

## Homework Assignment # 2

Due: Sunday, October 16, 2022, 11:59 p.m. Mountain time

Total marks: 100

### Question 1. [15 MARKS]

Your goal in this question is to find the closed-form solution for the following constrained optimization problem, for fixed non-negative coefficients  $c_1, \dots, c_m \geq 0$ .

$$\min_{\mathbf{w} \in [0,1]^m, \sum_{k=1}^m w_k = 1} \sum_{k=1}^m c_k \ln w_k$$

#### (a) [5 MARKS]

What is the Lagrangian  $L$  for this optimization? Be clear about what variables are given to the Lagrangian function.

#### (b) [10 MARKS]

Derive the closed-form solution to the above optimization problem. See Section 6.3, to see the general way to approach this problem. Show your work.

### Question 2. [10 MARKS]

Using prototype features, using kernels, provide a simple approach for non-linear data representation. In this question you will reason about some practical choices when using prototype representations. Later in this assignment you will implement kernel regression with a subset of prototypes.

#### (a) [5 MARKS]

As a practitioner using kernel regression, a key decision is the choice of kernel similarity. Go to [https://en.wikipedia.org/wiki/Kernel\\_method#Popular\\_kernels](https://en.wikipedia.org/wiki/Kernel_method#Popular_kernels) and read about one of the linked kernel functions other than RBFs. Provide a 2-3 sentence description of the kernel and what it is used for. **Do not copy and paste from Wikipedia**; explain the kernel in your own words.

#### (b) [5 MARKS]

You need to decide on how many prototypes  $m$  to select. In particular, your goal is to have good generalization performance. Imagine you pick  $m = 100$  prototypes out of  $n = 1000$  samples. First, do you expect your training error to be lower or higher if you add another 100 prototypes from the dataset? Explain your answer.

Second, do you expect your generalization error to be lower or higher if you add another 100 prototypes from the dataset? Explain your answer.

### Question 3. [50 MARKS]

Please complete the A2.jl notebook provided in the code.zip file. This file contains two datasets and the Pluto notebook for this assignment. You will need to do the following.

(a) [5 MARKS] Complete the parts of the code needed to create a kernel representation.

(b) [10 MARKS] Implement Lasso regression for a linear model.

(c) [5 MARKS] Fill in the necessary implementation details for mini-batch gradient descent using

the mean squared error (MSE) loss.

(d) [15 MARKS] Implement the three optimizers for linear regression with mini-batch gradient descent. These include *ConstantLR* which uses a constant stepsize and two adaptive stepsize approaches, *RMSProp*, and *LineSearch*.

(e) [5 MARKS] Implement the prototype selection strategy that uses Lasso regression.

(f) [10 MARKS] Implement (external) cross validation using Monte Carlo CV (random data splits).

#### Question 4. [25 MARKS]

Now you will run three experiments, in the same Pluto notebook, over various datasets using the algorithms you implemented in the rest of the notebook. The algorithms are compared using your implementation of external cross validation, and the results are visualized using two types of plots. Note that we do not use internal cross validation in this assignment, and instead simply fix the hyperparameters to reasonable values that we found.

(a) [10 MARKS]

The first experiment uses a synthetic dataset with correlated features and a linear model. The features are generated from a multivariate Gaussian distribution with covariance matrix built using  $\Sigma_{ij} = \rho^{|i-j|}$  for some  $\rho \in [0, 1.0)$ . The weights  $\beta$  are randomly generated from a Gaussian distribution and fixed. The targets are generated using  $y = \mathbf{x}^\top \beta + \epsilon$ , where the noise  $\epsilon$  is generated from a zero-mean Gaussian distribution with standard deviation  $\sigma = 3$ .

The goal of this experiment is to compare OLS, Ridge Regression, and Lasso on a dataset with highly correlated features. What conclusions do you draw with the default setting of  $\rho = 0.7$ ? What would you expect to happen as  $\rho$  decreases? As  $\rho$  increases? Please explain your reasoning. After you make your hypotheses, change the entry for  $\rho$  in the settings dictionary what do you observe.

(b) [10 MARKS]

The second experiment uses two datasets. The first is a synthetic dataset and the second is known as the Boston Housing dataset (see the link in the notebook for more details). The goal of the experiment is to compare the different data representations we implemented.

1. What conclusions can you draw? Specifically, what do you notice about the difference between the linear and kernel representations.
2. If you were to decide to use a prototype selection strategy in the future, which would you choose? Explain your reasoning. When might we expect the Lasso ( $\ell_1$ ) prototype selection strategy to perform better than random selection?

(c) [5 MARKS]

The final experiment is comparing the optimizers we implemented: *ConstantLR*, *RMSProp*, and *LineSearch*. To do this we use a subset of the Susy dataset, which is detailed in the link provided in the notebook. What conclusions can you draw about the three algorithms?

#### Homework policies:

Your assignment should be submitted on eClass as a single pdf document and a zip file containing: the code (a .jl file), a .html file of the pluto notebook with all the cells run. The answers must be written legibly and scanned or must be typed (e.g., Latex). All code should be turned in when you submit your assignment. This means submitting the completed Pluto notebook, where you took the Pluto notebook with `todos` and completed them with your implementation. You are not allowed to change any of the imports in the notebook.

Because assignments are more for learning, and less for evaluation, grading will be based on coarse bins. **The grading is atypical.** For grades between (1) 80-100, we round-up to 100; (2) 60-80, we round-up to 80; (3) 40-60, we round-up to 60; and (4) **0-40, we round down to 0.** The last bin is to discourage quickly throwing together some answers to get some marks. The goal for the assignments is to help you learn the material, and completing less than 50% of the assignment is ineffective for learning.

**We will not accept late assignments.** Plan for this and aim to submit at least a day early. If you know you will have a problem submitting by the deadline, due to a personal issue that arises, please contact the instructor as early as possible to make a plan. If you have an emergency that prevents submission near the deadline, please contact the instructor right away. Retroactive reasons for delays are much harder to deal with in a fair way.

All assignments are individual. All the sources used for the problem solution must be acknowledged, e.g. web sites, books, research papers, personal communication with people, etc. Academic honesty is taken seriously; for detailed information see the University of Alberta Code of Student Behaviour.

**Good luck!**