

# CAP 6610 Machine Learning, Spring 2022

## Homework 2 Solution

1. (10 points) What is the distance between two parallel hyperplanes  $\{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{a}^\top \mathbf{x} = b_1\}$  and  $\{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{a}^\top \mathbf{x} = b_2\}$ ? *Hint.* Let  $\mathbf{a}^\top \mathbf{x}_1 = b_1$ ,  $\mathbf{a}^\top \mathbf{x}_2 = b_2$ , and minimize  $\|\mathbf{x}_1 - \mathbf{x}_2\|^2$ .

**Solution.** The distance between two sets is the smallest distance between two points from each sets. It can be formulated as the following optimization problem

$$\begin{aligned} & \underset{\mathbf{x}_1, \mathbf{x}_2}{\text{minimize}} && \|\mathbf{x}_1 - \mathbf{x}_2\|^2 \\ & \text{subject to} && \mathbf{a}^\top \mathbf{x}_1 = b_1, \mathbf{a}^\top \mathbf{x}_2 = b_2. \end{aligned} \tag{1}$$

The two constraints implies that

$$\mathbf{a}^\top (\mathbf{x}_1 - \mathbf{x}_2) = b_1 - b_2,$$

together with the Cauchy-Schwarz inequality

$$|\mathbf{a}^\top (\mathbf{x}_1 - \mathbf{x}_2)| \leq \|\mathbf{a}\| \|\mathbf{x}_1 - \mathbf{x}_2\|, \tag{2}$$

we see that

$$\|\mathbf{x}_1 - \mathbf{x}_2\| \geq \frac{|b_1 - b_2|}{\|\mathbf{a}\|},$$

for any  $\mathbf{x}_1$  and  $\mathbf{x}_2$  that satisfy  $\mathbf{a}^\top \mathbf{x}_1 = b_1$  and  $\mathbf{a}^\top \mathbf{x}_2 = b_2$ .

Furthermore, if we let

$$\mathbf{x}_1 = \mathbf{a} \frac{b_1}{\|\mathbf{a}\|^2}, \quad \mathbf{x}_2 = \mathbf{a} \frac{b_2}{\|\mathbf{a}\|^2}, \tag{3}$$

then

$$\|\mathbf{x}_1 - \mathbf{x}_2\| = \frac{|b_1 - b_2|}{\|\mathbf{a}\|},$$

which attains the lowerbound in (2) This means (3) is a solution to Problem (1), and the distance is

$$\frac{|b_1 - b_2|}{\|\mathbf{a}\|}$$

2. (10 points) Let  $x$  be a real-valued random variable with sample space  $\{a_1, \dots, a_k\}$  where  $a_1 \leq a_2 \leq \dots \leq a_k$ . This can be view as a categorical random variable with each category assigned a real value. Let  $\Pr[x = a_i] = p_i$ , then the vector  $\mathbf{p}$  satisfies  $\mathbf{p} \geq 0$  and  $\mathbf{1}^\top \mathbf{p} = 1$ , i.e., it lies in the probability simplex  $\Delta$ . For each of the following functions of  $\mathbf{p}$  on the probability simplex, determine if the function is convex, concave, or neither.

- (a)  $E[x]$
- (b)  $\Pr[x > \alpha]$
- (c)  $\Pr[\alpha < x < \beta]$
- (d)  $-\sum_{i=1}^k p_i \log p_i$ , the entropy of this distribution
- (e)  $\text{var}(x)$

**Solution.**

(a)

$$\mathbb{E}[x] = \sum_{i=1}^k a_i p_i = \mathbf{a}^\top \mathbf{p},$$

where  $\mathbf{a} = (a_1, \dots, a_k)$ . This is a linear function of  $\mathbf{p}$ , and thus both convex and concave.

