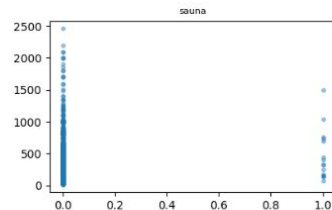
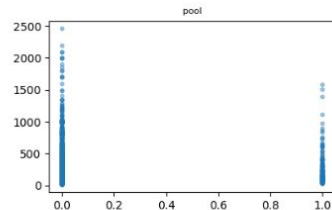
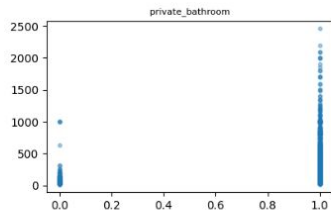
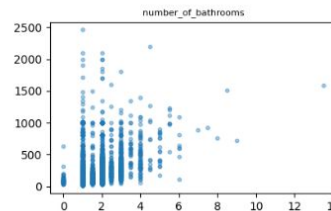
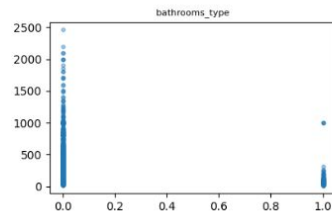
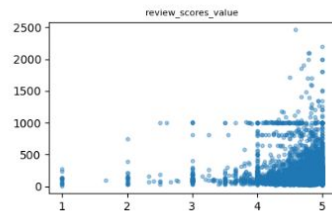
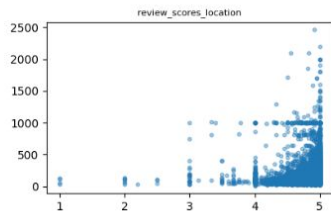
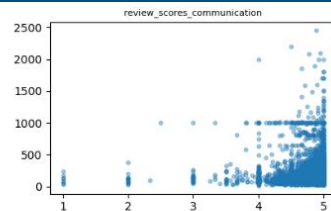
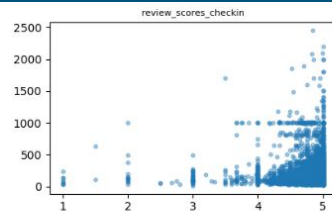
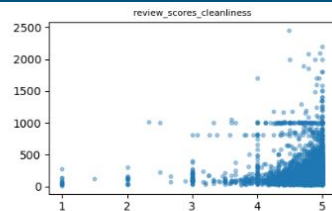
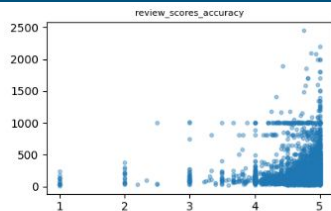
Amount of listings in the neighbourhoods



