# CS7643: Deep Learning
## Fall 2019
## Problem Set 0

Instructor: Dhruv Batra
TAs: Harsh Agrawal, Neha Jain, Ashwin Kalyan, Harish Kamath,
Anishi Mehta, Nirbhay Modhe, Michael Piseno, Viraj Prabhu, Sarah Wiegreffe
Discussions: https://piazza.com/gatech/fall2019/cs48037643

Due: Thursday, August 22, 11:00am

**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:

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

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?

○ True   ● **False**

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

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

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 $C(x) = Pr(X \leq x)$.]

- ● $C(x) = 2x^2$ **for** $0 \leq x \leq 1/2$
- ○ $C(x) = -2x^2 + 4x - 3/2$ for $1/2 \leq x \leq 1$
- ○ All of the above
- ○ 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.]
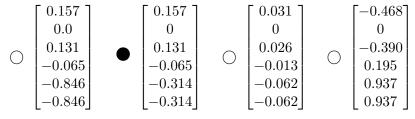
○ a + b + c   ○ c   ● **a + c**   ○ b + c

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

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

where $\sigma$ 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 $\hat{\mathbf{x}} = (5, -1, 6, 12, 7, -5)$.

$\bigcirc$ $\begin{bmatrix} 0.157 \\ 0.0 \\ 0.131 \\ -0.065 \\ -0.846 \\ -0.846 \end{bmatrix}$
$\bullet$ $\begin{bmatrix} 0.157 \\ 0 \\ 0.131 \\ -0.065 \\ -0.314 \\ -0.314 \end{bmatrix}$
$\bigcirc$ $\begin{bmatrix} 0.031 \\ 0 \\ 0.026 \\ -0.013 \\ -0.062 \\ -0.062 \end{bmatrix}$
$\bigcirc$ $\begin{bmatrix} -0.468 \\ 0 \\ -0.390 \\ 0.195 \\ 0.937 \\ 0.937 \end{bmatrix}$

5. (2 points) Which of the following functions are convex?

$\bigcirc$ $||\mathbf{x}||_{\frac{1}{2}}$

$\bigcirc$ $\min_i \mathbf{a}_i^T \mathbf{x}$ for $\mathbf{x} \in \mathbb{R}^n$

$\bigcirc$ $\log(1 + \exp(\mathbf{w}^T \mathbf{x}_i))$ for $\mathbf{w} \in \mathbb{R}^d$

$\bullet$ **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:

$\bullet$ **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: $\hat{y} = \frac{1}{n} \sum_{i=1}^n x_i$

$\bigcirc$ Both A & C

$\bigcirc$ 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.]

Let $f(x) = log(x)$ and $g(x) = x - 1$.

Then, $f'(x) = \frac{1}{x}$ and $g'(x) = 1$.

Therefore, $\forall x \in (0, 1)$, both first-derivatives are monotonically increasing and $f(x) < g(x)$.

We know $\lim_{x \to 0^+} f(x) = -\infty$ and g(0) = -1. We also know f'(x) < g'(x) $\forall x \in (0, 1)$, so we can conclude f(x) < g(x) $\forall x \in (0, 1)$.

Finally, we want to show the equality iff $x = 1$.

First, we show $x = 1 \implies log(x) = x - 1$: $log_e(x) = log_e(1) = 0 = (1 - 1) = (x - 1)$.

Next, we show $log_e(x) = (x - 1) \implies x = 1$:

$log_e(x) = (x - 1)$

$e^{log_e(x)} = e^{x-1}$

$x = e^{x-1}$

$1 = e^{1-1}$

$1 = e^0$

$1 = 1$ so x $= 1$

$\therefore$ We have shown $log_e x \le x - 1$ $\forall x > 0$ and $log_e x = x - 1$ iff $x = 1$.  ■

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

$$\sum_{i=1}^{k} p_i = \sum_{i=1}^{k} q_i = 1 \tag{8a}$$

$$p_i > 0, q_i > 0, \quad \forall i \in \{1, \ldots, k\} \tag{8b}$$

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

$$KL(p, q) = \sum_{i=1}^{k} p_i \log\left(\frac{p_i}{q_i}\right) \tag{9}$$

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.

