

Rent Analysis and Prediction

FOR APARTMENTS IN
GERMANY

FEMI ONAFALUJO

SPRINGBOARD DATA SCIENCE

SEPTEMBER 26, 2021



“The more you know about the past, the better prepared you are for the future.”

- THEODORE ROOSEVELT (1858-1919)

Project Objectives

3

Problem

- Help landlords set the rent for their apartments

Solution

- Develop a predictive model

It is a Journey: Solutions Areas & Scopes

Data Wrangling

Collect and organize data

Clean data

Exploratory Data Analysis

Relationships between rent and numerical features

Relationship between rent and categories

Machine Learning

Training data development

Metrics and testing

Model development

Model selection

Model application

Data Wrangling

COLLECTING DATA | CLEANING DATA

Primary Dataset – Rental Information

6

Source: Immoscout24 – largest real estate website in Germany

Information: Record of apartments for rent

Size: 267,859 records, 49 features

Examples of features

- Heating cost
- Rent (€)
- Living space (square meters)
- Number of rooms
- Interior quality
- Location information (State, city / town, municipality, zip code, street address, house number)
- Facilities (Balcony, garden, cellar, etc.)

Secondary Dataset – State Information

7

Source: Wikipedia

Information: State macro-economic data

Size: 16 states x 5 features

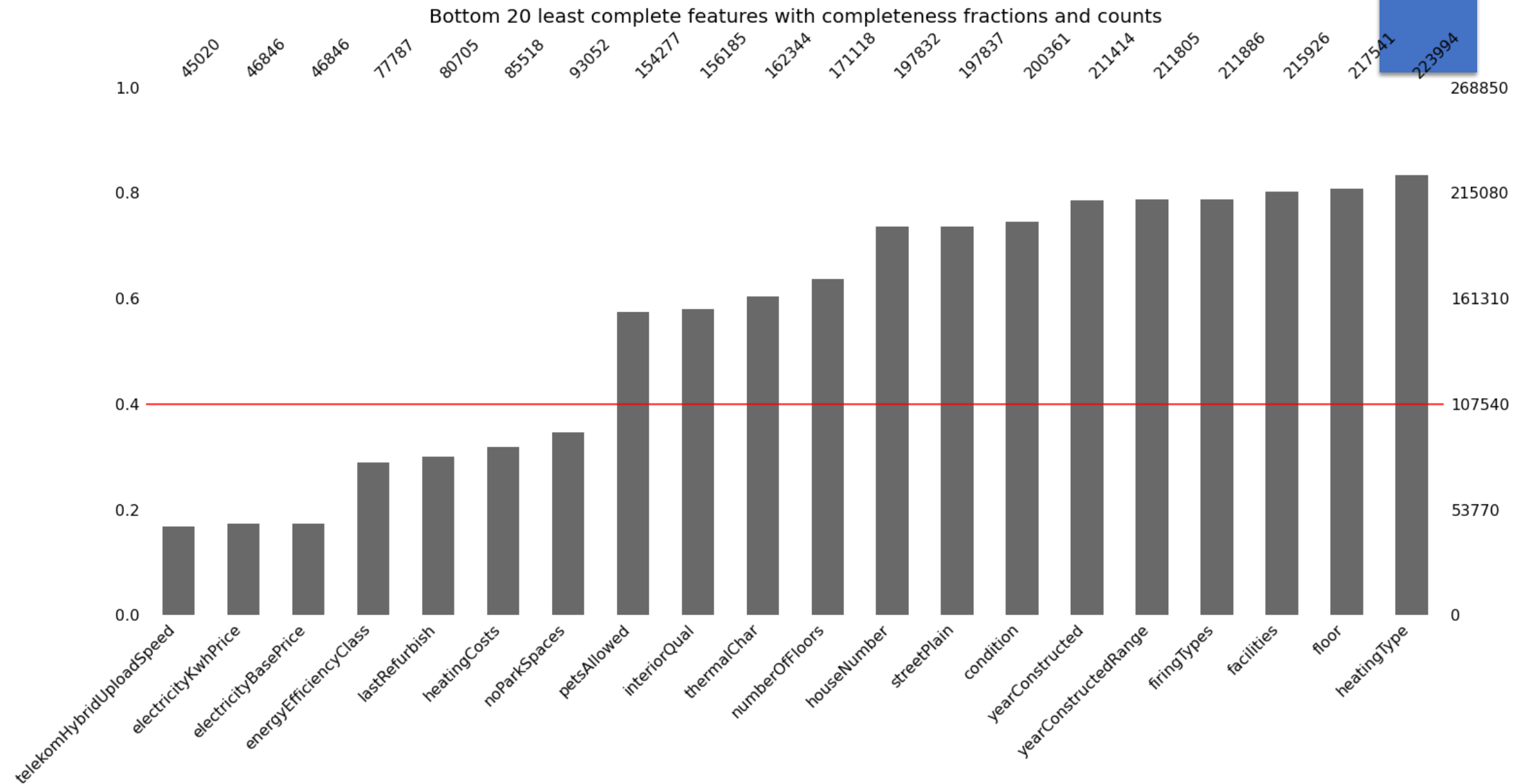
Features

- State area
- State population
- State population per area
- GDP per capita
- Human development index (HDI)

