

Assignment 1: k -nearest neighbours and linear regression

CS480/680 – Spring 2019

Out: May 13, 2019

Due: May 24 (11:59 pm), 2019.

Submit an electronic copy of your assignment via LEARN. Late submissions incur a 2% penalty for every rounded up hour past the deadline. For example, an assignment submitted 5 hours and 15 min late will receive a penalty of $\text{ceiling}(5.25) * 2\% = 12\%$.

Be sure to include your name and student number with your assignment.

1. **[30 pts]** Classification. Implement k -nearest neighbours using any programming language of your choice. Do not use any library such as scikit-learn that already has k -nearest neighbour or cross validation implemented. Implementing k -nearest neighbour and cross validation from scratch will be a good exercise to make sure that you fully understand those algorithms. Feel free to use general libraries for array and matrix operations such as numpy. Feel free to verify the correctness of your implementation with existing libraries such as scikit-learn.

Download the dataset posted on the course web page. Classify each input x according to the most frequent class amongst its k nearest neighbours as measured by the Euclidean distance (L2-norm). Break ties at random. Determine the best number of neighbours k by 10-fold cross validation.

What to hand in:

- Your code for k -nearest neighbours and cross validation.
- Find the best k by 10-fold cross validation. Draw a graph that shows the cross validation accuracy as k increases from 1 to 30.
- Report the best number of neighbours k and the accuracy on the test set of k -nearest neighbours with the best k .

2. **[30 pts]** Regression. Using any programming language, implement linear least square regression with the penalty term $0.5\lambda w^T w$. Do not use any library such as scikit-learn that already has linear regression or cross validation implemented. Implementing linear regression and cross validation from scratch will be a good exercise to make sure that you fully understand those algorithms. Feel free to use general libraries for array and matrix operations such as numpy. Feel free to verify the correctness of your implementation with existing libraries such as scikit-learn.

Download the dataset posted on the course web page. The output space is continuous (i.e., $y \in \mathbb{R}$). Determine the best λ by 10-fold cross validation.

What to hand in:

- Your code for linear regression and cross validation.
- Find the best λ by 10-fold cross validation. Draw a graph that shows the cross validation accuracy as λ increases from 0 to 4 in increments of 0.1.
- Report the best λ and the accuracy of linear regression on the test set for this best λ .

3. **[40 pts]** Theory. In class, we discussed several loss functions for linear regression. However all the loss functions that we discussed assume that the error contributed by each data point have the same importance. Consider a scenario where we would like to give more weight to some data points. Our goal is to fit the data points (x_n, y_n) in proportion to their weights r_n by minimizing the following objective:

$$L(w, b) = \sum_{n=1}^m r_n (y_n - wx_n + b)^2$$

where w and b are the model parameters, the training data pairs are (x_n, y_n) . To simplify things, feel free to consider 1D data (i.e., x_n and w are scalars).

- (a) **[20 pts]** Derive a closed-form expression for the estimates of w and b that minimize the objective. Show the steps along the way, not just the final estimates.
- (b) **[20 pts]** Show that this objective is equivalent to the negative log-likelihood for linear regression where each data point may have a different Gaussian measurement noise. What is the variance of each measurement noise in this model?