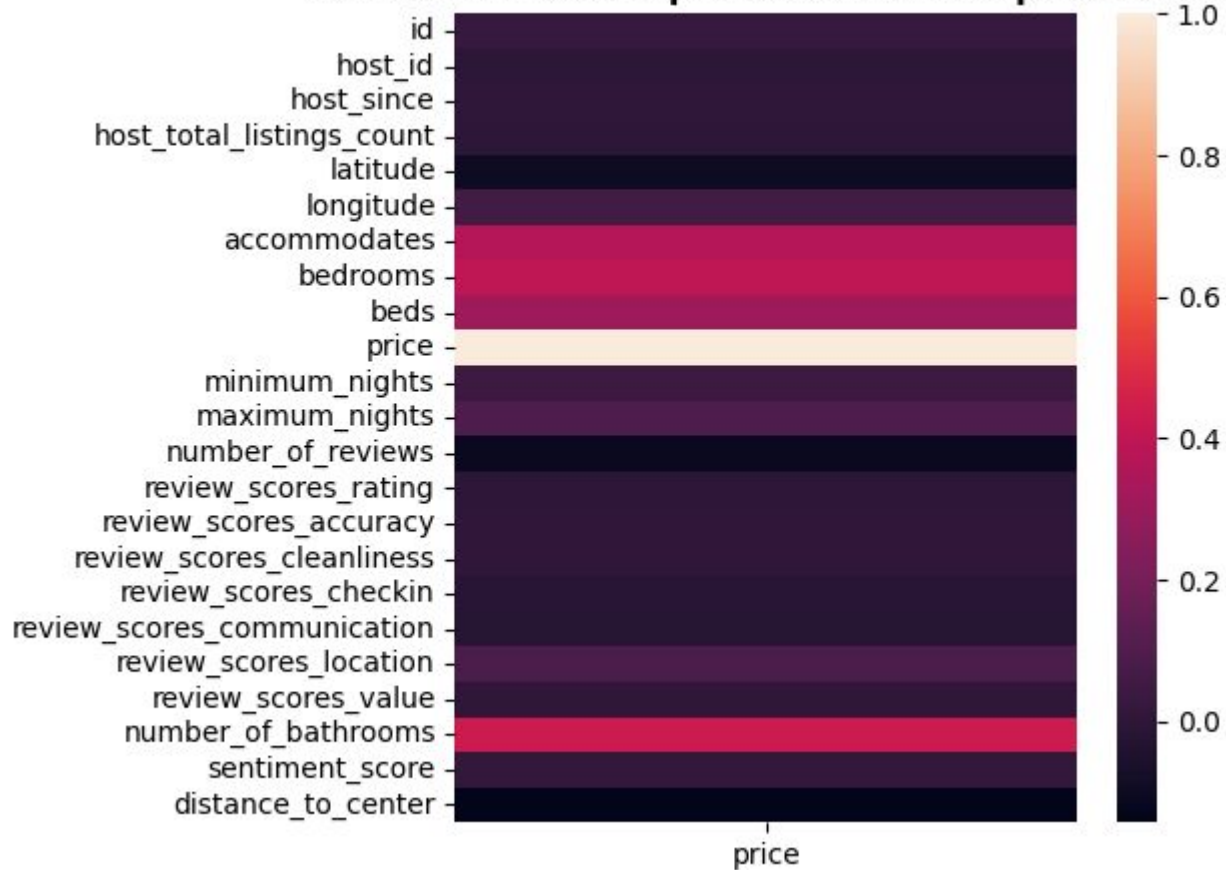


## AirBnB attributes vs. Price

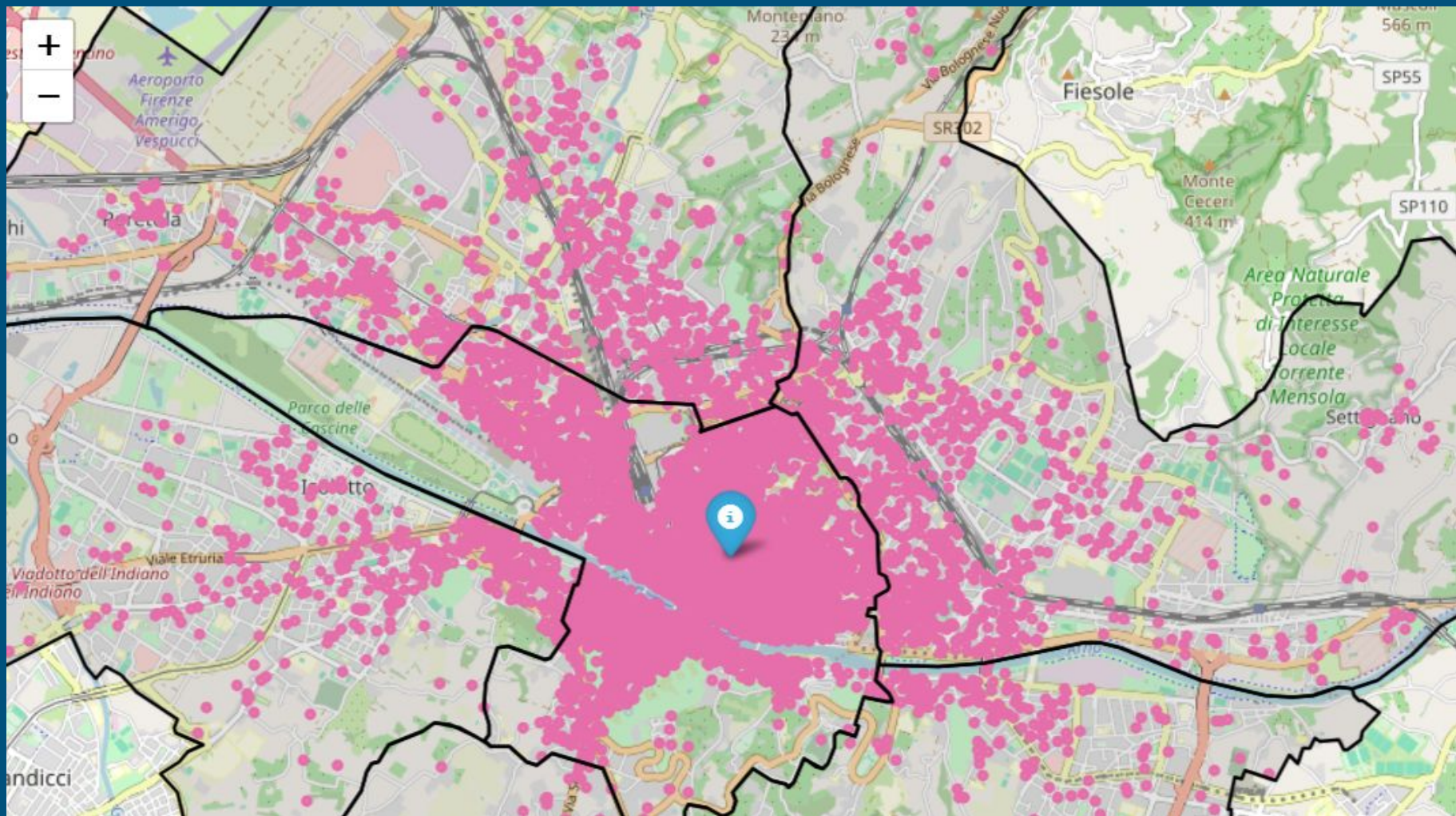




**Correlation Heatmap between Airbnb parameters**







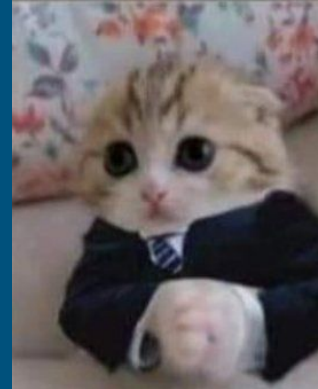
# Model training





# Background Knowledge

---



## Loss functions:

- **L1 Loss (MAE):** average absolute difference between predicted and true values -> robust to outliers, harder to optimize in gradient based methods
- **L2 Loss (MSE):** penalizes large errors more than small ones -> sensitive to outliers, for regression task
- **RMSE Loss:** same unit as targets -> easier to interpret, still sensitive to outliers
- **Huber Loss:** combines L1 and L2 Loss -> less sensitive to outliers than MSE, easier to optimize than MAE