Issues: Cleaning the data

8

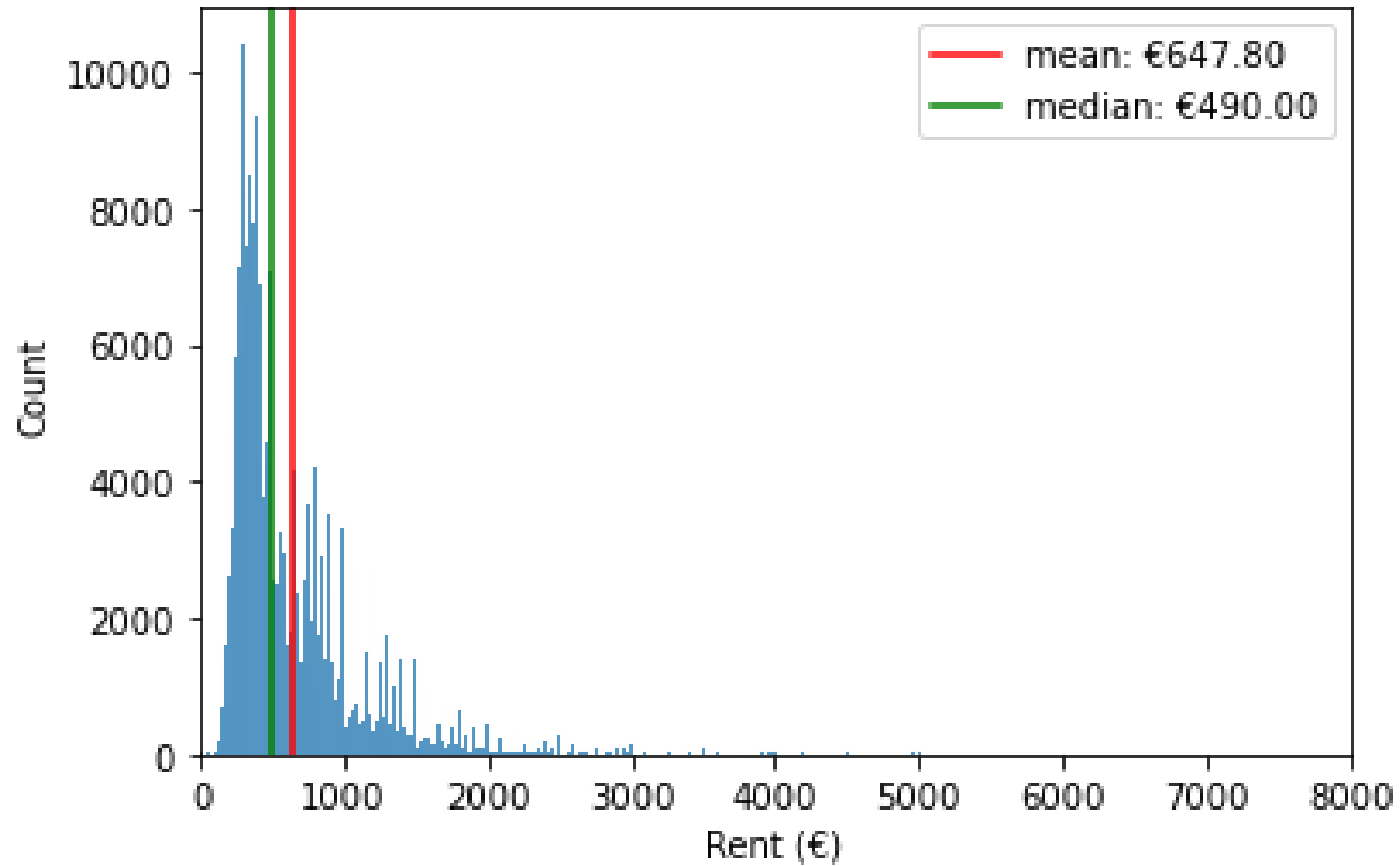


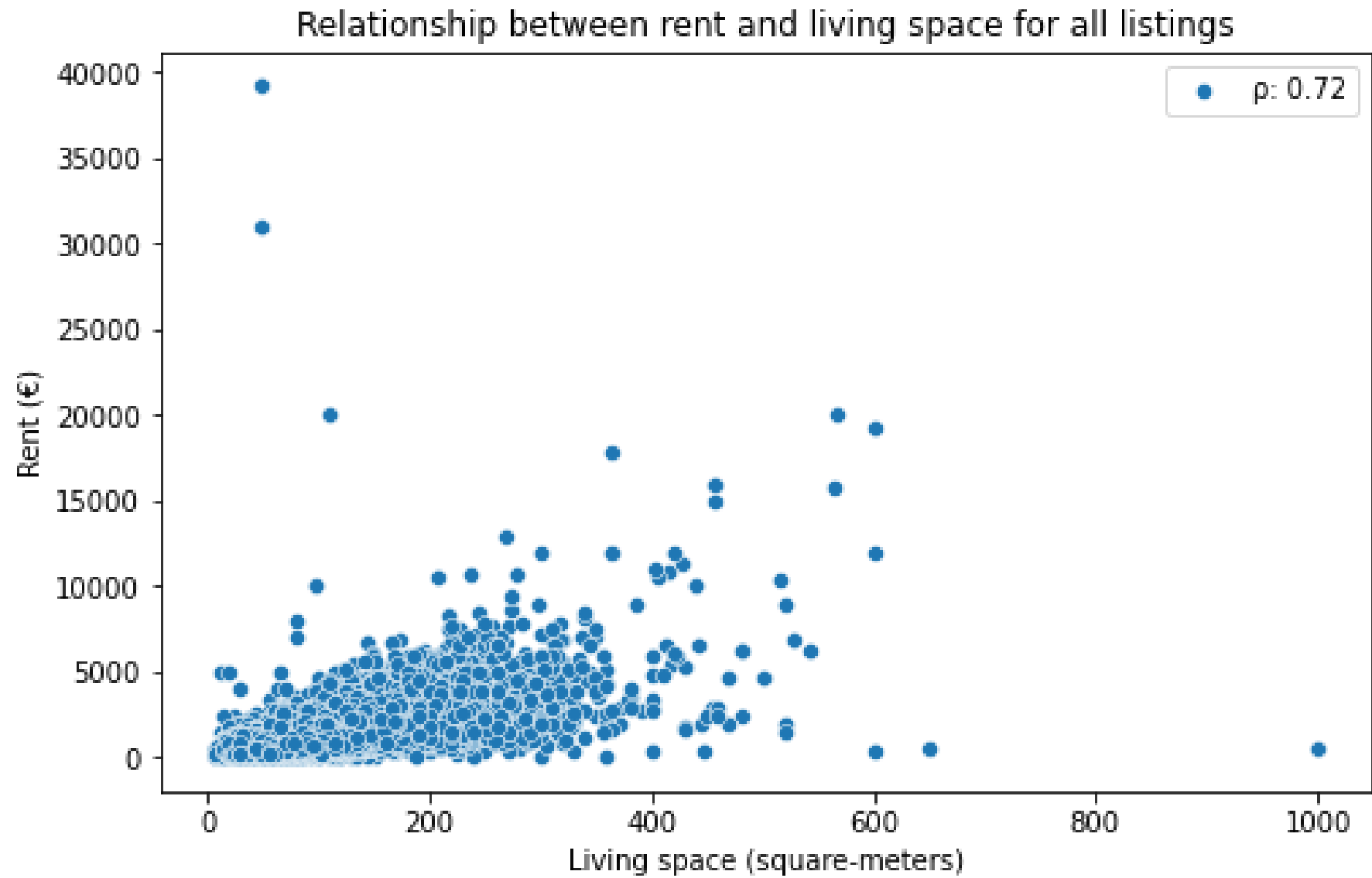


Exploratory Data Analysis

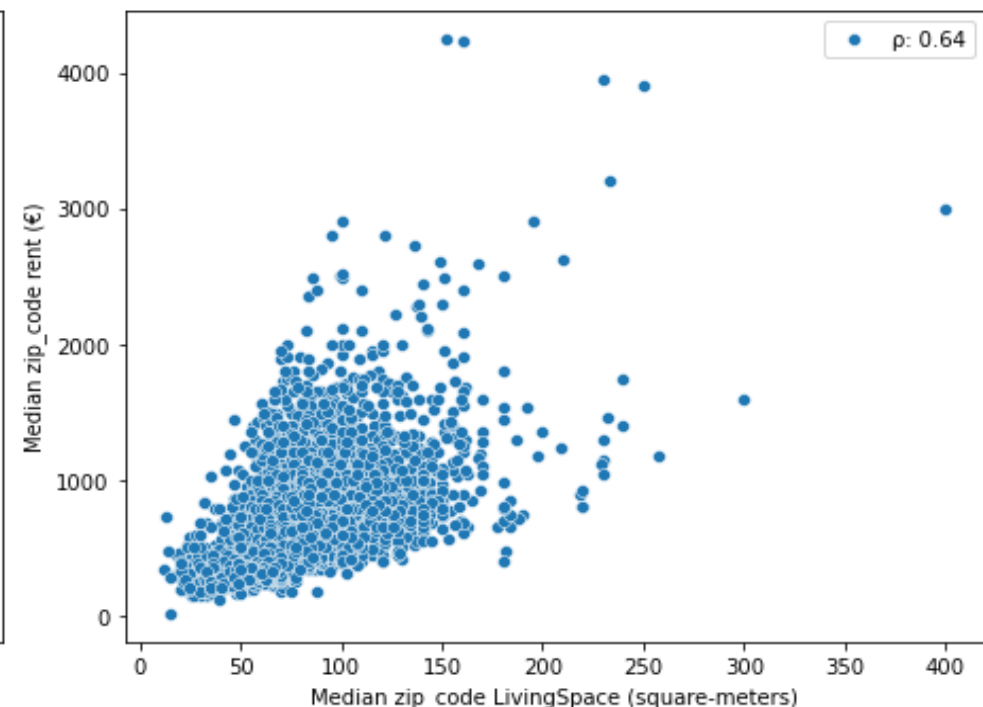
