

PURDUE

BME 64600 – 001 and ECE 60146 – 001

Midterm #2, Spring 2023

NAME _____

PUID _____

Exam instructions:

- You have 75 minutes to work the exam.
- This is a closed-book and closed-note exam. You may not use or have access to your book, notes, any supplementary reference, a calculator, or any communication device including a cell-phone or computer.
- You may not communicate with any person other than the official proctor during the exam.

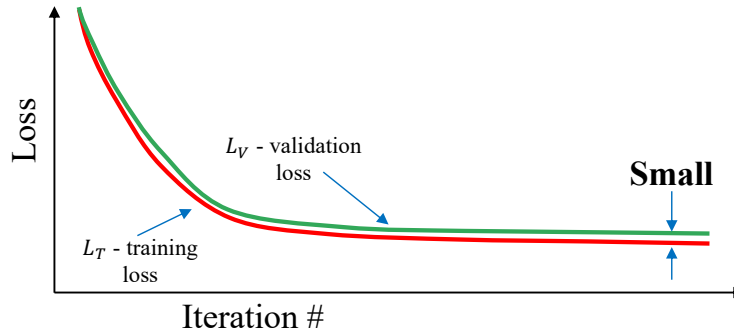
To ensure Gradescope can read your exam:

- Write your full name and PUID above and on the top of every page.
- Answer all questions in the area designated for each problem.
- Write only on the front of the exam pages.
- DO NOT run over to the next question.

Name/PUID: _____ **Key**

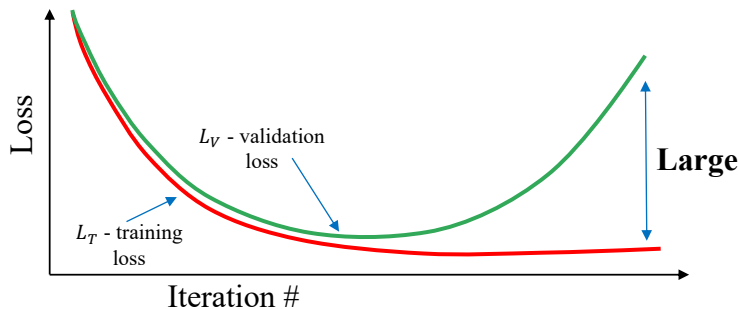
Problem 1. (18pt) Training Convergence

For each of the plots below, answer the associated questions.



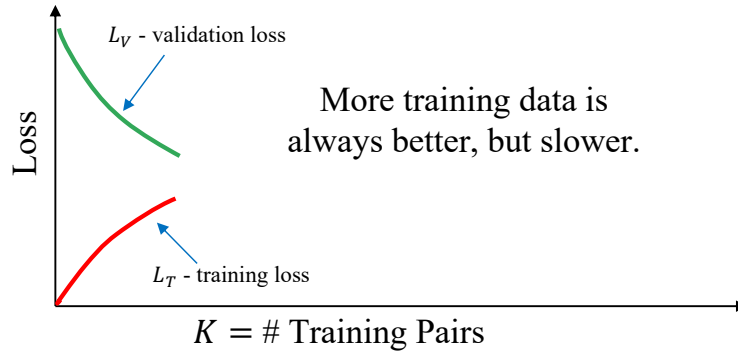
a) What is the data telling you? Based on this, what action(s) might you take?

Solution: The capacity of the network may be too small. You might want to increase the capacity by adding more layers, structures or parameters.



b) What is the data telling you? Based on this, what action(s) might you take?

Solution: The capacity of the network may be too large. You might want to decrease the capacity by a) removing layers or parameters, b) increasing the amount of training data, or c) introducing some regularization of the network parameters.



c) What is the data telling you? Based on this, what action(s) might you take?

Solution: You probably do not have enough training data. You might want to increase the amount of training data, or reduce the number of parameters in the network, or introduce some regularization of the network parameters.

