

- Unlike the rest of the module coursework you must do this assignment entirely yourself - you must not discuss or collaborate on the assignment with other students in any way, you must write answers in your own words and write code entirely yourself. If you use any online or other external content in your report you should take care to cite the source. It is mandatory to complete the declaration that the work is entirely your own and you have not collaborated with anyone - the declaration form is available on Blackboard. All submissions will be checked for plagiarism.
- Reports must be typed and submitted as a separate pdf on Blackboard (not as part of a zip file).
- Include the source of code written for the assignment as an appendix in your submitted pdf report. Also include a separate zip file containing the executable code and any data files needed. Programs should be running code written in Python, and should load data etc when run so that we can unzip your submission and just directly run it to check that it works. Keep code brief and clean with meaningful variable names etc.
- Important: For each problem, your primary aim is to articulate that you understand what you're doing - not just running a program and quoting numbers it outputs. Generally most of the credit is given for the explanation/analysis as opposed to the code/numerical answer.
- If you use machine learning models not covered in the course then you must take care to show that you understand them and are not just running code in a "black box" fashion (so explain how predictions are generated from an input, what the cost function is, what the model parameters and hyperparameters are and how they affect the predictions etc).
- Reports should typically be about 5 pages, with 10 pages the upper limit (excluding appendix with code).

## Downloading Dataset

- Go to the Dublin Bikes Open Data web site at <https://data.gov.ie/dataset/dublinbikes-api>. Download the "Dublinbikes 2020 Q1 usage data". This should be a large-ish (250MB) csv file. The file contains the number of bikes at each of the bike stations in Dublin, updated every 5 mins from 1st Jan to 24th Feb 2020. Pick two bike stations to study, choosing two stations with different patterns of behaviour e.g. one located in the city centre and one in the suburbs. In your submitted report please state which two stations you chose and why.

## Assignment

1. Write a short report evaluating the feasibility of predicting bike station occupancy 10mins, 30mins and 1 hour in the future. Appropriate feature selection is likely to be important so give this due attention. Select two machine learning approaches

(justify your choice), apply them to the dataset and critically evaluate their prediction performance. Remember it's v important to clearly explain/justify any design choices that you make and any conclusions you arrive at. Include any code you use in an appendix. [75 marks: indicative breakdown (i) feature engineering 20 marks, (ii) machine learning methodology 20 marks, (iii) evaluation 25 marks, (iv) report presentation 10 marks]

2.
  - (i) What is a ROC curve. How can it be used to evaluate the performance of a classifier. [5 marks]
  - (ii) Give two examples of situations when a linear regression would give inaccurate predictions. Explain your reasoning. [5 marks]
  - (iii) Discuss three pros/cons of an SVM classifier vs a neural net classifier. [5 marks]
  - (iv) Describe the operation of a convolutional layer in a convNet. Give a small example to illustrate. [5 marks]
  - v) In k-fold cross-validation a dataset is resampled multiple times. What is the idea behind this resampling i.e. why does resampling allow us to evaluate the generalisation performance of a machine learning model. Give a small example to illustrate. [5 marks]