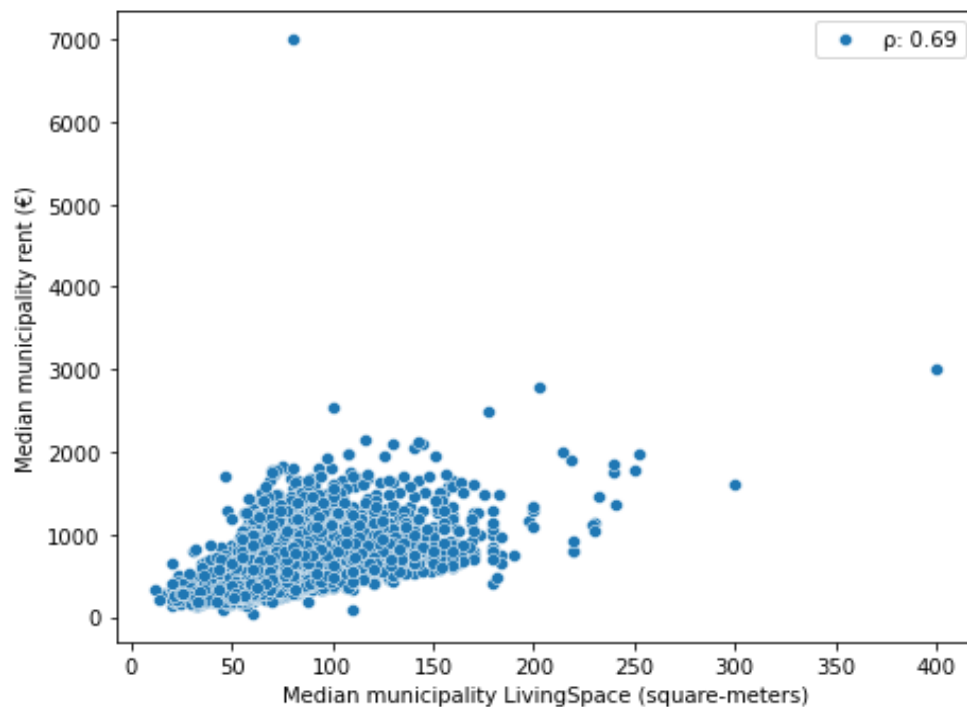
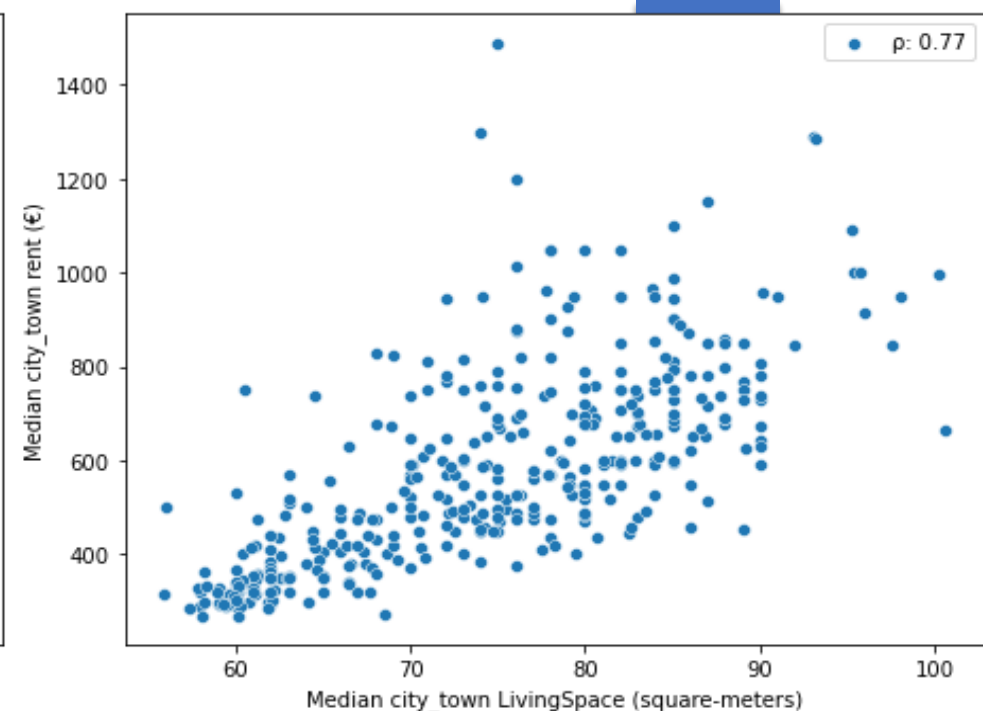
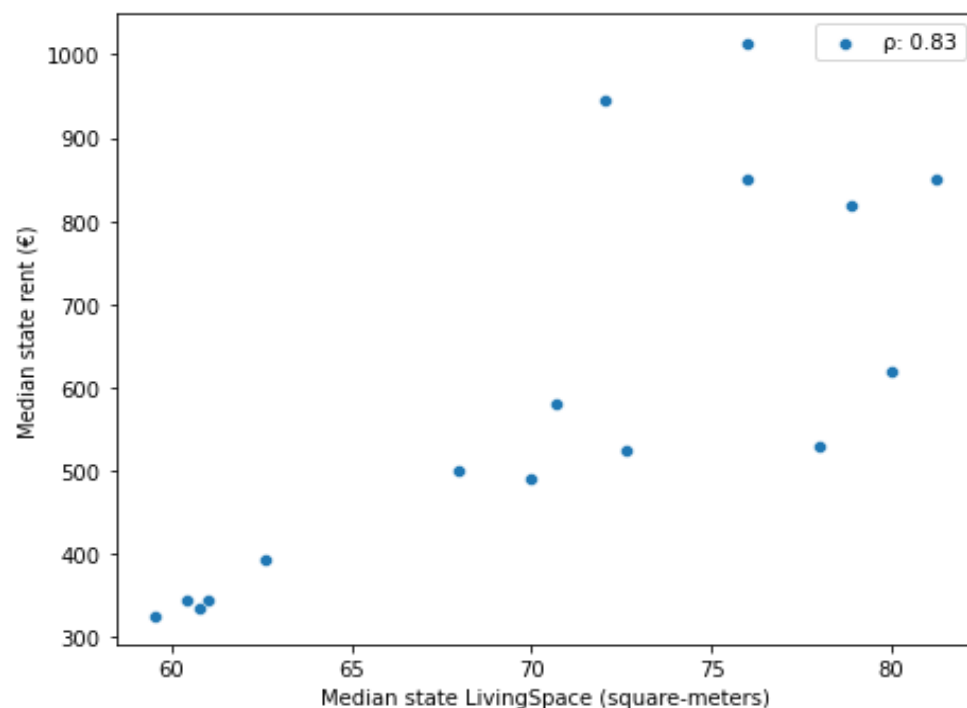
RENT VS. NUMERICAL AND CATEGORICAL FEATURES

