

Assignment 3

General information

Please give short (2-3 sentences) interpretations / explanations to your answers, not only the program code and outputs. Be concise and focused (less could be more ;)).

Grades will be distributed with the following rule: from the points you earn, you get 100% if you submit until the due date (**2025-03-23 21:00 CET**), 50% within 24 hours past due date, and 0% after that.

Predict real estate value

In this exercise you will predict property prices in New Taipei City, Taiwan, using [this dataset](#). (I have uploaded the data to the repo for you with cleaned up variable names. You can find it in the `real_estate` folder, [here](#).) Let's say you want to build a simple web app where potential buyers and sellers could rate their homes, and the provided `.csv` contains the data you have collected.

Similarly to what we did in the class, let's just work with a 20% subsample of the original data first. Put aside 30% of that sample for the test set. (*Hint: Extend the snippet below.*)

```
In [ ]: import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split

prng = np.random.RandomState(20250317)

url_data_on_github = {#TODO}
real_estate_data = pd.read_csv(url_data_on_github)
real_estate_sample = real_estate_data.sample(frac=0.2)

outcome = real_estate_sample['house_price_of_unit_area']
features = {#TODO}
X_train, X_test, y_train, y_test = train_test_split(
    features, outcome, test_size=0.3, random_state=prng
)

print(f"Size of the training set: {#TODO}, size of the test set: {#TODO}")
```

Tasks

1. Think about an appropriate loss function you can use to evaluate your predictive models. What is the risk (from a business perspective) that you would have to take by making a wrong prediction?
2. Build a simple benchmark model and evaluate its performance on the hold-out set (using your chosen loss function).
3. Build a simple linear regression model using a chosen feature and evaluate its performance. Would you launch your evaluator web app using this model?
4. Build a multivariate linear model with all the meaningful variables available. Did it improve the predictive power?
5. Try to make your model (even) better using the following approaches:
 - A. Feature engineering: e.g. including squares and interactions or making sense of latitude&longitude by calculating the distance from the city center, etc.
 - B. Training more flexible models: e.g. random forest or gradient boosting
6. Rerun three of your previous models (including both flexible and less flexible ones) on the full train set. Ensure that your test result remains comparable by keeping that dataset intact. (*Hint*: extend the code snippet below.) Did it improve the predictive power of your models? Where do you observe the biggest improvement? Would you launch your web app now?

```
In [ ]: real_estate_full = real_estate_data.loc[~real_estate_data.index.isin(X_test.index)]
print(f"Size of the full training set: {#TODO}")
```