Name/PUID: _____

Problem 2. (48pt) Stochastic Gradient Descent

Consider a problem in which you are doing supervised training of an inference network, $f_\theta : \mathfrak{R}^{N_y} \rightarrow \mathfrak{R}^{N_x}$, for $\theta \in \mathfrak{R}^p$. Let the loss function for each training pair, (x_k, y_k) , be given by

$$L_k(\theta) = \frac{1}{2} \|x_k - f_\theta(y_k)\|^2 ,$$

and let the total loss function be given by

$$L(\theta) = \sum_{k=0}^{K-1} L_k(\theta) .$$

Also assume the gradient of the inference network is given by

$$A_k = \nabla_\theta f_\theta(y_k) ,$$

where $A_k \in \mathfrak{R}^{N_x \times p}$, and define $g_k \in \mathfrak{R}^p$ given by

$$g_k = \nabla_\theta L_k(\theta) ,$$

to be a column vector equal to the gradient of the loss function for the k^{th} training pair.

Let G be a p -dimensional random column vector produced by selecting the gradient of the loss for a single training pair at random with a uniform distribution. Then the probability density of G is given by

$$p(g) = \frac{1}{K} \sum_{k=0}^{K-1} \delta(g - g_k) ,$$

where $\delta(x)$ is a multi-dimensional delta function.

Furthermore, let G_0, \dots, G_{K_b-1} be K_b independently sampled (with replacement) gradients, and let

$$\hat{G} = \frac{1}{K_b} \sum_{k=0}^{K_b-1} G_k ,$$

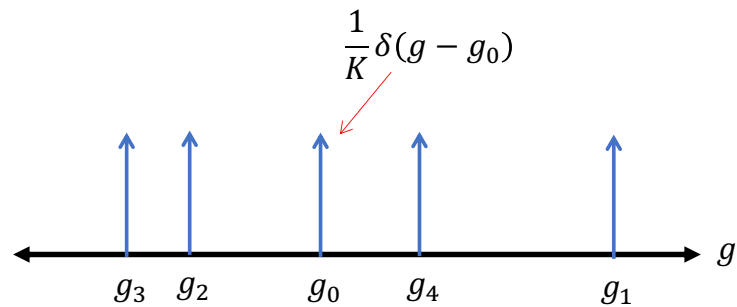
be the average of their gradients.

a) Calculate an expression for $g_k = \nabla_{\theta} L_k(\theta)$.

Solution:

$$g_k = [A_k]^t (x_k - f_{\theta}(y_k))$$

b) Sketch the probability density $p(g)$ for a typical case. Label your sketch to explain the meaning of its structure. For your sketch, you can assume $p = 1$.



Solution:

c) Calculate μ , the mean of G .

Solution:

$$\begin{aligned} \mu &= E[G] \\ &= \int_{\mathbb{R}^p} g \frac{1}{K} \sum_{k=0}^{K-1} \delta(g - g_k) dg \\ &= \frac{1}{K} \sum_{k=0}^{K-1} \int_{\mathbb{R}^p} g \delta(g - g_k) dg \\ &= \frac{1}{K} \sum_{k=0}^{K-1} g_k \end{aligned}$$

d) Calculate R , the covariance of G .

Solution:

$$\begin{aligned} R &= E[(G - \mu)(G - \mu)^t] \\ &= \int_{\mathbb{R}^p} (g - \mu)(g - \mu)^t \frac{1}{K} \sum_{k=0}^{K-1} \delta(g - g_k) dg \\ &= \frac{1}{K} \sum_{k=0}^{K-1} \int_{\mathbb{R}^p} (g - \mu)(g - \mu)^t \delta(g - g_k) dg \\ &= \frac{1}{K} \sum_{k=0}^{K-1} (g_k - \mu)(g_k - \mu)^t \end{aligned}$$

e) Calculate $\hat{\mu}$, the mean \hat{G} .

Solution:

$$\begin{aligned} \hat{\mu} &= E[\hat{G}] \\ &= E\left[\frac{1}{K_b} \sum_{k=0}^{K_b-1} G_k\right] \\ &= \frac{1}{K_b} \sum_{k=0}^{K_b-1} E[G_k] \\ &= \frac{1}{K_b} \sum_{k=0}^{K_b-1} \mu \\ &= \mu \end{aligned}$$

f) Calculate an expression for \hat{R} , the covariance of \hat{G} .