Distribution of rent

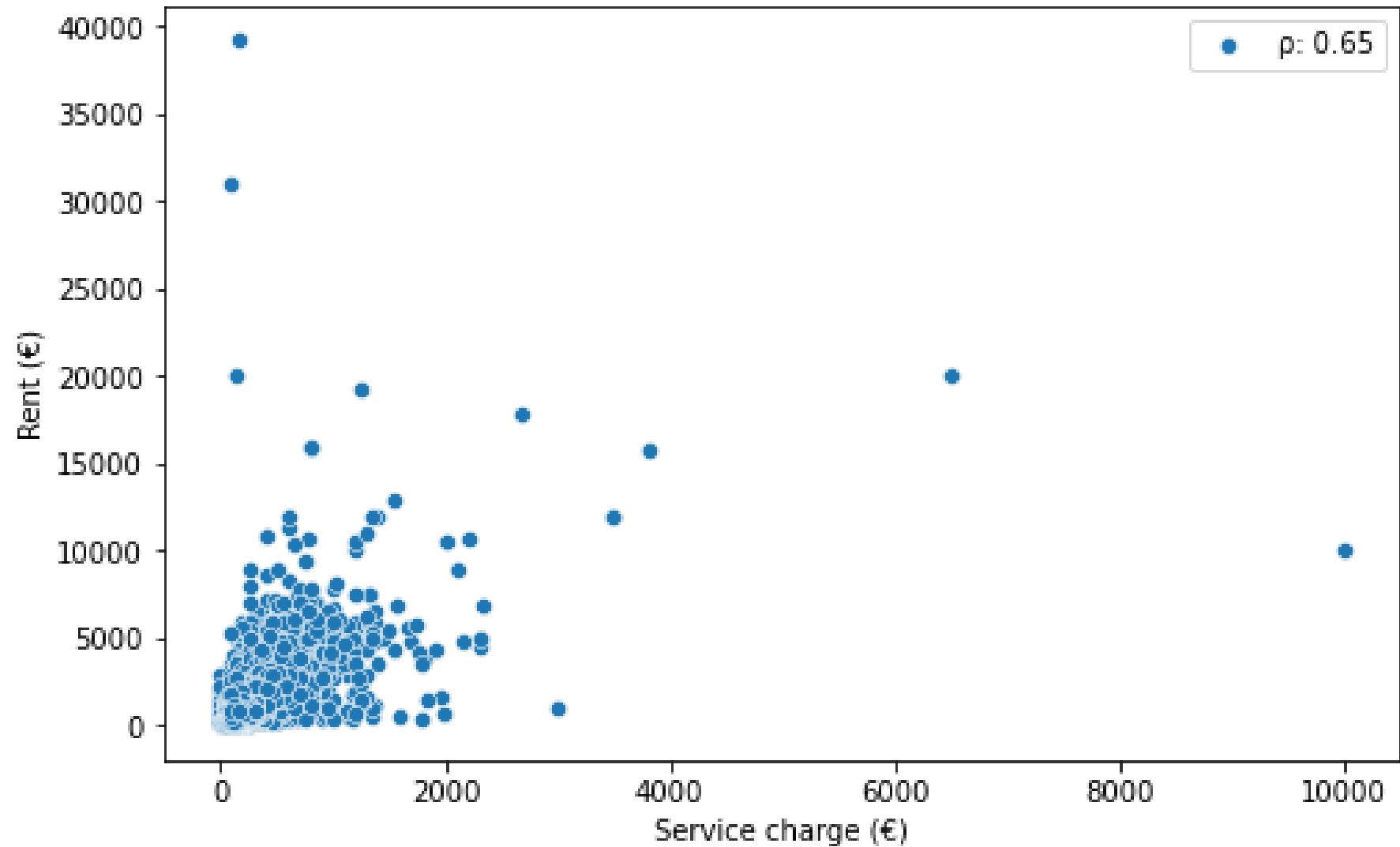




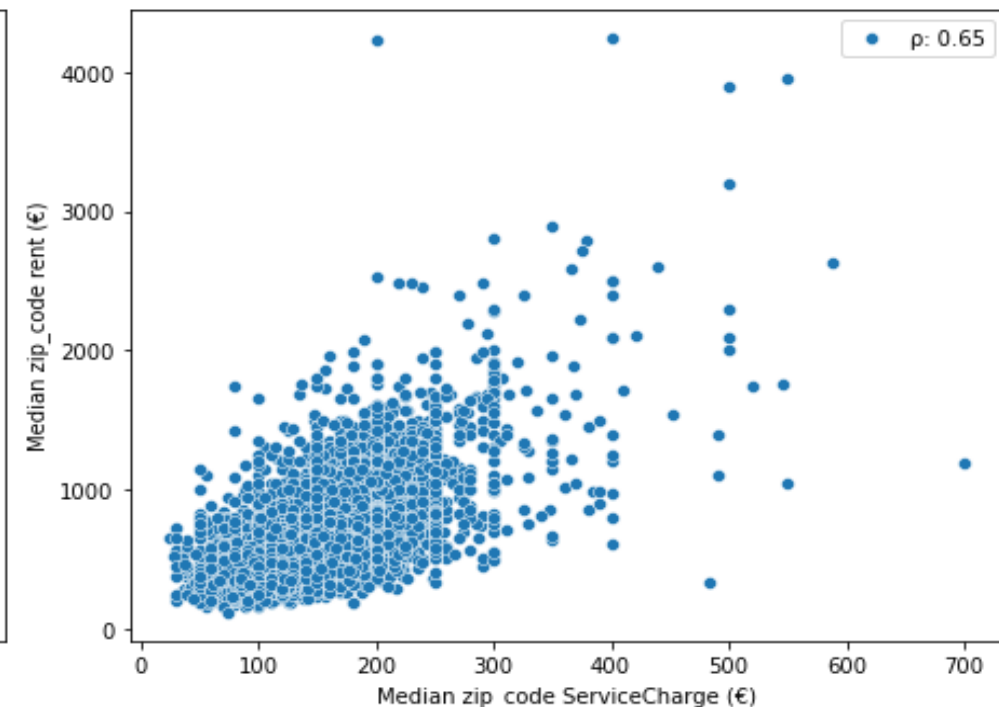
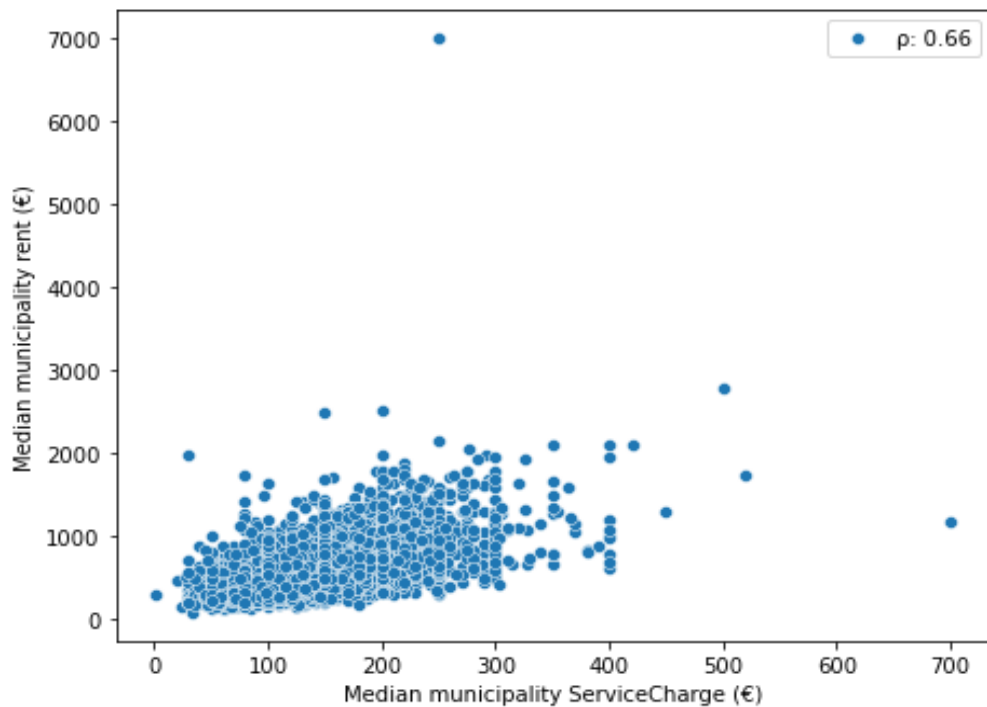
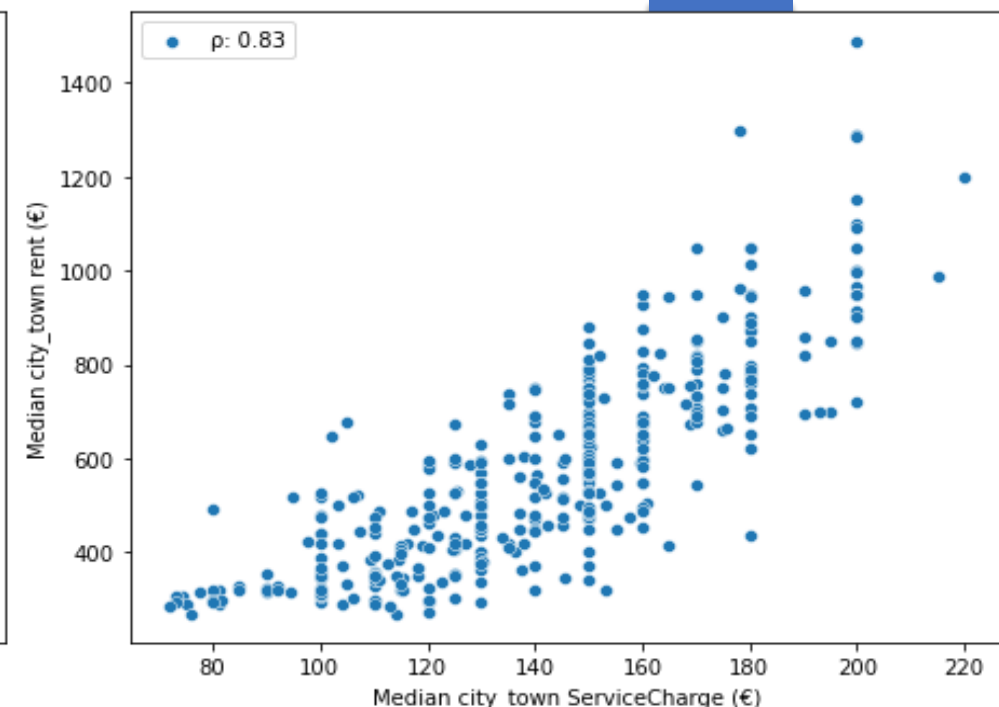
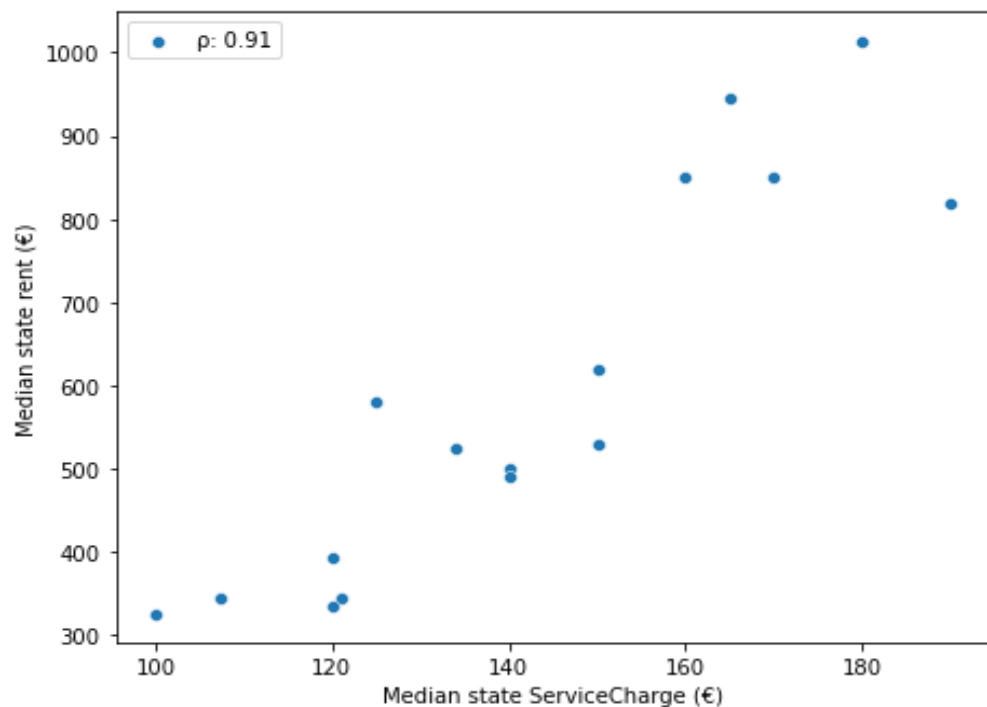
Plot of median rent versus median living space at the state, city / town, municipality, and zip code level