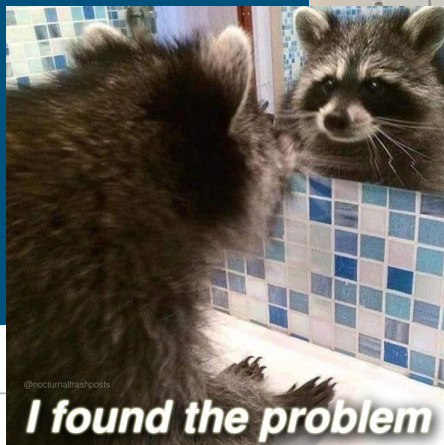
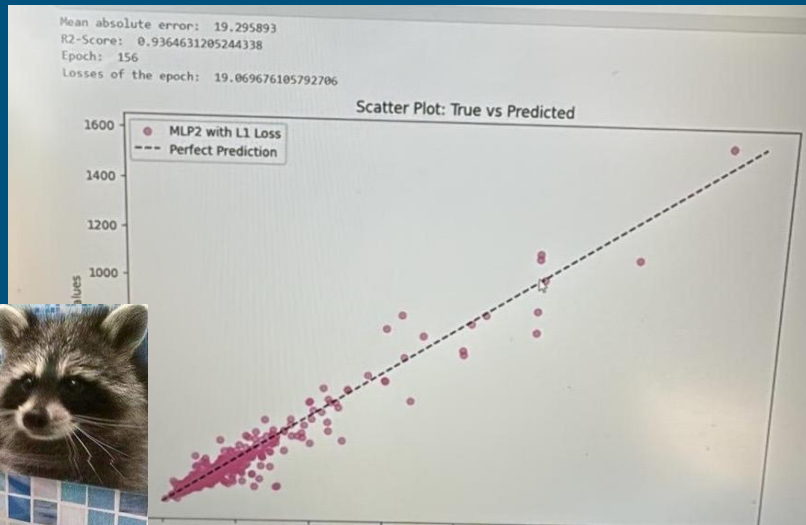
## Parameters:

- **learning rate:** controls how fast the model learns,
  - too high -> overshooting, unstable, too low -> slow convergence
- **weight decay:** adds penalty on model weights to prevent overfitting



# Data Leakage Problem

- We thought we had the perfect MLP
- $R^2$ : 0.9+
- Checked everything
- normalization on the whole data frame
- Solution:



## Ostrich algorithm

[Article](#) [Talk](#)

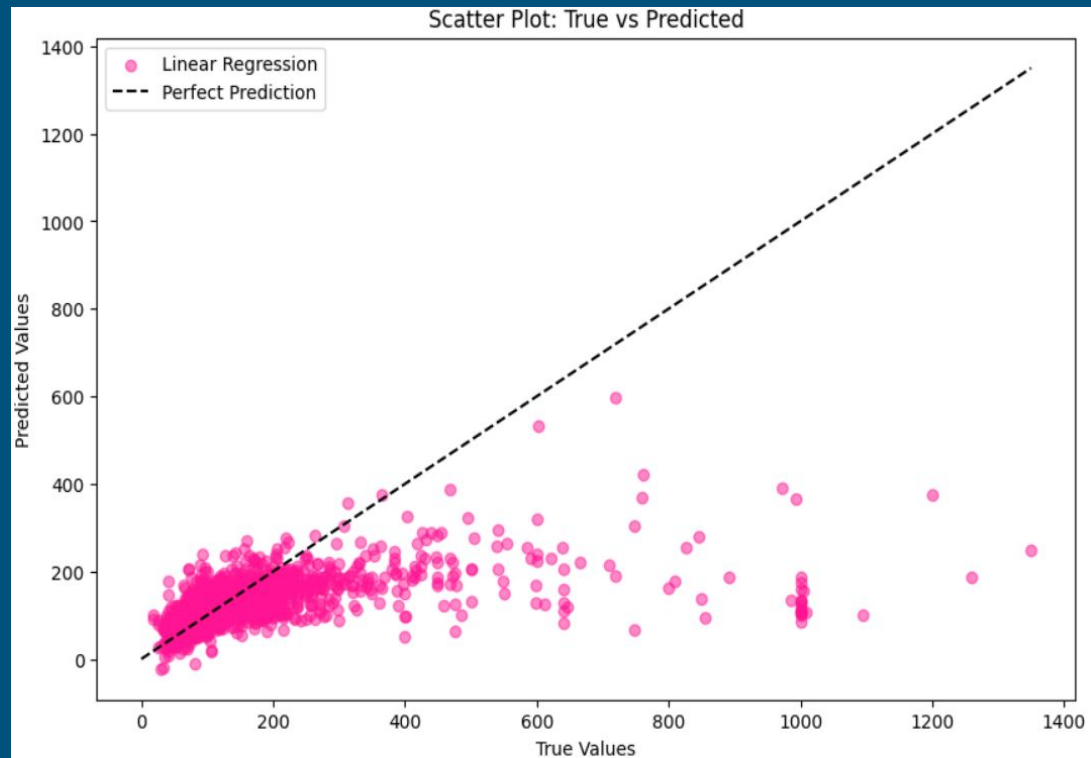
[14 languages](#)

[Read](#) [Edit](#) [View history](#) [Tools](#)