(b)

$$\Pr[x > \alpha] = \sum_{i: a_i > \alpha} p_i = \mathbf{b}^\top \mathbf{p},$$

where  $\mathbf{b}$  is defined as

$$b_i = \begin{cases} 0 & a_i \leq \alpha \\ 1 & a_i > \alpha. \end{cases}$$

This is a linear function of  $\mathbf{p}$ , and thus both convex and concave.

(c)

$$\Pr[x > \alpha] = \sum_{i: \alpha < a_i < \beta} p_i = \mathbf{c}^\top \mathbf{p},$$

where  $\mathbf{c}$  is defined as

$$c_i = \begin{cases} 1 & \alpha < a_i < \beta \\ 0 & \text{otherwise.} \end{cases}$$

This is a linear function of  $\mathbf{p}$ , and thus both convex and concave.

(d)  $u \log u$  is a convex function of  $u$ , therefore  $-\sum_{i=1}^k p_i \log p_i$  is a concave function of  $\mathbf{p}$ .

(e) We have

$$\text{var}(x) = \mathbb{E}[x^2] - (\mathbb{E}[x])^2 = \sum_{i=1}^k a_i^2 p_i - (\mathbf{a}^\top \mathbf{p})^2.$$

This is a quadratic function with the Hessian matrix  $-\mathbf{a}\mathbf{a}^\top$ , which is negative semidefinite. Therefore it is a concave function of  $\mathbf{p}$ .

3. (10 points) *Log-concavity of Gaussian cumulative distribution function.* The cumulative distribution-function of a Gaussian random variable,

$$\Phi(t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-x^2/2} dx$$

is log-concave. This follows from the general result that the convolution of two log-concave functions is log-concave. In this problem we guide you through a simple self-contained proof that  $\Phi$  is log-concave. A useful fact is that  $\Phi$  is log-concave if and only if  $\Phi''(t)\Phi(t) \leq (\Phi'(t))^2$ .

(a) Verify that  $\Phi''(t)\Phi(t) \leq (\Phi'(t))^2$  for  $t \geq 0$ . That leaves us the hard part, which is to show the inequality for  $t < 0$ .

(b) Verify that for any  $t$  and  $x$  we have  $x^2/2 \geq -t^2/2 + tx$ .

(c) Using part (b) to show that  $e^{-x^2/2} \leq e^{t^2/2-tx}$ . Conclude that

$$\int_{-\infty}^t e^{-x^2/2} dx \leq e^{t^2/2} \int_{-\infty}^t e^{-tx} dx.$$

(d) Use part (c) to verify that  $\Phi''(t)\Phi(t) \leq (\Phi'(t))^2$  for  $t < 0$ .

Let us first verify the condition that  $\Phi$  is log-concave if and only if  $\Phi''(t)\Phi(t) \leq (\Phi'(t))^2$ . This comes from the definition that  $(\log \Phi(t))'' \leq 0$  everywhere if  $\Phi$  is log-concave. Applying the chain rule, we have

$$(\log \Phi(t))' = \frac{\Phi'(t)}{\Phi(t)}.$$

Its second derivative is

$$(\log \Phi(t))'' = \frac{\Phi''(t)}{\Phi(t)} - \frac{(\Phi'(t))^2}{(\Phi(t))^2}.$$

Letting it nonpositive gives the condition  $\Phi''(t)\Phi(t) \leq (\Phi'(t))^2$ .

**Solution.** The first and second derivative of  $\Phi$  are

$$\Phi'(t) = \frac{1}{\sqrt{2\pi}} e^{-t^2/2}, \quad \Phi''(t) = -\frac{t}{\sqrt{2\pi}} e^{-t^2/2}.$$

(a)  $\Phi''(t) \leq 0$  when  $t \geq 0$ .

(b) Since  $t^2/2$  is convex, we have

$$t^2/2 \geq x^2/2 + x(t-x) = tx - x^2/2.$$

This is the general inequality for any differentiable convex function  $f(t)$

$$f(t) \geq f(x) + f'(x)(t-x)$$

applied to  $f(t) = t^2/2$ .

