CS7643: Deep Learning Spring 2020 Problem Set 0

Instructor: Zsolt Kira
TAs: Rahul Duggal, Jiachen Yang, Sameer Dharur, Yinquan Lu

Patrick Grady, Anishi Mehta

Discussions: https://piazza.com/gatech/spring2020/cs4803d17643a/home

Due: Tuesday, Jan 14, 11:55pm

Instructions

- 1. We will be using Gradescope to collect your assignments. Please read the following instructions for submitting to Gradescope carefully! Failure to follow these instructions may result in parts of your assignment not being graded. We will not entertain regrading requests for failure to follow instructions.
 - For Section 1: Multiple Choice Questions, it is mandatory to use the LATEX template provided on the class webpage (https://www.cc.gatech.edu/classes/AY2020/cs7643_fall/assets/ps0.zip). For every question, there is only one correct answer. To mark the correct answer, change \choice to \CorrectChoice
 - For Section 2: Proofs, each problem/sub-problem is in its own page. This section has 5 total problems/sub-problems, so you should have 5 pages corresponding to this section. Your answer to each sub-problem should fit in its corresponding page.
 - For Section 2, LATEX'd solutions are strongly encouraged (solution template available at https://www.cc.gatech.edu/classes/AY2020/cs7643_fall/assets/ps0.zip), but scanned handwritten copies are acceptable. If you scan handwritten copies, please make sure to append them to the pdf generated by LATEX for Section 1.
- 2. Hard copies are **not** accepted.
- 3. We generally encourage you to collaborate with other students. You may talk to a friend, discuss the questions and potential directions for solving them. However, you need to write your own solutions and code separately, and *not* as a group activity. Please list the students you collaborated with.

Exception: PS0 is meant to serve as a background preparation test. You must NOT collaborate on PS0.

1 Multiple Choice Questions

1. (1 point) true/false We are machine learners with a slight gambling problem (very different from gamblers with a machine learning problem!). Our friend, Bob, is proposing the following payout on the roll of a dice:

$$payout = \begin{cases} \$1 & x = 1\\ -\$1/4 & x \neq 1 \end{cases}$$
 (1)

where $x \in \{1, 2, 3, 4, 5, 6\}$ is the outcome of the roll, (+) means payout to us and (-) means payout to Bob. Is this a good bet i.e are we expected to make money?

- 2. (1 point) X is a continuous random variable with the probability density function:

$$p(x) = \begin{cases} 4x & 0 \le x \le 1/2 \\ -4x + 4 & 1/2 \le x \le 1 \end{cases}$$
 (2)

Which of the following statements are true about equation for the corresponding cumulative density function (cdf) C(x)?

[Hint: Recall that CDF is defined as $C(x) = Pr(X \le x)$.]

- $C(x) = 2x^2$ for $0 \le x \le 1/2$
- $C(x) = -2x^2 + 4x 3/2 \text{ for } 1/2 \le x \le 1$
- None of the above
- 3. (2 point) A random variable x in standard normal distribution has following probability density

$$p(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \tag{3}$$

Evaluate following integral

$$\int_{-\infty}^{\infty} p(x)(ax^2 + bx + c)dx \tag{4}$$

[*Hint:* We are not sadistic (okay, we're a little sadistic, but not for this question). This is not a calculus question.]

 \bigcirc a + b + c \bigcirc c \bigcirc a + c \bigcirc b + c

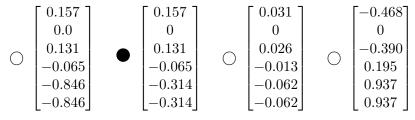
4. (2 points) Consider the following function of $\mathbf{x} = (x_1, x_2, x_3, x_4, x_5, x_6)$:

$$f(\mathbf{x}) = \sigma \left(\log \left(5 \left(\max\{x_1, x_2\} \cdot \frac{x_3}{x_4} - (x_5 + x_6) \right) \right) + \frac{1}{2} \right)$$
 (5)

where σ is the sigmoid function

$$\sigma(x) = \frac{1}{1 + e^{-x}} \tag{6}$$

Compute the gradient $\nabla_{\mathbf{x}} f(\cdot)$ and evaluate it at at $\hat{\mathbf{x}} = (5, -1, 6, 12, 7, -5)$.



- 5. (2 points) Which of the following functions are convex?
 - $\bigcup ||\mathbf{x}||_{\frac{1}{2}}$
 - $\bigcirc \min_{i} \mathbf{a}_{i}^{T} \mathbf{x} \text{ for } \mathbf{x} \in \mathbb{R}^{n}$
 - \bullet log $(1 + \exp(\mathbf{w}^T \mathbf{x}_i))$ for $\mathbf{w} \in \mathbb{R}^d$
 - All of the above
- 6. (2 points) Suppose you want to predict an unknown value $Y \in \mathbb{R}$, but you are only given a sequence of noisy observations $x_1...x_n$ of Y with i.i.d. noise $(x_i = Y + \epsilon_i)$. If we assume the noise is I.I.D. Gaussian $(\epsilon_i \sim N(0, \sigma^2))$, the maximum likelihood estimate (\hat{y}) for Y can be given by:
 - \bigcirc A: $\hat{y} = \operatorname{argmin}_{y} \sum_{i=1}^{n} (y x_i)^2$
 - \bigcirc B: $\hat{y} = \operatorname{argmin}_{y} \sum_{i=1}^{n} |y x_{i}|$
 - \bigcirc C: $\hat{y} = \frac{1}{n} \sum_{i=1}^{n} x_i$
 - Both A & C
 - O Both B & C

2 Proofs

7. (3 points) Prove that

$$\log_e x \le x - 1, \qquad \forall x > 0 \tag{7}$$

with equality if and only if x = 1.