From Wikipedia, the free encyclopedia

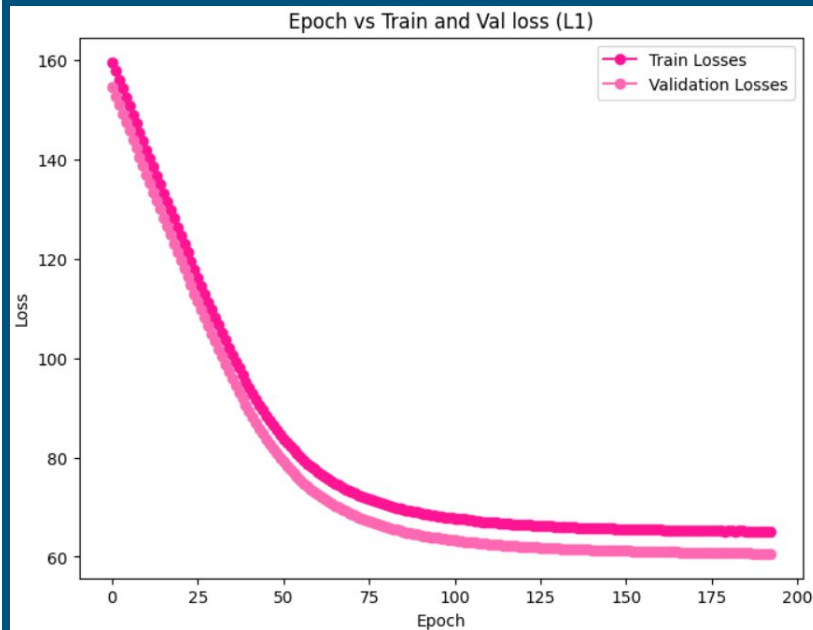
In [computer science](#), the **ostrich algorithm** is a strategy of ignoring potential problems on the basis that they may be exceedingly rare. It is named after the [ostrich effect](#) which is defined as "to stick one's head in the sand and pretend there is no problem". It is used when it appears the situation may be more cost-effectively managed by allowing the problem to continue to occur rather than to attempt its prevention.

# Linear Regression



## Hyperparameters

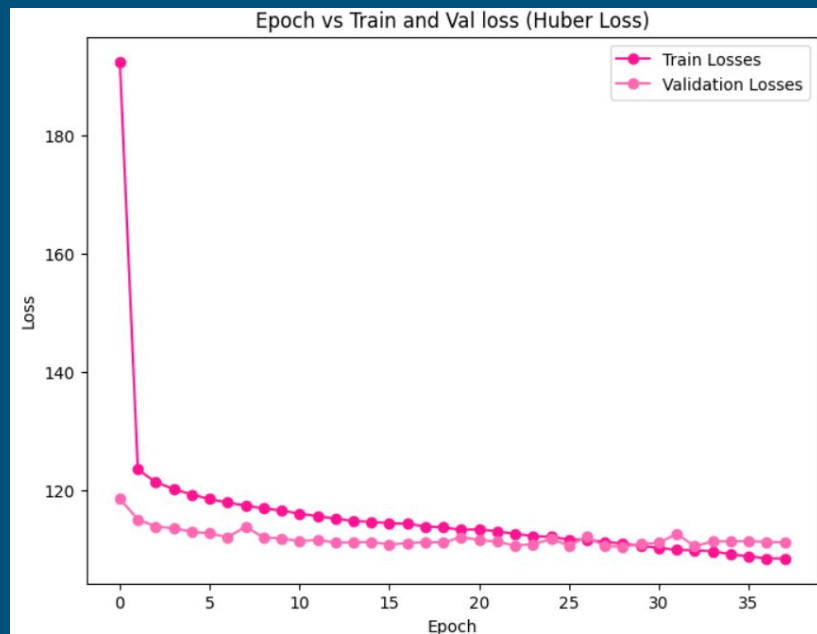
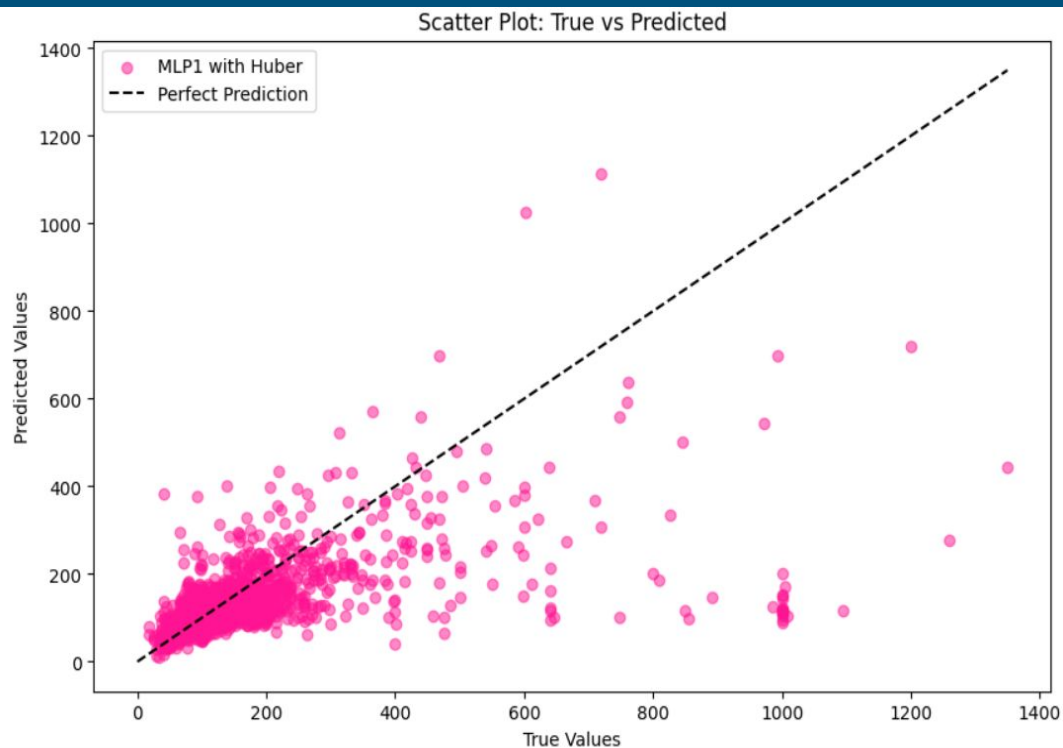
- Loss Function: L1 Loss
- Learning Rate: 0.008
- Optimizer: SGD