(Hint: This calculation is a bit tricky, so you might just want to guess at the answer for partial credit.)

Solution:

$$\begin{aligned}\hat{R} &= E \left[\left(\hat{G} - \mu \right) \left(\hat{G} - \mu \right)^t \right] \\ &= E \left[\left(\frac{1}{K_b} \sum_{k=0}^{K_b-1} G_k - \mu \right) \left(\frac{1}{K_b} \sum_{k=0}^{K_b-1} G_k - \mu \right)^t \right] \\ &= E \left[\left(\frac{1}{K_b} \sum_{k=0}^{K_b-1} (G_k - \mu) \right) \left(\frac{1}{K_b} \sum_{k=0}^{K_b-1} (G_k - \mu) \right)^t \right] \\ &= E \left[\left(\frac{1}{K_b} \sum_{i=0}^{K_b-1} (G_i - \mu) \right) \left(\frac{1}{K_b} \sum_{j=0}^{K_b-1} (G_j - \mu) \right)^t \right] \\ &= E \left[\frac{1}{K_b^2} \sum_{i=0}^{K_b-1} \sum_{j=0}^{K_b-1} (G_i - \mu)(G_j - \mu)^t \right] \\ &= \frac{1}{K_b^2} \sum_{i=0}^{K_b-1} \sum_{j=0}^{K_b-1} E [(G_i - \mu)(G_j - \mu)^t] \\ &= \frac{1}{K_b^2} \sum_{i=0}^{K_b-1} \sum_{j=0}^{K_b-1} E [(G_i - \mu)(G_i - \mu)^t] \delta(i - j) \\ &= \frac{1}{K_b^2} \sum_{i=0}^{K_b-1} E [(G_i - \mu)(G_i - \mu)^t] \\ &= \frac{1}{K_b^2} \sum_{i=0}^{K_b-1} R \\ &= \frac{1}{K_b} R\end{aligned}$$

g) What is/are advantage(s) of choosing K_b to be larger?

Solution: By choosing K_b large, you can reduce the variance in your estimate of the gradient. This is better for exactly getting at the precise minimum of the loss function.

h) What is/are advantage(s) of choosing K_b to be smaller?

Solution: By choosing K_b smaller, you increase the variance in your estimate of the gradient. This can be advantageous in certain situations since it can help to jump out of local minima of the loss function. It is also faster to do each update since the number of training pairs in the batch is smaller.

Name/PUID: _____

Problem 3. (48pt) Optimization

Consider a problem where $\tilde{Y} \sim \tilde{p}(y)$ and $Y \sim p(y)$ are two random vectors in \mathfrak{R}^N , and define the likelihood ratio as

$$R(y) = \frac{\tilde{p}(y)}{p(y)} .$$

Also, assume the technical condition that $\tilde{p}(y)$ is absolutely continuous with respect to $p(y)$ i.e., that $\forall y \in \mathfrak{R}^N, R(y) < \infty$.¹

Also define the function

$$C(R) = E [(1 + R(Y)) \log(1 + R(Y))] ,$$

and let R_a and R_b be any two valid likelihood ratios.

a) Is $R(y) = 2$ a valid likelihood ratio? Why or why not?

Solution: No, $R(y) = 2$ is not a valid likelihood ratio. To see this, notice that then $\tilde{p}(y) = 2p(y)$, but this implies that $\int_{\mathfrak{R}^N} \tilde{p}(y) dy = \int_{\mathfrak{R}^N} 2p(y) dy = 2$. However, this is not possible, so the likelihood ratio must not be valid.

b) Show that for every valid likelihood ratio, $E [R(Y)] = 1$.