[Hint: Consider differentiation of $\log(x) - (x - 1)$ and think about concavity/convexity and second derivatives.]

Proof:

Let us consider

$$f(x) = \log_e x - (x - 1) \tag{8}$$

Then,

$$\frac{df}{dx} = \frac{1}{x} - 1\tag{9}$$

And,

$$\frac{d^2f}{dx^2} = -\frac{1}{x^2} \tag{10}$$

We know that the second derivative of a function determines the concavity/convexity of a function. We see that the second derivative of f(x) here is always negative. This means that the first derivative of f(x) gives the local maxima at f'(x) = 0. Equating the first derivative obtained in equation (9) above to zero, we get that f(x) will have its maxima at x=1. Substituting the value x=1 in f(x), we get f(1) = 0. This means that:

$$f(x) \le 0 \tag{11}$$

Hence,

$$\log_e x \le x - 1, \qquad \forall x > 0 \tag{12}$$

8. (6 points) Consider two discrete probability distributions p and q over k outcomes:

$$\sum_{i=1}^{k} p_i = \sum_{i=1}^{k} q_i = 1 \tag{13a}$$

$$p_i > 0, q_i > 0, \quad \forall i \in \{1, \dots, k\}$$
 (13b)

The Kullback-Leibler (KL) divergence (also known as the *relative entropy*) between these distributions is given by:

$$KL(p,q) = \sum_{i=1}^{k} p_i \log \left(\frac{p_i}{q_i}\right)$$
(14)

It is common to refer to KL(p,q) as a measure of distance (even though it is not a proper metric). Many algorithms in machine learning are based on minimizing KL divergence between two probability distributions. In this question, we will show why this might be a sensible thing to do.

[Hint: This question doesn't require you to know anything more than the definition of KL(p,q) and the identity in Q7]

(a) Using the results from Q7, show that KL(p,q) is always non-negative.

Proof:

Given:

$$KL(p,q) = \sum_{i=1}^{k} p_i \log\left(\frac{p_i}{q_i}\right)$$
 (15)

Then,

$$-KL(p,q) = -\sum_{i=1}^{k} p_i \log\left(\frac{p_i}{q_i}\right) = \sum_{i=1}^{k} p_i \log\left(\left(\frac{p_i}{q_i}\right)^{-1}\right) = \sum_{i=1}^{k} p_i \log\left(\frac{q_i}{p_i}\right)$$
(16)

As proved in Q.7. above,

$$\log_e x \le x - 1, \qquad \forall x > 0 \tag{17}$$

Substituting $x = \frac{q_i}{p_i}$ in the above equation, we get:

$$\sum_{i=1}^{k} p_i \log \left(\frac{q_i}{p_i}\right) \le \sum_{i=1}^{k} p_i \left(\frac{q_i}{p_i} - 1\right) \tag{18}$$

$$\sum_{i=1}^{k} p_i \log \left(\frac{q_i}{p_i}\right) \le \sum_{i=1}^{k} (q_i - p_i) \tag{19}$$

$$\sum_{i=1}^{k} p_i \log \left(\frac{q_i}{p_i}\right) \le \sum_{i=1}^{k} q_i - \sum_{i=1}^{k} p_i \tag{20}$$

From eqn. (13a) above, RHS becomes zero as $\sum_{i=1}^k p_i = \sum_{i=1}^k q_i = 1$. Hence,

$$\sum_{i=1}^{k} p_i \log \left(\frac{q_i}{p_i} \right) \le 0 \tag{21}$$

Multiplying both LHS and RHS by -1, we get

$$\sum_{i=1}^{k} p_i \log \left(\frac{p_i}{q_i} \right) \ge 0 \tag{22}$$

(b) When is KL(p,q) = 0?

Solution: From Q.8a above, we know that,

$$KL(p,q) = \sum_{i=1}^{k} p_i \log \left(\frac{p_i}{q_i}\right)$$
 (23)

$$KL(p,q) = 0$$
 when $\sum_{i=1}^{k} p_i \log \left(\frac{p_i}{q_i}\right) = 0$.

KL(p,q)=0 when $\sum_{i=1}^k p_i \log \left(\frac{p_i}{q_i}\right)=0$. This will happen when $\log \left(\frac{p_i}{q_i}\right)=0$ for all i, that is $p_i=q_i$ for all i. This means that all the points of the two probability distributions overlap, and this means that they are identical. Hence, the KL divergence, or the distance between the two probability distributions, is zero.