# MLP1 - 4 Layer Perceptron

## Hyperparameters

- Loss Function: Huber Loss
- Learning Rate: 0.001
- Optimizer: Adam
- weight decay: 0.0005

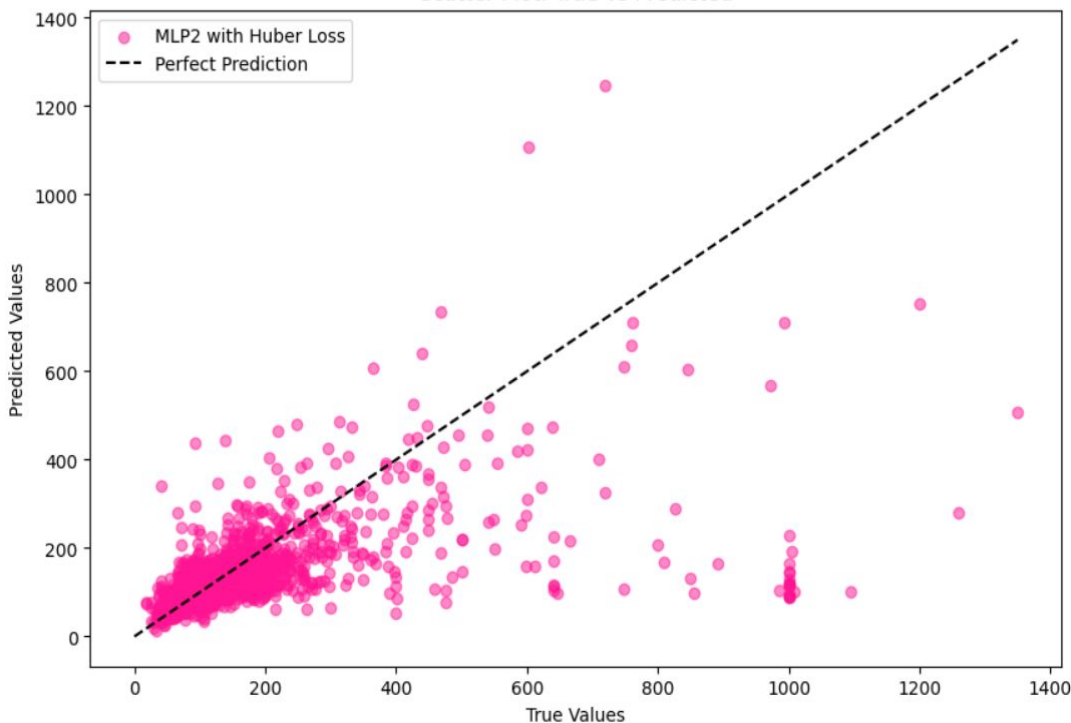


# MLP2 - 5 Layer Perceptron

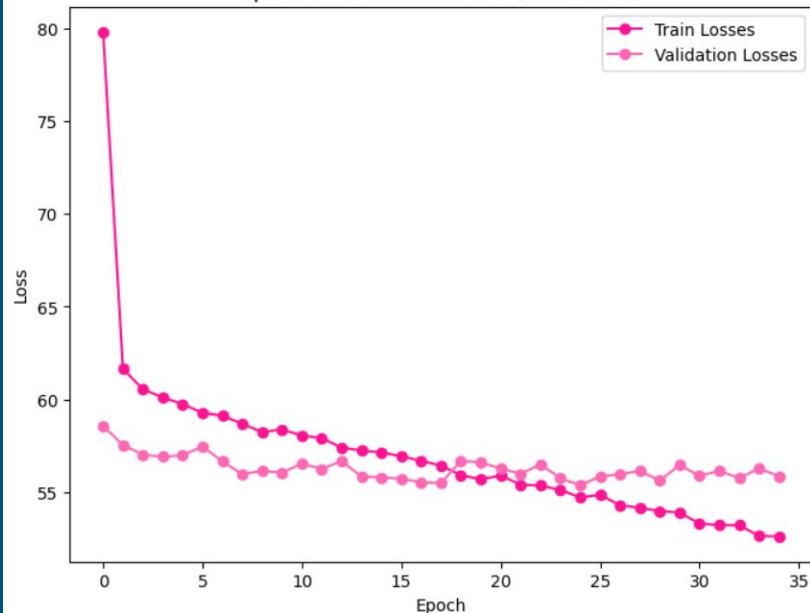
## Hyperparameters

- Loss Function: L1 Loss
- Learning Rate: 0.001
- Optimizer: Adam
- weight decay: 0.0005

Scatter Plot: True vs Predicted



Epoch vs Train and val loss (Huber Loss)