(c) Rearrange and taking exponential on both sides gives

$$e^{-x^2/2} \leq e^{t^2/2-tx}.$$

Since it holds for any  $x$ , the inequality still holds if we sum over different values of  $x$ . Let  $x$  take any value between  $-\infty$  and  $t$ , the limit of the sum becomes the integral

$$\int_{-\infty}^t e^{-x^2/2} dx \leq e^{t^2/2} \int_{-\infty}^t e^{-tx} dx.$$

(d) We can evaluate the integral

$$\int_{-\infty}^t e^{-tx} dx = -\frac{1}{t} e^{-tx} \Big|_{-\infty}^t = -\frac{1}{t} e^{-t^2}.$$

Notice that we only consider  $t < 0$ , so  $\lim_{x \rightarrow -\infty} e^{-tx} = 0$ . Plugging it back to part (c) gives

$$\int_{-\infty}^t e^{-x^2/2} dx \leq -\frac{1}{t} e^{-t^2} e^{t^2/2} \implies -te^{-t^2/2} \int_{-\infty}^t e^{-x^2/2} dx \leq e^{-t^2}.$$

The inequality does not change direction because, again, we only consider  $t < 0$  here. Multiply both sides by  $1/(2\pi)$  shows  $\Phi''(t)\Phi(t) \leq (\Phi'(t))^2$ .

We established  $\Phi''(t)\Phi(t) \leq (\Phi'(t))^2$  for both  $t \geq 0$  in part (a) and  $t < 0$  in part (b)–(d). Hence  $\Phi$  is log-concave.

4. (10 points) Show that the following two convex problems are equivalent. Carefully explain how the solution of (b) is obtained from the solution of (a).

(a) The robust least squares problem

$$\underset{\boldsymbol{\theta}}{\text{minimize}} \quad \sum_{i=1}^n h(\boldsymbol{\phi}_i^\top \boldsymbol{\theta} - \psi_i),$$

where  $h : \mathbb{R} \rightarrow \mathbb{R}$  is the Huber function defined (with a constant  $M$ ) as

$$h(t) = \begin{cases} t^2 & |t| \leq M \\ M(2|t| - M) & |t| > M. \end{cases}$$

(b) The quadratic program

$$\begin{aligned} &\underset{\boldsymbol{\theta}, \mathbf{u}, \mathbf{v}}{\text{minimize}} \quad \sum_{i=1}^n (u_i^2 + 2Mv_i) \\ &\text{subject to} \quad -\mathbf{u} - \mathbf{v} \leq \boldsymbol{\Phi}\boldsymbol{\theta} - \boldsymbol{\psi} \leq \mathbf{u} + \mathbf{v} \\ &\quad \quad \quad 0 \leq \mathbf{u} \leq M\mathbf{1}, \quad \mathbf{v} \geq 0. \end{aligned}$$

**Solution.** Suppose we fix  $\boldsymbol{\theta}$  in Problem (b). First we notice that at the optimum of (b) we must have  $u_i + v_i = |\boldsymbol{\phi}_i^\top \boldsymbol{\theta} - \psi_i|$ , because otherwise we can further decrease the objective function without violating the constraints. Therefore  $v_i = |\boldsymbol{\phi}_i^\top \boldsymbol{\theta} - \psi_i| - u_i$  at optimum.