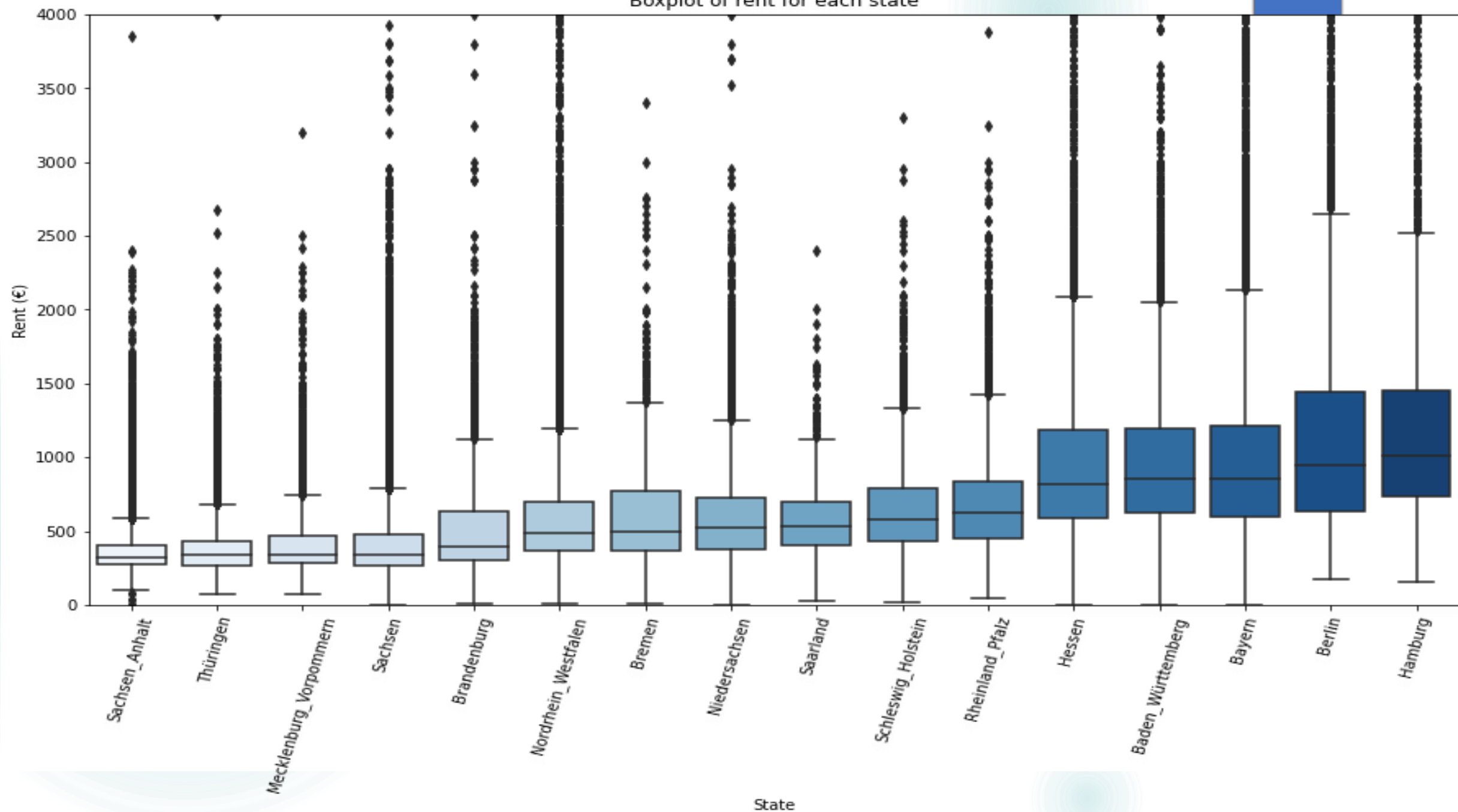
Relationship between rent and service charge for all listings



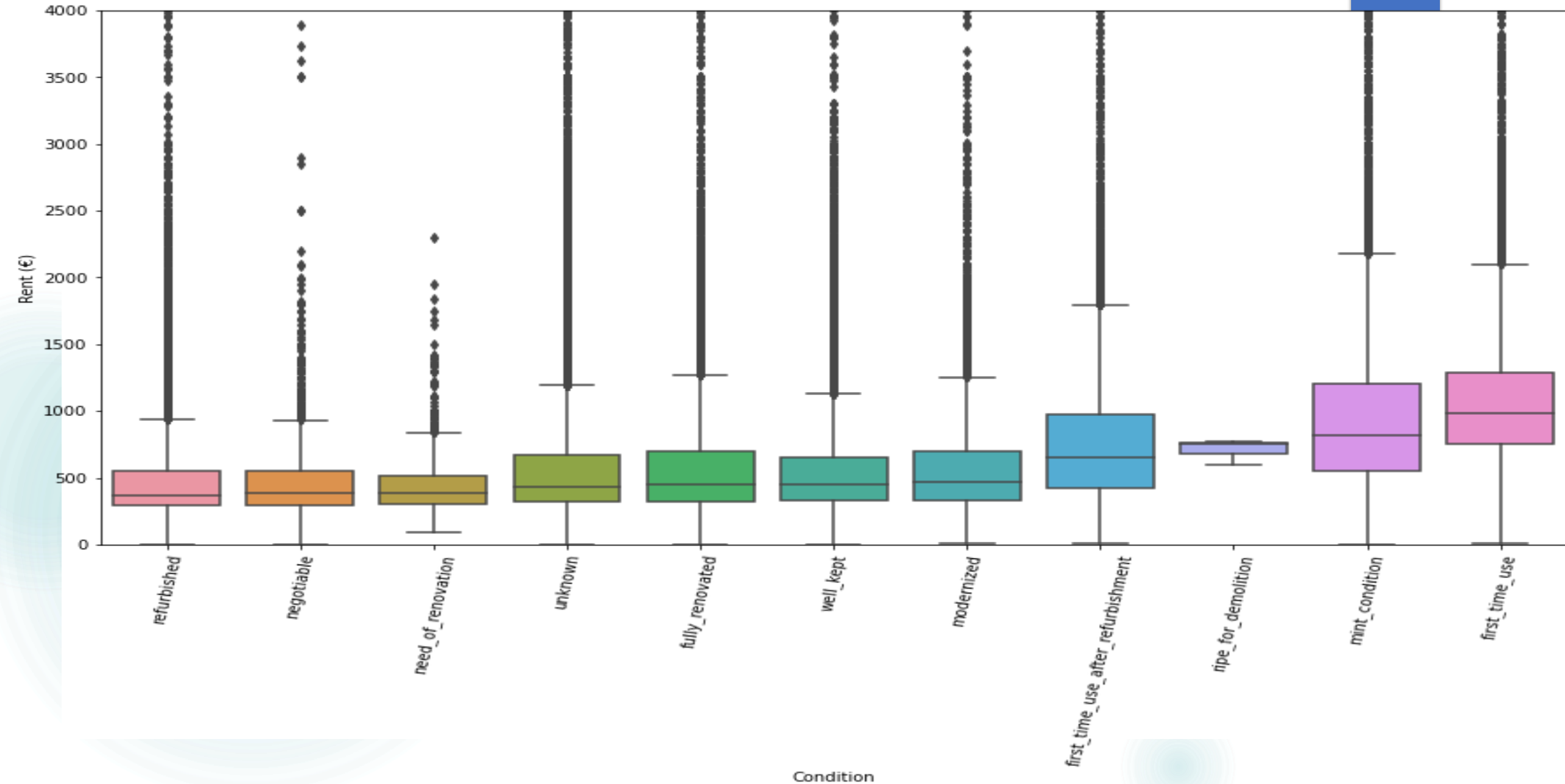
Plot of median rent versus median service charge at the state, city / town, municipality, and zip code level



