

Lab 5

(Deadline 18:00 10/01/2020)

Task

You will be provided with a machine learning benchmark dataset (details below). The task focuses on the implementation and critical analysis of multiple regression methods: You should implement, evaluate, and analyse each of the following algorithms:

- Nearest neighbour
- Linear regression
- Regression forest
- Gaussian process

The results should be presented as a report with a 3000 word limit in a **pdf document**. Please also include your code as well (Jupyter workbook and/or normal .py files).

Mark scheme

This project is worth 60% of your marks for the unit:

- 10 marks for the method,
- 20 marks for designing and validating on a toy problem (see below),
- 20 marks for the experiments,
- 10 marks for the analysis,

for a total of 60 marks.

If a piece of work is submitted after the submission date (and no extension has been explicitly granted by the Director of Studies), the maximum possible mark will be 40% of the full mark. If work is submitted more than five working days after the submission date, the student will receive zero marks.

Data set

Included on Moodle is the SARCOS data set, a regression problem where the task is to predict the torque of one motor of a robotic arm given physical joint details. Specifically, position, velocity and acceleration for 7 degrees of freedom.

The data set is provided as one csv file - you are responsible for splitting it for train/test and hyperparameter learning. Each row is an exemplar, and each column a feature. Your task is to predict the last column (#22) given the first 21 columns. Be aware that you will not want to use all of the exemplars when training a Gaussian process as it will get too slow.

Suggested document structure

Please structure your report sensibly. Below is a suggested structure which we would strongly suggest you follow.

Section 1: Introduction

Explain the problem.

Section 2: Methods

Explain regression forests and Gaussian process regression. Include a discussion of why (or why not) you would use these algorithms. Note that you do not need to discuss nearest neighbours or linear regression.

Section 3: Validation on a toy problem

Design a toy problem and provide a discussion on validating your implementations on this problem. Please explain the rationale of your toy problem design and how it does, or does not, demonstrate that you have implemented your algorithm correctly.

A toy problem is used to debug, test and validate a machine learning algorithm; it will usually have the following properties:

- Low dimensional (2 or 3 dimensions), so it can be visualised.
- Synthetic, so the correct answer is known and there is no noise.
- Designed so that the behaviour of the algorithms being tested is predictable, so you will know if something is wrong.
- Very simple, e.g. you could type it into a calculator in a few seconds.

Note that a single toy problem can cover all four algorithms.

Section 4: Experiments and analysis

Please report your experimental results. This could include

- 1) Your hyperparameter selection strategies.
- 2) Performance comparison of different algorithms in terms of accuracy, computational complexity, etc.
- 3) Analysis: How the results are influenced by different hyper-parameter choices? Why does one algorithm perform better than the others on the given dataset?

Source

The data set was originally obtained from <http://gaussianprocess.org/gpml/data/>