Eliminating  $\mathbf{v}$  yields the following problem

$$\begin{aligned} &\underset{\mathbf{u}}{\text{minimize}} \quad \sum_{i=1}^n (u_i^2 - 2Mu_i + 2M|\boldsymbol{\phi}_i^\top \boldsymbol{\theta} - \psi_i|) \\ &\text{subject to} \quad 0 \leq u_i \leq \min(M, |\boldsymbol{\phi}_i^\top \boldsymbol{\theta} - \psi_i|), \quad i = 1, \dots, n. \end{aligned}$$

The problem is separable over each  $u_i$ , and we rewrite it as

$$\begin{aligned} &\underset{u_i}{\text{minimize}} \quad (u_i - M)^2 - M^2 + 2M|\boldsymbol{\phi}_i^\top \boldsymbol{\theta} - \psi_i| \\ &\text{subject to} \quad 0 \leq u_i \leq \min(M, |\boldsymbol{\phi}_i^\top \boldsymbol{\theta} - \psi_i|). \end{aligned}$$

It is easy to see that if  $M < |\boldsymbol{\phi}_i^\top \boldsymbol{\theta} - \psi_i|$ , then we should choose  $u_i = M$  to minimize the objective; otherwise we should choose  $u_i$  to be as close to  $M$  as possible, which is  $|\boldsymbol{\phi}_i^\top \boldsymbol{\theta} - \psi_i|$ . Thus, we conclude that for a fixed  $\boldsymbol{\theta}$  in Problem (b), the optimal value is given by the Huber function

$$\sum_{i=1}^n h(\boldsymbol{\phi}_i^\top \boldsymbol{\theta} - \psi_i).$$

5. (30 points) We test the performance of three regression methods on the wine data set <http://archive.ics.uci.edu/ml/datasets/Wine+Quality>. We will only consider the red wine data set, with 1599 samples. We use the first 1400 samples for training, and the last 199 samples for testing. The goal is to build a linear model of the first 11 features (together with a constant term) to predict the quality of the wine. All models are trained by solving the following optimization problem

$$\underset{\mathbf{w}, \beta}{\text{minimize}} \quad \sum_{i=1}^n \ell(\mathbf{x}_i^\top \mathbf{w} + \beta - y_i),$$

where the loss functions are

- least squares loss  $\ell(t) = t^2$
- Huber loss defined in the previous problem, with  $M = 1$
- hinge (deadzone-linear) loss

$$\ell(t) = \begin{cases} 0 & |t| \leq 0.5 \\ |t| - 0.5 & |t| > 0.5 \end{cases}$$

The least squares loss can be directly solved by the command `Phi\y` for some properly defined `Phi`. For the latter two, you will use the `cvx` package found on Prof. Boyd’s website <https://web.stanford.edu/~boyd/software.html>. Report their prediction performance on the test set using a different metric, mean absolute error (MAE), defined as  $(1/n) \sum_{i=1}^n |y_i - \hat{y}_i|$ .

**Solution.** The returned MAEs on the test set are 0.5330, 0.5327, and 0.5481, respectively. We see that their performances are essentially the same. Since the given scores are only integers, an average deviation of  $\pm 0.5$  is reasonably good, but not particularly impressive.

6. (30 points) We test the performance of three classification methods on the ionosphere data set <https://archive.ics.uci.edu/ml/datasets/ionosphere>. There are 351 samples. We use the first 300 samples for training, and the last 51 samples for testing. The goal is to build a linear model of the 34 features (together with a constant term) to predict the binary ( $\pm 1$ ) outcome. All models are trained by solving the following optimization problem

$$\underset{\mathbf{w}, \beta}{\text{minimize}} \quad \sum_{i=1}^n \ell(\mathbf{x}_i^\top \mathbf{w} + \beta, y_i),$$

where the loss functions are

- least squares loss  $\ell(t, y) = (yt - 1)^2$
- logistic loss  $\ell(t, y) = \log(1 + \exp(-yt))$
- hinge loss  $\ell(t, y) = \max(0, 1 - yt)$

Again, you will use the backslash command to solve for the first model, and `cvx` to solve for the latter two. Report their prediction accuracy on the test set.

**Solution.** The returned prediction accuracies on the test set are all 100% correct. This is perhaps because of the fact that all of the last 51 samples are in the “good” category, which makes it somewhat easier to guess.