(c) Provide a counterexample to show that the KL divergence is not a symmetric function of its arguments: $KL(p,q) \neq KL(q,p)$

Proof:

We have:

$$KL(p,q) = \sum_{i=1}^{k} p_i \log\left(\frac{p_i}{q_i}\right)$$
 (24)

and

$$KL(q,p) = \sum_{i=1}^{k} q_i \log \left(\frac{q_i}{p_i}\right)$$
 (25)

Then,

$$KL(p,q) - KL(q,p) = \sum_{i=1}^{k} p_i \log\left(\frac{p_i}{q_i}\right) - \sum_{i=1}^{k} q_i \log\left(\frac{q_i}{p_i}\right)$$
 (26)

$$= \sum_{i=1}^{k} p_i \log \left(\frac{p_i}{q_i}\right) + \sum_{i=1}^{k} q_i \log \left(\frac{p_i}{q_i}\right)$$
 (27)

$$= \sum_{i=1}^{k} (p_i + q_i) \log \left(\frac{p_i}{q_i}\right) \tag{28}$$

Taking an example of two probability distribution functions where $p_1 = \frac{1}{2}$ and $p_2 = \frac{1}{2}$, and $q_1 = \frac{3}{5}$ and $q_2 = \frac{2}{5}$, by substituting values we get:

$$KL(p,q) - KL(q,p) = \sum_{i=1}^{k} (p_i + q_i) \log\left(\frac{p_i}{q_i}\right)$$
(29)

$$= (p_1 + q_1) \log \left(\frac{p_1}{q_1}\right) + \sum_{i=1}^k (p_2 + q_2) \log \left(\frac{p_2}{q_2}\right)$$
 (30)

$$= (\frac{1}{2} + \frac{3}{5})\log\frac{5}{6} + (\frac{1}{2} + \frac{2}{5})\log\frac{5}{4}$$
 (31)

$$= -\frac{11}{10}0.081 + \frac{9}{10}0.097 \tag{32}$$

$$=0.0018$$
 (33)

From the above example, we can see that $KL(p,q)-KL(q,p)\neq 0$. Hence, $KL(p,q)\neq KL(q,p)$.

9. (6 points) In this question, you will prove that cross-entropy loss for a softmax classifier is convex in the model parameters, thus gradient descent is guaranteed to find the optimal parameters. Formally, consider a single training example (\mathbf{x}, y) . Simplifying the notation slightly from the implementation writeup, let

$$\mathbf{z} = W\mathbf{x} + \mathbf{b},\tag{34}$$

$$p_j = \frac{e^{z_j}}{\sum_k e^{z_k}},\tag{35}$$

$$L(W) = -\log(p_u) \tag{36}$$

Prove that $L(\cdot)$ is convex in W.

[*Hint:* One way of solving this problem is "brute force" with first principles and Hessians. There are more elegant solutions.]

Proof:

Given:

$$L(W) = -\log(p_y) = -\log\frac{e^{z_j}}{\sum_k e^{z_k}} = \log\frac{\sum_k e^{z_k}}{e^{z_j}}$$
(37)

Finding the partial derivative of L(W) with respect to W,

$$\frac{\partial L}{\partial W} = \frac{\partial s}{\partial t} \tag{38}$$

We need to prove that $L(\cdot)$ is convex in W. So, if the second derivative of L with respect to W is greater than of equal to zero on an interval, then L is convex in W. Thus, by chain rule, we get:

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial p_y} \frac{\partial p_y}{\partial z_y} \frac{\partial z_y}{\partial W} \tag{39}$$

$$= -\frac{\partial(\log p_y)}{p_y} \frac{\partial p_y}{\partial z_y} \frac{\partial z_y}{\partial W} \tag{40}$$

$$= -\frac{1}{p_y} \frac{\partial p_y}{\partial z_y} \frac{\partial z_y}{\partial W} \tag{41}$$

Since p_y is a softmax classifier, we know that the derivative of softmax is $p_y(1-p_y)$. Substituting this value in the above equation, we get:

$$\frac{\partial L}{\partial W} = -\frac{1}{p_y} p_y (1 - p_y) \frac{\partial z_y}{W} \tag{42}$$

$$= -\frac{1}{p_y} p_y (1 - p_y) \mathbf{x} \tag{43}$$

$$= -(1 - p_y)\mathbf{x} \tag{44}$$

Now, finding the double derivative of L with respect to W,

$$[2]LW = \frac{\partial(-(1-p_y)\mathbf{x})}{\partial p_y} \frac{\partial p_y}{\partial z_y}$$
(45)

$$= -(-\mathbf{x})p_y(1-p_y)\mathbf{x} \tag{46}$$

$$= \mathbf{x}p_y(1 - p_y)\mathbf{x} \tag{47}$$

Since $p_y \ge 0$ and $p_y \le 1$, we can see that the second derivative of L with respect to W will always be zero. This means that L is convex in W.