Solution:

$$\begin{aligned} E [R(Y)] &= \int_{\mathfrak{R}^N} R(y) p(y) dy \\ &= \int_{\mathfrak{R}^N} \frac{\tilde{p}(y)}{p(y)} p(y) dy \\ &= \int_{\mathfrak{R}^N} \tilde{p}(y) dy \\ &= 1 \end{aligned}$$

¹This is not exactly the correct definition of absolutely continuous, but it's good enough.

c) Show that for any valid likelihood ratio, R , and integrable function, $f : \mathfrak{R}^N \rightarrow \mathfrak{R}$, then

$$E [f(\tilde{Y})] = E [f(Y)R(Y)] .$$

Solution:

$$\begin{aligned} E [f(\tilde{Y})] &= \int_{\mathfrak{R}^N} f(y) \tilde{p}(y) dy \\ &= \int_{\mathfrak{R}^N} f(y) \frac{\tilde{p}(y)}{p(y)} p(y) dy \\ &= \int_{\mathfrak{R}^N} f(y) R(y) p(y) dy \\ &= E [f(Y)R(Y)] \end{aligned}$$

d) Show that the function $h(x) = (1 + x) \log(1 + x)$ is strictly convex for all $x \geq 0$.

Solution:

$$\begin{aligned} \frac{dh(x)}{dx} &= \log(1 + x) + \frac{1 + x}{1 + x} \\ &= \log(1 + x) + 1 \end{aligned}$$

So then

$$\frac{d^2h(x)}{dx^2} = \frac{1}{1 + x} .$$

So then $\forall x \geq 0$, we have that $\frac{d^2h(x)}{dx^2} > 0$, which implies that $h(x)$ is strictly convex.

e) Show that $\forall \lambda \in (0, 1)$, then $R = \lambda R_a + (1 - \lambda)R_b$ is a valid likelihood ratio.

Solution:

$$\begin{aligned} E[R(Y)] &= E[\lambda R_a(Y) + (1 - \lambda)R_b(Y)] \\ &= \lambda E[R_a(Y)] + (1 - \lambda)E[R_b(Y)] \\ &= \lambda + (1 - \lambda) = 1 \end{aligned}$$

Furthermore, if R_a and R_b are valid, then $R(y) < \infty$.

f) Show that C is a strictly convex function by showing that $\forall \lambda \in (0, 1)$, then

$$C(R) < \lambda C(R_a) + (1 - \lambda)C(R_b) .$$

Solution:

$$\begin{aligned} C(R) &= E[(1 + R(Y)) \log(1 + R(Y))] \\ &= E[h(R(Y))] \\ &= E[h(\lambda R_a + (1 - \lambda)R_b)] \\ &\leq E[\lambda h(R_a) + (1 - \lambda)h(R_b)] \\ &= \lambda E[h(R_a)] + (1 - \lambda)E[h(R_b)] \\ &= \lambda C(R_a) + (1 - \lambda)C(R_b) \end{aligned}$$

g) Use these results to argue that if $C(R)$ has a local minimum, then the local minimum must be the unique global minimum.

Solution: Since it is a strictly convex function, any local minimum must be the unique global minimum.

h) Use these results to argue that the unique global minimum of $C(R)$ must occur for the likelihood ratio given by $R(y) = 1$.

Solution: First notice that for $x = 1$, $\frac{dh(x)}{dx} = 1 + \log 2$.

To do this, we need to your the method of variations. However, we get the correct answer if we differentiate with respect to R as if it was a finite dimensional vector. We also need to enforce the constraint that $E[R(Y)] = 1$, so we will minimize the following Lagrangian

$$L(R) = C(R) + \alpha E[R(Y)] ,$$

by differentiating with respect to R .

$$\begin{aligned} \nabla_R L(R) &= \nabla_R [C(R) + \alpha E[R(Y)]] \\ &= \nabla_R [E[h(R(Y))] + \alpha E[R(Y)]] \\ &= \nabla_R E[h(R(Y)) + \alpha R(Y)] \\ &= \nabla_R \int_{\mathfrak{R}^N} [h(R(y)) + \alpha R(y)] p(y) dy \\ &= \int_{\mathfrak{R}^N} \nabla_R [h(R(y)) + \alpha R(y)] p(y) dy \\ &= \int_{\mathfrak{R}^N} [1 + \log 2 + \alpha] p(y) dy \\ &= 1 + \log 2 + \alpha \end{aligned}$$

If we choose $\alpha = \frac{-1}{1+\log 2}$, then we find that

$$\nabla_R L(R) = 0 ,$$

which means that $R = 1$ is a minimum to the constrained optimization problem.