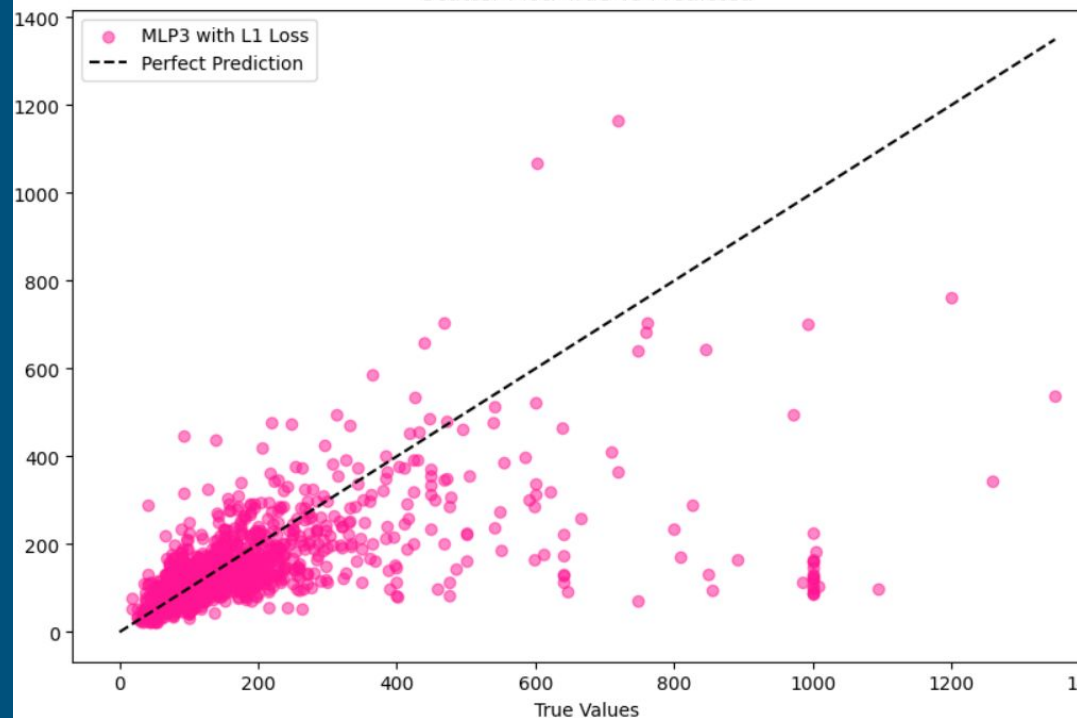


# MLP3 - 6 Layer Perceptron

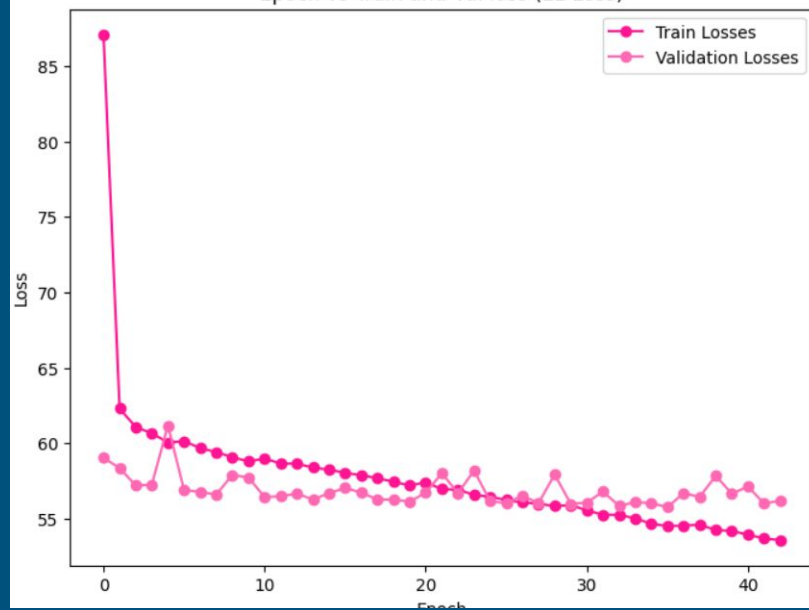
## Hyperparameters

- Loss Function: L1 Loss
- Learning Rate: 0.001
- Optimizer: Adam
- weight decay: 0.005