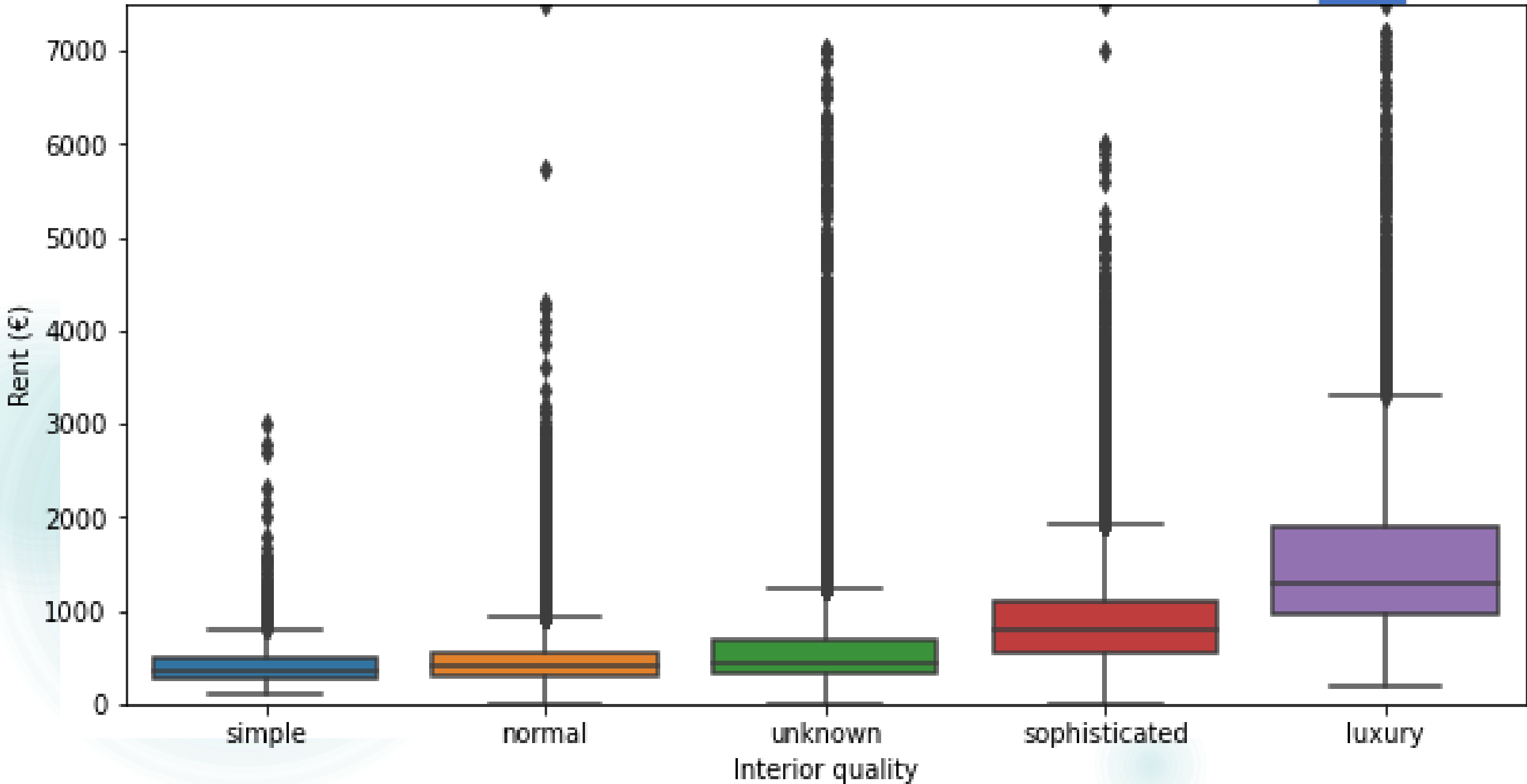
Boxplot of rent for each state



Relationship between rent and condition



Relationship between rent and quality



Machine Learning

TRAINING DATA DEVELOPMENT | METRICS | TESTING | MODEL
DEVELOPMENT | SELECTION | APPLICATION

Training data development

20

Application-set:

- 5 random samples

Training-set:

- 70% of remaining observations

Test-set:

- 30% of remaining observations

R-squared
score
(R²)

$$1 - \frac{\text{Residual sum of square errors}}{\text{Total sum of square errors}}$$

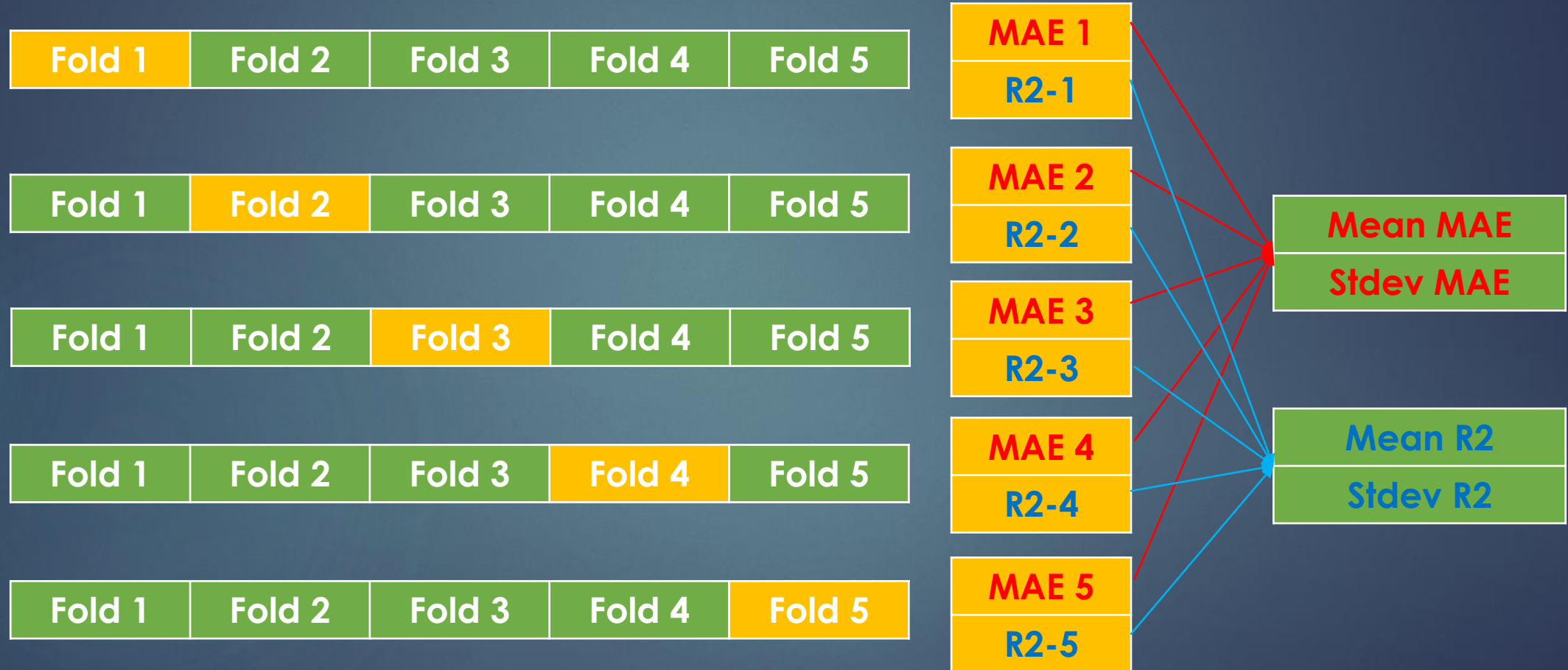
Mean
Absolute
Error
(MAE)

$$\sum \frac{|y - \hat{y}|}{N}$$

Testing strategy on training-set (70% of data)

