

# 6.790 Homework 5

Revision: 11/6/24 12:30PM

Please hand in your work via Gradescope via the link at <https://gradml.mit.edu/info/homeworks/>. If you were not added to the course automatically, please use Entry Code R7RGGX to add yourself to Gradescope. Make sure to assign the problems to the corresponding pages in your solution when submitting via Gradescope.

1. Latex is not required, but if you are hand-writing your solutions, please write clearly and carefully. You should include enough work to show how you derived your answers, but you don't have to give careful proofs.
2. Homework is due on Tuesday November 12 at 11PM.
3. Lateness and extension policies are described at [https://gradml.mit.edu/info/class\\_policy/](https://gradml.mit.edu/info/class_policy/).

## Contents

## Robust Binary Classification (30 points)

1. In this problem we will study **robust** binary classification. For any input  $x \in \mathbb{R}^d$ , the hypothesis function is

$$h_\theta(x) = w^\top x + b, \text{ where } \theta = \{w \in \mathbb{R}^d, b \in \mathbb{R}\}.$$

The class label  $y \in \{+1, -1\}$ . The prediction probability of class +1 is given as

$$p(y = +1|x, \theta) = \frac{1}{1 + \exp(-h_\theta(x))}$$

and  $p(y = +1|x, \theta) = 1 - p(y = -1|x, \theta)$ .

We will use the following loss function:

$$\ell(h_\theta(x), y) = -\log P(y|x, \theta) = \log(1 + \exp(-yh_\theta(x)))$$

- (a) (5 points) The loss  $\ell(h_\theta(x), y)$  can be written as  $L(z) = \log(1 + \exp(-z))$ , where  $z = yh_\theta(x)$ . What can we say about the monotonicity of  $L(z)$  with respect to  $z$ ?

The robust binary classification problem can now be defined as:

$$\min_{\theta} \max_{\delta: \|\delta\| \leq \epsilon} \ell(h_\theta(x + \delta), y) \quad (1)$$

where  $\|\cdot\|$  denotes a valid norm.

To understand how this formulation works, we need an intermediate step of understanding so called dual norms. They are defined as

$$\|\theta\|_* = \max_{\|\delta\| \leq 1} \delta^\top \theta \quad (2)$$

where, for example,  $\|\cdot\|$  could be the 2-norm  $\|\delta\|_2 = (\sum_{j=1}^d \delta_j^2)^{1/2}$  or the  $\ell$ -infinity norm  $\|\delta\|_\infty = \max_{j=1}^d |\delta_j|$ . As a warm-up exercise, let's solve what the corresponding dual norms are.

- (b) (5 points) Find what norms the dual norms  $\|\theta\|_* = \max_{\|\delta\|_2 \leq 1} \delta^\top \theta$  and  $\|\theta\|_* = \max_{\|\delta\|_\infty \leq 1} \delta^\top \theta$  are.
- (c) (5 points) The inner maximization problem in (1) i.e.  $\max_{\delta: \|\delta\| \leq \epsilon} L(yh_\theta(x + \delta))$ , using  $L(\cdot)$  can be written as:

$$\max_{\delta: \|\delta\| \leq \epsilon} L(y(w^\top(x + \delta) + b))$$

Find the closed form solution of this inner maximization problem when:

- 1)  $\|\cdot\|$  is  $\infty$  - norm.
- 2)  $\|\cdot\|$  is 2 - norm.

Hint: [Use the previous parts]

- (d) (5 points) Using the solution from part c, provide the simplified min-max problem. Is this problem same as solving the nominal classification problem with a regularizer term? If yes, then explain why is this the case, and if not, then explain the differences.

We now consider a dataset of  $(x, y)$  pairs with  $d = 1$ :

$$\mathcal{D} = \{(-1, -1), (-1, -1), (-1, -1), (-1, -1), (-1, -1), (-1, -1), \\ (-1, -1), (-1, -1), (-1, -1), (-1, -1), (-1, -1), (-1, -1), \\ (-1, -1), (-1, -1), (-1, -1), (-1, -1), (-1, -1), (1, -1), (2, 1)\} \quad (3)$$

- (e) (2 points) What parameter vector  $\theta$  minimizes the regular logistic regression objective? If more than one, explain. (You should do this without code).
- (f) (2 points) What parameter vector  $\theta$  minimizes l2-regularized logistic regression with  $\lambda = 10^{-12}$  (i.e. adding  $\lambda w^2$  to the objective)? (This question and the rest of the parts of this question can be done using code.)
- (g) (2 points) What class probability distribution does logistic regression with these parameters assign to the point  $x = 1.51$ ?
- (h) (2 points) What parameter vector minimizes the robust regression objective, if we want to be correct even if the input is perturbed by  $\epsilon = 0.48$ ?
- (i) (2 points) What class probability distribution does logistic regression with these parameters assign to the point  $x = 1.51$ ?

## Learning Theory (10 points)

2. You work for Googolog and plan to train a lot of neural networks!

- You are going to try 10 different architectures
- Each with 10 different learning-rate schedules
- And measure the performance at 100 epochs (each) to avoid overtraining

You have a big pile of data. You'd like to use as much as possible for training, but you know you need to save a held-out set for validation in order to pick the best architecture, learning-rate schedule and number of training epochs.

You would like to be able to predict the risk on unseen points by computing the risk on the validation set, and guarantee that, with probability  $> 0.99$  it is within 0.01 of the true risk. Suppose the risk on any point is bounded in  $[0, 1]$ .