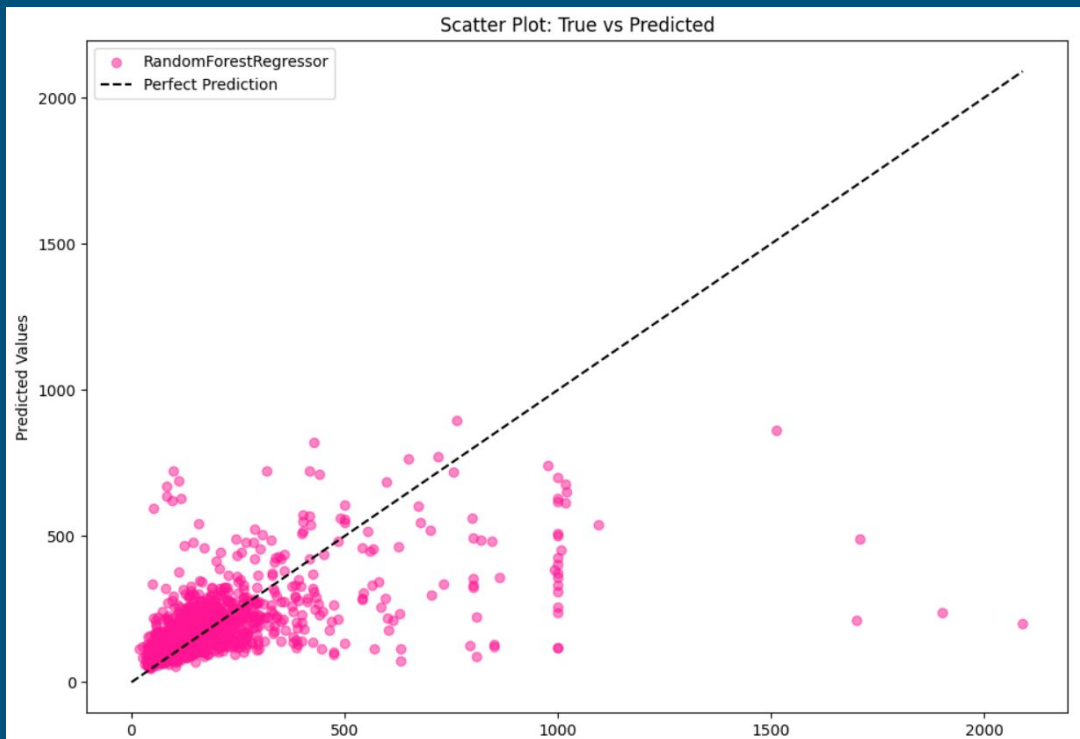
Scatter Plot: True vs Predicted



Epoch vs Train and Val loss (L1 Loss)



# RandomForestRegressor



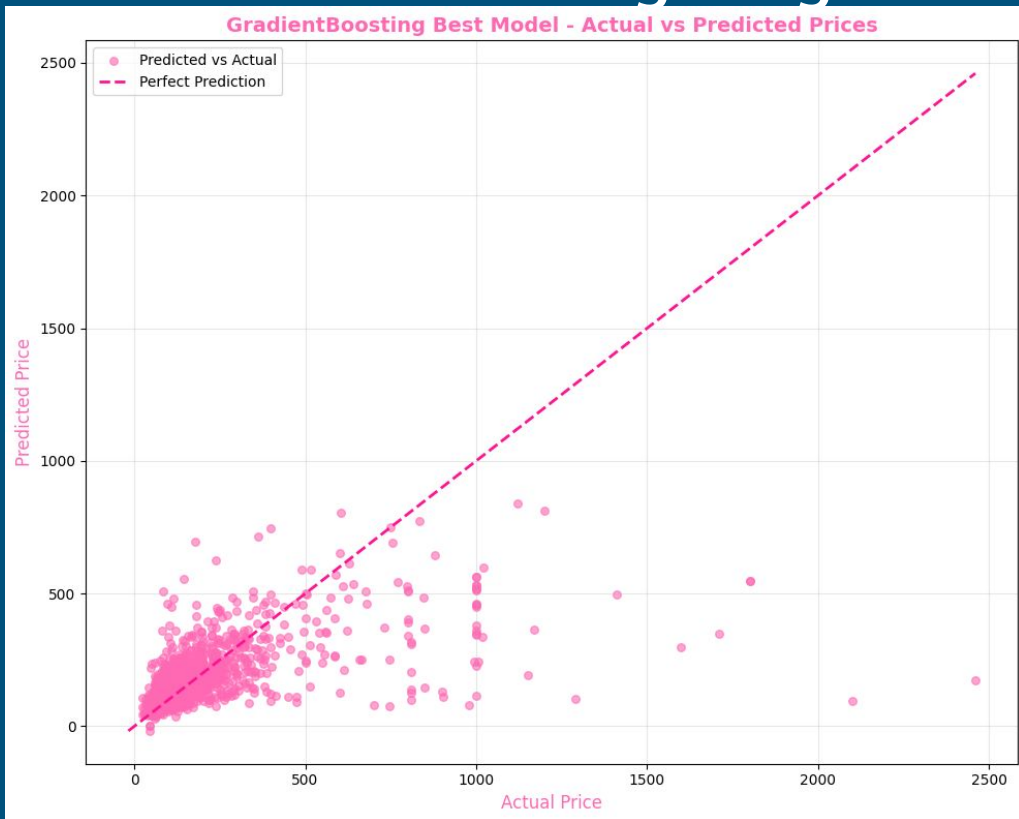
## Hyperparameters:

- **max\_depth = 10**
- **max\_features = 30**
- **min\_samples\_leaf = 4**
- **n\_estimators = 500**

→ found with GridSearchCV()