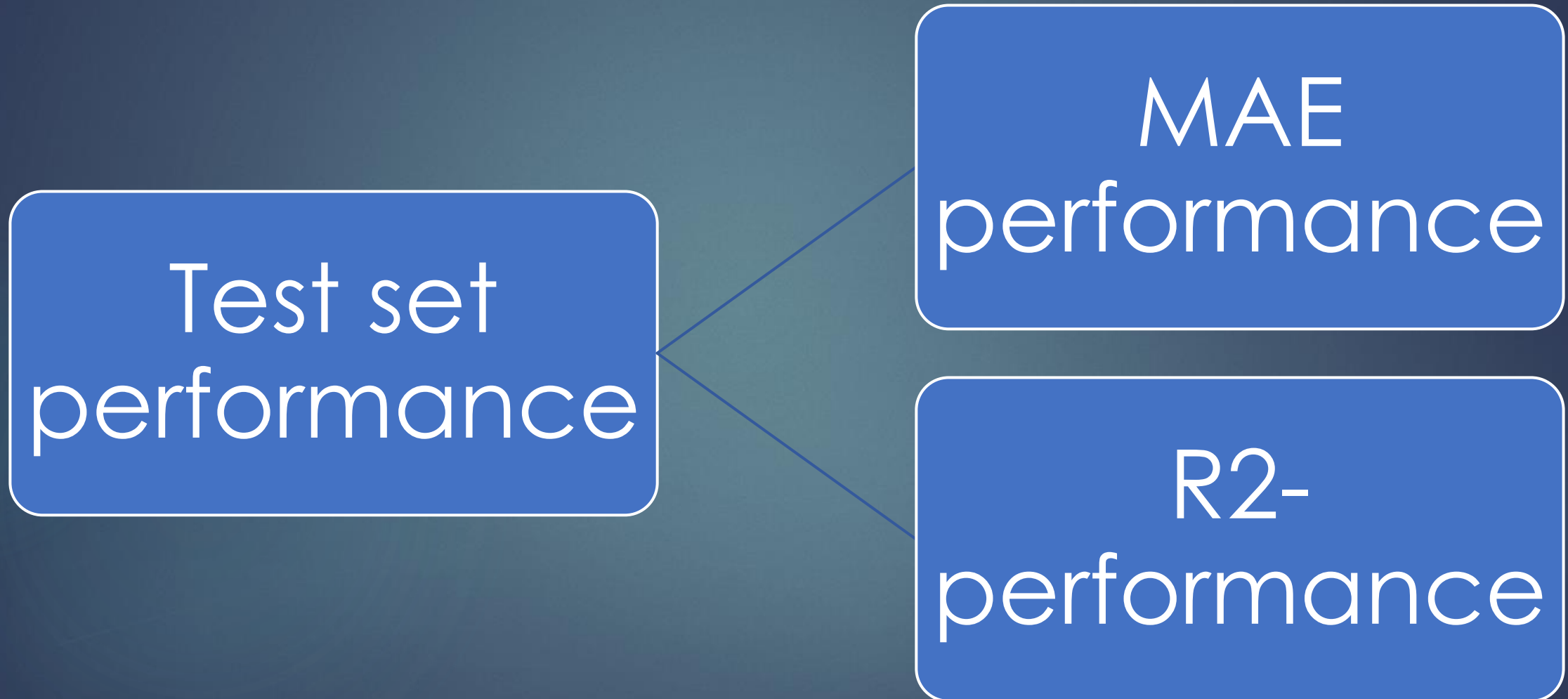
22

- Cross-validation on 70% training set (5 fold)



Testing strategy (on test-set (30% of data))

23



Baseline model – linear regression

	R ² score	Mae (€)
Train set	0.73	146.36
Test set	0.75	147.29

Linear

Lasso
regression

Ridge
regression

Tree induction

Random
forest

XGBoost

Feature selection

Used feature importance
for tree-induction model

Gini importance threshold
value of 0.001

Hyperparameter tuning

Random forest

- Number of estimators
- Maximum depth of tree

XGBoost

- Subsample
- Maximum depth of tree
- Column sample

Model Performance: r-squared

Model	Mean cv r2 score	Stdev of cv r2 scores	R2 test score
XGBoost	0.837	0.045	0.866
XGBoost w/ feat sel.	0.838	0.045	0.863
Random forest w/ hyper	0.819	0.047	0.861
XGBoost w/ hyper	0.827	0.047	0.849
Random forest w/ feat sel.	0.809	0.045	0.848
Random forest	0.807	0.049	0.841
Ridge	0.730	0.048	0.748
Lasso	0.730	0.048	0.747

Model performance: mean absolute error

28

Model	Mean cv mae score (€)	Stdev of cv mae scores (€)	mae test score (€)
Random forest w/ hyper	91.34	0.67	90.34
XGBoost	91.90	0.31	92.04
XGBoost w/ feat sel.	91.79	0.52	92.14
Random forest w/ feat sel.	97.44	0.66	97.19
Random forest	97.30	0.80	97.20
XGBoost w/ hyper	103.52	0.35	103.49
Lasso	145.63	0.81	146.49
Ridge	146.35	0.82	147.20

