

# DAT405 Introduction to Data Science and AI

2020-2021, Reading Period 1

## Assignment 2: Regression and classification

1. The following web page lists the selling prices of villas in Landvetter that were sold in the past 6 months. Find a linear regression model that relates the living area to the selling price. (You may transcribe the values from the web page into your program or into a data file by hand, or you can write a program to do this, but don't spend too much time doing this because "web scraping" is not the main objective of this assignment!)

[https://www.hemnet.se/salda/bostader?location\\_ids%5B%5D=940808&item\\_types%5B%5D=villa&sold\\_age=6m](https://www.hemnet.se/salda/bostader?location_ids%5B%5D=940808&item_types%5B%5D=villa&sold_age=6m)

- a. What are the values of the slope and intercept of the regression line?
  - b. Use this model to predict the selling prices of houses which have living area 100 m<sup>2</sup>, 150 m<sup>2</sup> and 200 m<sup>2</sup>.
  - c. Draw a residual plot.
  - d. Discuss the results, and how the model could be improved.
2. Use a confusion matrix to evaluate the use of logistic regression to classify the iris data set. Use the one-vs-rest option to use the same setup as in the lectures for multiclass regression.
  3. Use k-nearest neighbours to classify the iris data set with some different values for k, and with uniform and distance-based weights. What will happen when k grows larger for the different cases? Why?
  4. Compare the classification models for the iris data set that are generated by k-nearest neighbours (for the different settings from question 3) and by logistic regression. Calculate confusion matrices for these models and discuss the performance of the various models.
  5. Explain why it is important to use a separate test (and sometimes validation) set.

### Submitting work

In each file that you submit, give the names of the people submitting the work. On the first page of the report state how many hours each person spent working on the assignment.

If you upload a zip file, please also upload any PDF files separately (so that they can be viewed more conveniently in Canvas).

Deadline: Monday 14 September 2020 at 23:59.

### Self-check

Is all the required information on the front page? Have you answered all questions to the best of your ability? Anything else you can easily check? (details, terminology, arguments, clearly stated answers etc.?)

Do not submit an incomplete assignment! We are available to help you, and you can receive a short extension if you contact us.

Grading will be based on a qualitative assessment of each assignment. It is important to:

- i. Present clear arguments
- ii. Present the results in a pedagogical way
  - i. Should it be table/plot? What kind of plot? Is everything clear and easy to understand?
- iii. Show understanding of the topics
- iv. Give correct solutions.
- v. Make sure that the code is well commented.
  - i. Important parts of the code should be included in the running text and the full code uploaded to Canvas.