

German Apartment Rentals

Magnus Paal and Sander Nemvalts

The data

- A recently published Kaggle dataset
- Rental apartments in Germany from ImmobilienScout24
- 268 000 rows/listings
- 49 columns of varying amounts of data.
- Latitude and longitude added with postprocessing

name	
regio1	condition
serviceCharge	interiorQual
heatingType	petsAllowed
telekomTvOffer	street
telekomHybridUploadSpeed	streetPlain
newlyConst	lift
balcony	baseRentRange
picturecount	typeOfFlat
pricetrend	geo_plz
telekomUploadSpeed	noRooms
totalRent	thermalChar
yearConstructed	floor
scoutId	numberOfFloors
noParkSpaces	garden
firingTypes	livingSpaceRange
hasKitchen	regio2
geo_bln	regio3
cellar	description
yearConstructedRange	facilities
baseRent	heatingCosts
houseNumber	energyEfficiencyClass
livingSpace	lastRefurbished
geo_krs	electricityBasePrice
	electricityKwhPrice
	date

What problem needed to be solved?

- Can we predict the rental price?
 - Help both parties know if the price is fair for a given property.

- Can we predict the heating costs?
 - Help the renter understand the not-so-obvious costs of renting.
 - More advanced: Heating type, heating costs -> CO₂ intensity of apartment heating.

Predicting rent

- Regression (predict total rent)
 - Explored various models
 - Random forests, neural networks, linear regression, more general gradient boosting, XGBoost
 - Best results from XGBoost regressor
 - Mean rent of dataset (y) - ~802
 - Standard deviation of y - ~437
 - Mean absolute error on test set - ~87
 - R^2 of ~0.9
 - Exceeded expectations
- Classification (predict price class/quantile)
 - 4 equally sized classes based on total rent
 - Best results from random forest classifier
 - Accuracy of ~77%, random baseline is 25%
 - Results did not match expectations

Predicting heating and CO₂ output from heating

- Also regression (predicting monthly heating costs).
- Was confident in XGBoost, stuck only with it.
- Best results:
 - XGBRegressor
 - Dataset mean heating costs was 69, standard deviation 22
 - Mean absolute error of ~11 on test set, pretty good.
 - R² of ~0.5 on test set. Would've liked better results.
- CO₂ prediction
 - Very simple logic, without ML
 - Calculate CO₂ emitted from heating type, €/kWh and kgCO₂/kWh.
 - Calculations seem close to reality (1-3 tons of CO₂ per property per year)
 - Also created a model that takes in same data as heating model and predicts CO₂ emitted
 - R² of ~0.9

What we learned

- Preprocessing the data turned out to be the most important step
 - We were not getting the results we wanted
 - Discovered that we had made a mistake preprocessing the data
 - Didn't approach cleaning with enough methodology, there were some considerable outliers left when we started training the models.
 - Cleaned data even more
 - Fixed mistakes made in cleaning
 - Results got much better

GitHub Repository

<https://github.com/magnuspaal/germany-rental-ml>