Common important features

29

Service
charge

Median city
picture
count

Median zip
code service
charge

Luxurious
interior
quality

Living space

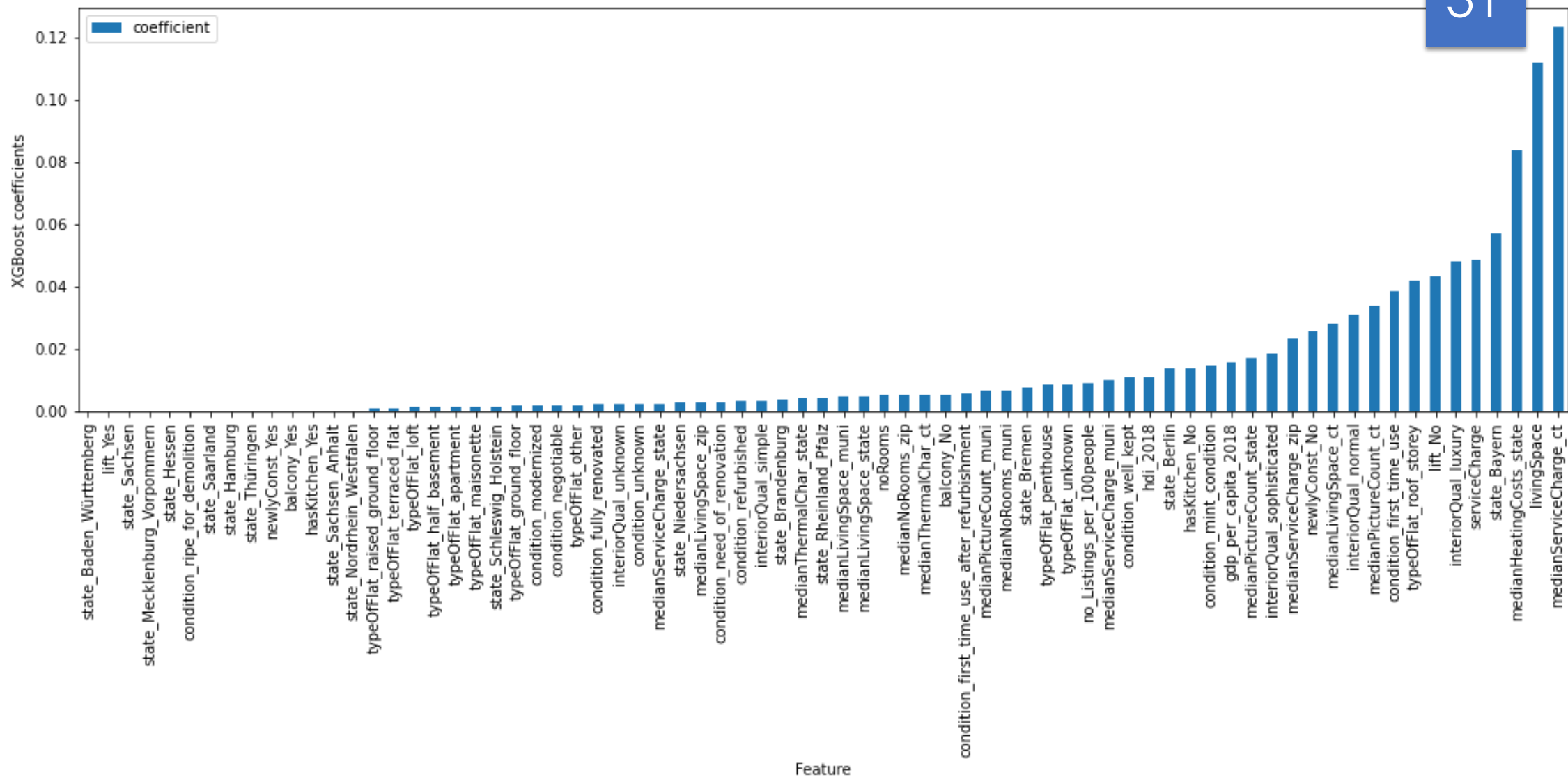
Median city
service
charge

Median city
living space

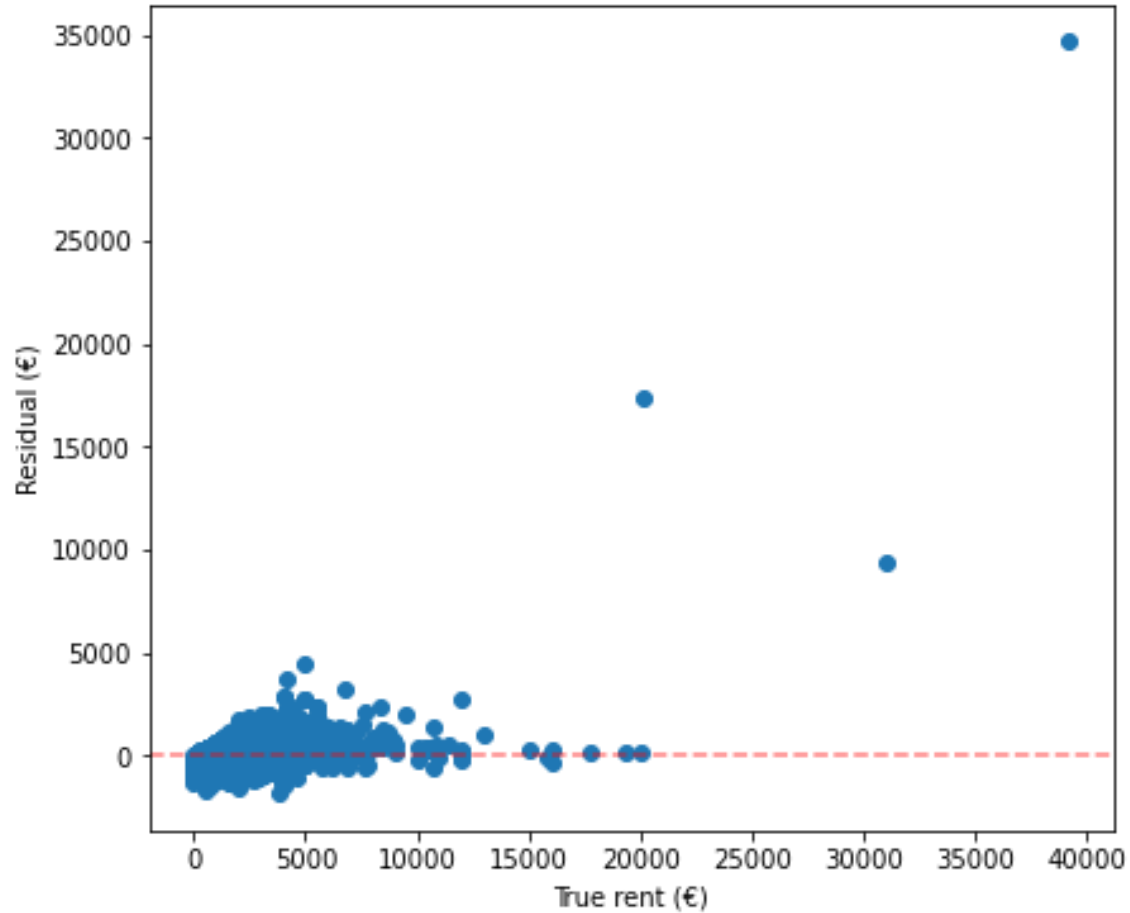
Model selection: XGBoost

Data	Mean cv r2 score	Stdev of cv r2 scores	Mean cv mae score (€)	Stdev of cv mae scores (€)
All	0.774	0.073	114.73	29.66
Training	0.837	0.045	91.90	0.31

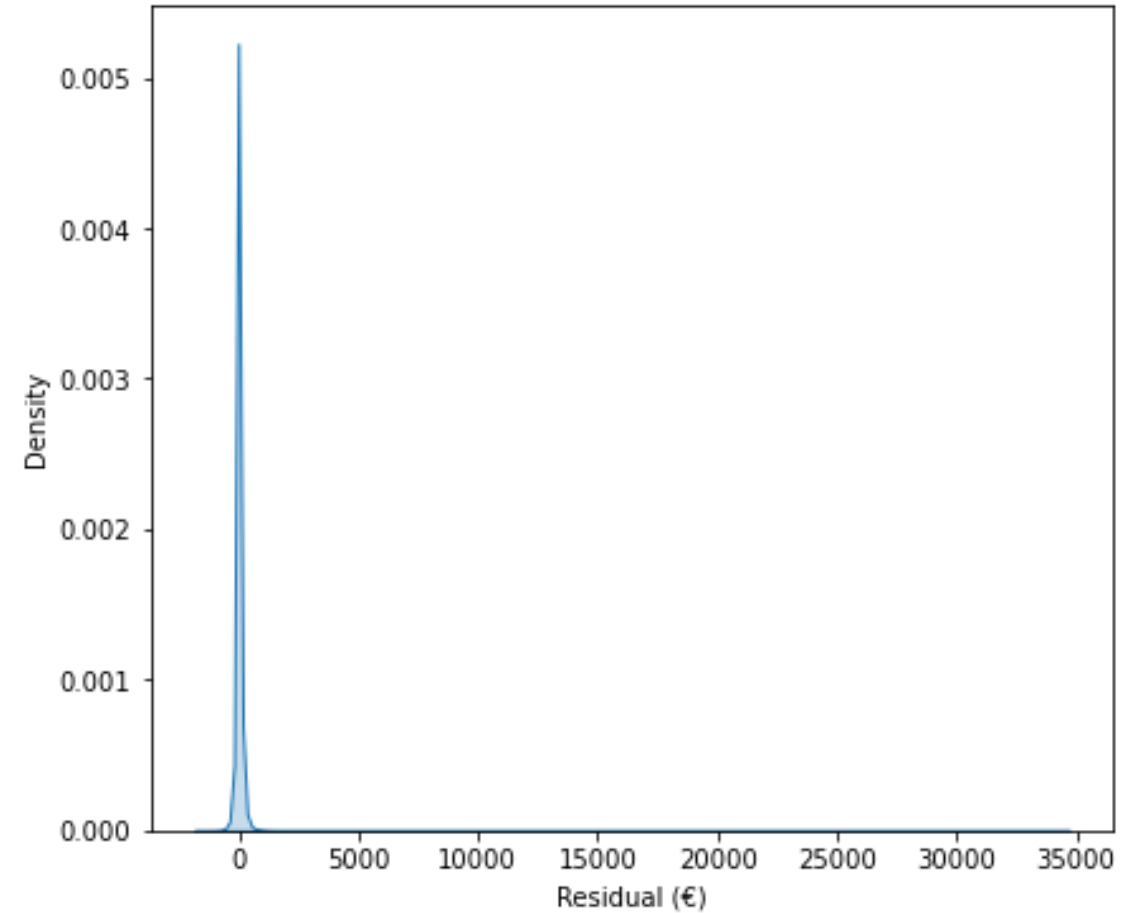
Plot of feature coefficient versus feature



Plot of true rents vs residuals



Histogram of residuals



Mean residual: €0 | Median residual: - €5 | Stdev: €157.5 | Skew: 46.5

Model application results

33

#	State	City/Town	Condition	rooms	Area (sqm)	Predicted rent (€)	Lower limit (€)	Set rent (€)	Higher limit (€)
1	Brandenburg	Uckermark	First time use	3.0	50.51	578.25	493.18	300.00	722.65
2	Schleswig Holstein	Dithmarschen	First time use after refurbishment	3.0	112.24	910.22	825.14	589.26	1054.61
3	Sachsen	Chemnitz	refurbished	2.0	51.67	641.11	556.03	518.00	785.50
4	Sachsen Anhalt	Halle Saale	negotiable	1.0	28.75	328.69	243.61	285.00	473.08
5	Hessen	Main Kinzig	Well kept	2.5	50.00	429.64	344.56	320.00	574.03

Conclusion / Recommendation

34

Determined most important factors for rent

- Living space
- Service charge
- Interior quality – luxurious

Predicted rent with XGBoost model

- R-squared: 0.77
- Mean absolute error: €114.73

Assumptions/Limitations/Opportunities

35

Analysis of text features (i.e. description and facilities)

Including time features in model

Uncertainty around sampling methodology

Higher compute capabilities for hyperparameter tuning

Questions

36

A 3D perspective view of a transparent rectangular box, possibly made of glass or acrylic, floating in a dark blue gradient background. The box is tilted, showing its top and side surfaces. In the center of the top surface, there is a large, white, stylized question mark.

?