We will show $KL(p, q) \geq 0$ by showing $-KL(p, q) \leq 0$.

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

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

$\leq \sum_{i=1}^{k} p_i(\frac{q_i}{p_i} - 1)$ since Q7 showed us $log(x) \leq x - 1$

$= \sum_{i=1}^{k} q_i - \sum_{i=1}^{k} p_i$

$= 1 - 1$

$= 0$

Since we have shown $-KL(p, q) \leq 0$, we can respectively conclude $KL(p, q) \geq 0$.

$\therefore KL(p, q)$ is always non-negative. ∎

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

When probability distributions $p$ and $q$ are identical.

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

Let $p$ $B(2, \frac{1}{2})$ and $q$ $B(2, \frac{1}{3})$ where $B(k, p)$ is the Bernoulli distribution with $k$ samples and true-class probability $p$.

$KL(p, q) = \frac{1}{2}log(\frac{\frac{1}{2}}{\frac{1}{3}}) + \frac{1}{2}log(\frac{\frac{1}{2}}{\frac{2}{3}}) = 0.088 + (-0.062) = 0.026$

$KL(q, p) = \frac{1}{3}log(\frac{\frac{1}{3}}{\frac{1}{2}}) + \frac{2}{3}log(\frac{\frac{2}{3}}{\frac{1}{2}}) = (-0.059) + 0.083 = 0.025$

So, by switching the arguments, we show the KL divergence is not symmetric.

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

$$\mathbf{z} = W\mathbf{x} + \mathbf{b}, \tag{10}$$

$$p_j = \frac{e^{z_j}}{\sum_k e^{z_k}}, \tag{11}$$

$$L(W) = -\log(p_y) \tag{12}$$

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.]

Assume $b \in \mathbb{R}$ and $W, x \in \mathbb{R}^n$.

First, we will show $\forall j \ p_j \in (0, 1]$.

We know $z \in \mathbb{R}$ since it's the addition of two real scalars.

(*) We also know $\forall k \ e^{z_k} > 0$ since $\forall x \in \mathbb{R} \ e^x > 0$ ($e^x$ is strictly increasing).

Due to (*), $\sum_k e^{z_k}$ is strictly positive and $\forall j \ e^{z_j} < \sum_k e^{z_k}$.

Since $\forall j \ e^{z_j} < \sum_k e^{z_k}$ and both values are positive, real numbers, we get $p_j = \frac{e^{z_j}}{\sum_k e^{z_k}} \in (0, 1]$.

$\therefore \forall j \ p_j \in (0, 1]$.

Now, we will show $L(W) = -log(p_y)$ is a strictly decreasing function.

By taking its derivative, we see $L'(W) = -\frac{1}{p_y}$.

Since we have already shown $p_y$ is always positive, we conclude $L'(W)$ is always negative.

Additionally, as a side note, we observe $\lim_{x \to 0^+} -log(x) = +\infty$ and $-log(1) = 0$.

So, since $L'(W)$ is always negative and $L(W)$'s range is bounded by $(+\infty, 0]$, we can conclude $L(W)$ is a strictly decreasing function with range $(+\infty, 0]$ in the domain of its input $p_y$.

$\therefore L(W)$ is strictly decreasing.

Finally, we will show $L(W)$ is concave up.

We already showed $L'(W) = -\frac{1}{p_y}$, so the second derivative is $L''(W) = \frac{1}{p_y^2}$.

Since $p_y \in (0, 1]$, $L''(W)$ is strictly positive.

The second derivative measures convexity, so $L''(W) > 0 \implies L(W)$ is concave up.

$\therefore L(W)$ is convex.

We have shown given arbitrary $W$, $p_j$ is limited to the range $(0, 1]$, limiting $L(W)$ to the range $(+\infty, 0]$ as a strictly decreasing function, resulting in a convex function.

So, we have $L(\cdot)$ is convex in $W$. ∎