- (a) (8 points) How big does your validation set need to be?
- (b) (2 points) Your friend says that if you increase the training set size, you'll be able to make this bound tighter. Is your friend right or wrong? Explain.

## Causality (10 points)

3. One idea from learning causal models (which we may explore later in the class) is that the relationships in the "causal" direction are more generally useful and transferrable, than the other way around. So, for example, if we are interested in the relationship between a disease and its symptoms (which are caused by the disease), we might prefer to model

$$P(\text{symptoms} \mid \text{disease}) \quad (10)$$

rather than

$$P(\text{disease} \mid \text{symptoms}) \quad (11)$$

But, of course, we are usually called upon to predict the disease from the symptoms, and make a classifier that produces  $P(\text{disease} \mid \text{symptoms})$ .

So, for example, if we trained a classifier for malaria based on data in a tropical country, we might end up with  $P_{\text{tropical}}(\text{malaria} \mid \text{symptoms})$ .

- (a) (8 points) It doesn't seem good to use this same classifier in the US. If we know the incidence of malaria in our tropical population is 1%, and the incidence in US is 0.001%, how could we make use of our  $P_{\text{tropical}}(\text{malaria} \mid \text{symptoms})$  at Mass General?
- (b) (2 points) If a patient came in with symptoms that would, using the tropical model, made us predict that they had malaria with probability 0.75, what prediction would this US-adjusted model make?

## Day-Night (15 points)

4. You are training your robot outside in the daylight and come up with a good classifier for whether the terrain in front of it is grass. Now, you get data from night-time and everything looks different! Luckily, there were some examples in your original training data with substantial shadows, and some gathered at dusk. Letting  $D_{\text{day}} = (x_i, y_i)_{i=1}^n$  be the first, labeled data set, and  $D_{\text{night}} = (x_i)_{i=n+1}^m$  be the night-time data that you are going to have to make predictions on, let's explore using importance reweighting to address this problem.

- (a) (2 points) We need to begin by training a classifier. What training data would we use (inputs and targets)?
- (b) (3 points) We would typically use logistic regression for this classifier. Would it work to use a decision tree?
- (c) (3 points) Given the trained classifier, explain what problem we need to solve next. Is there a way of using the weights we learned on the daytime-only data to predict on night-time data, or do we need to retrain? If not, how do we make predictions?
- (d) (5 points) Assuming we do retrain, explain a strategy for approximating the importance-reweighted objective, while using a standard pre-existing logistic regression package.
- (e) (2 points) Now let's consider a fake domain in which  $x \in \mathbb{R}^2$ , and we have

$$y = f(x_1) \quad \text{if } x_2 = 1 \quad (13)$$

$$y = g(x_1) \quad \text{otherwise} \quad (14)$$

Furthermore, assume that our labeled data set  $D_1$  only has  $x$  values in which  $x_2 = 1$ . Now we have a new set of *unlabeled* examples,  $D_2$ , which only has  $x$  values in which  $x_2 = 0$ .

Explain what the importance re-weighting method would do in this case.

## Fairness (20 points)

5. Read the linked Zafar et al. paper, Sections 1–4: <https://www.jmlr.org/papers/volume20/18-262/18-262.pdf>. It's okay to skip parts about SVMs

In the following, we will focus on Overall Misclassification Rate (Equation 3.3)

- (a) (1 points) When is  $g_\theta(y, x)$  non-zero? (it's defined after Equation 4.4)
- (b) (5 points) The paper says (just after Equation 4.4) “if a decision boundary satisfies Eq. (3.3), then the (empirical) covariance defined above will be (approximately) zero (for a sufficiently large training set).”

Write a formal version of this statement, prove it, and state any additional assumptions that need to be made for this to hold.

- (c) (5 points) If we are using a classifier where the signed distance to decision boundary  $d_\theta(x) = \frac{\theta^T x}{\|\theta\|}$ , what is a value of  $c$  in Equation 4.5 that allows a gap of  $\epsilon$  in the error rate between the two groups ( $z = 1$  and  $z = 0$ )? Assume that the groups have equal proportions. Also, assume that all misclassifications are exactly distance 1 from the decision boundary (do not use any additional assumptions you made in part b).
- (d) (2 points) Why are there no  $z$  values in Equation 4.9?
- (e) (2 points) If we optimize Equation 4.13 to obtain classifier  $\theta$ , will we need  $z$  values at prediction time? Why would it be better if we did not?
- (f) (5 points) Suppose that after running extensive experiments, we find that if we train a classifier on a population of  $n_0$  examples with  $z = 0$  and  $n_1$  examples with  $z = 1$ , the classification error rate ( $P(y \neq \hat{y})$ ) on unseen points with  $z = 0$  will be  $\frac{c}{n_0}$  and the error rate on unseen points with  $z = 1$  will be  $\frac{c}{n_1}$ . We know that the underlying data distribution has  $z = 0$  and  $z = 1$  in equal proportion.

Unfortunately, we have lost our labels for which points are in group  $z = 0$  and  $z = 1$ . Suppose we wish to guarantee (with probability 99%) our classifier will be fair on test points up to some small error  $\epsilon$ :

$$|P(y \neq \hat{y} \mid z = 0) - P(y \neq \hat{y} \mid z = 1)| \leq \epsilon \quad (23)$$

Suppose we draw a training set of  $n$  pairs  $(x, y)$  to train our classifier. We can show that in order to guarantee the above property, the number of training points must scale as  $\epsilon^{-p}$ . What is  $p$ ?