# Gradientboosting Regressor

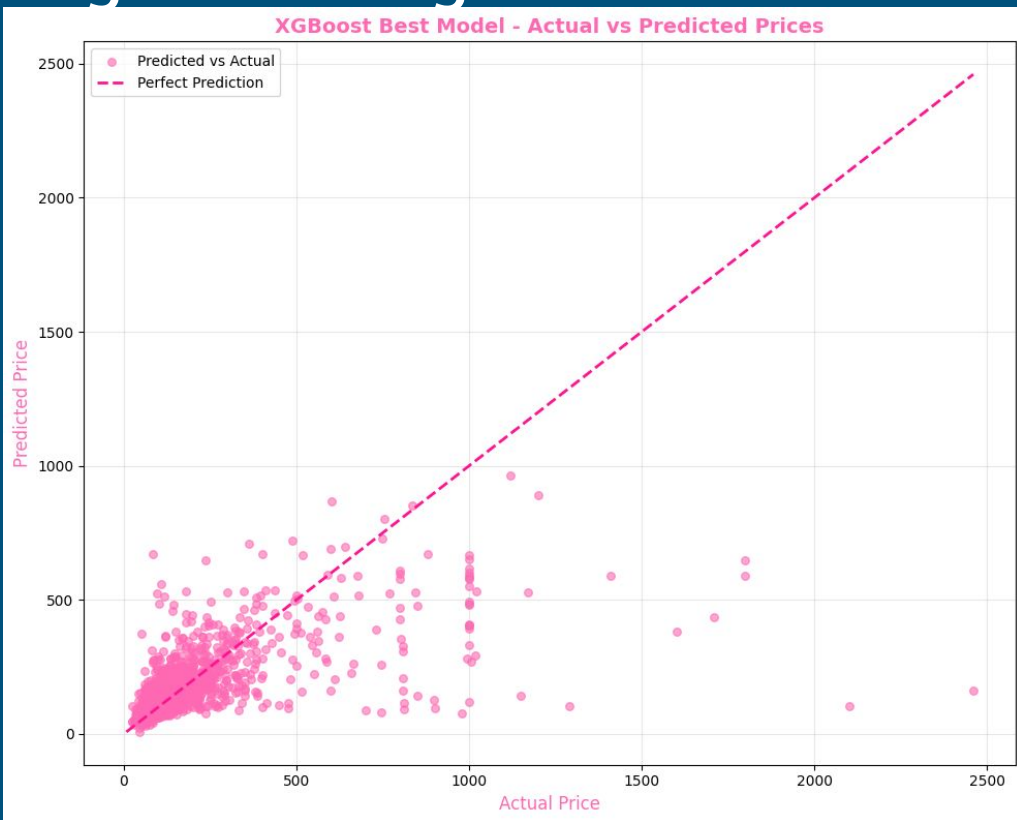


## Hyperparameters:

- **n\_estimators = 500**
- **min\_samples\_split = 10**
- **min\_samples\_leaf = 15**
- **max\_depth = 4**
- **learning\_rate = 0.05**
- **subsamples = 0.8**

→ found with RandomizedSearchCV()

# XgboostRegressor



## Hyperparameters:

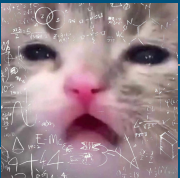

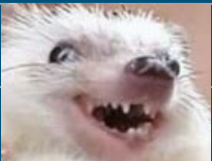
- **n\_estimators = 500**
- **min\_child\_weight = 5**
- **max\_depth = 5**
- **learning\_rate = 0.03**
- **reg\_alpha = 0.1**
- **reg\_lambda = 10**

→ found with RandomizedSearchCV()

# Scores Comparison



# Comparison

Model Name	$R^2$		MAE	$R^2$ with PCA	MAE with PCA
Linear Regression	0.18		61.23	0.024	58.83
MLP1	0.27		56.81	0.26	64
MLP2	0.27		56.37	0.24	51.41
MLP3	0.29		56.31	0.24	51.1
RandomForestRegressor	0.37		62.9	-	-
GradientBoosting	0.36		65.38	-	-
XGBoost	0.40		62.30	-	-

→results from the latest model run

Evaluation



# Evaluation

Problems we've encountered:

- normalized with the wrong mean and standard deviation
- skewed results due to specific outliers
- missing price values
- urls that lead to nothing
- the models tend to overfit → change loss function or early stopping patience
- after data cleaning we only had 9000 rows left

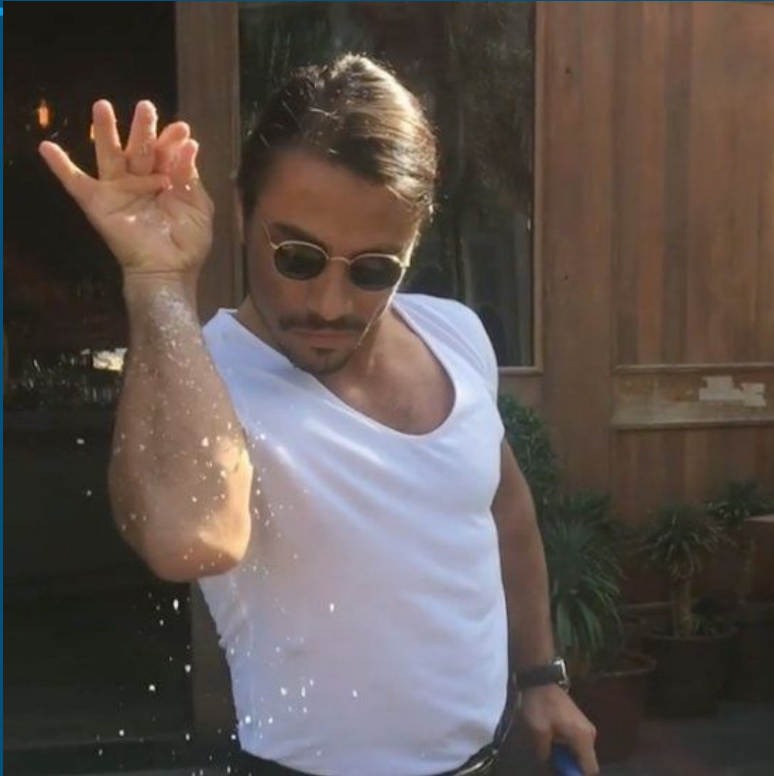
Our insights:

- small difference between using PCA or not, since we do not have a high dimensionality
- we didn't use PCA for the trees, since it doesn't result in a huge difference
- different optimizers and losses
- reduced the amount of amenities because of suspected low influence on the price



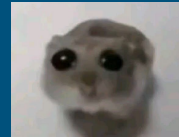


# Conclusion



We cooked

or we were cooked...



# Sources

---

- [3 Big Benefits of Data Cleansing - Teraflow.ai](#) (Data Cleaning Picture)
- [Forbes Model Training](#) (Model Training Picture)
- [Ostrich algorithm - Wikipedia](#) (Ostrich)