

Problem Set 3

Oke

CEE 616: Probabilistic Machine Learning

10.18.2025

Due October 25, 2025 at 11:59PM. Submit on Moodle.

The standard problems are worth a total of **53 points**.

Problem 1 *Batch learning (5 pts)*

In training a neural network via batch learning, the weight updates are given by:

$$\boldsymbol{\theta}_{\ell,t+1} = \boldsymbol{\theta}_{\ell,t} - \rho \frac{1}{N} \sum_{n=1}^N \frac{\partial \mathcal{L}_n(\boldsymbol{\theta}_{\ell,t})}{\partial \boldsymbol{\theta}_{\ell,t}} \quad (1)$$

Now, imagine you are training an MLP with $N = 4000$ observations.

- (a) If you use a batch learning approach (i.e. all observations are used to compute a single weight update), how many iterations/weight updates are obtained in one training epoch? [1]
- (b) If you use a stochastic minibatch approach with a batch size $B = 32$, how many training iterations are obtained in one epoch? (Show your calculation.) [2]
- (c) What is the maximum number of iterations possible in 50 epochs if your batch size must be no less than 10 samples (note that sampling for SGD approaches is typically done without replacement)? [2]

Problem 2 *Subdifferential of ReLU (7 pts)*

The **subdifferential** $\partial f(\mathbf{x})$ of a function f at a point \mathbf{x} is the *set* of all the **subgradients** \mathbf{g} at that point. Thus, for instance, if $f(x) = |x|$, then the subdifferential at $x = 0$ is given by

$$\partial f(0) = [-1, 1] \quad (2)$$

At $x = 1$, however, there is only one unique subgradient: $\partial f(1) = \{1\}$, which is the subdifferential (set of 1 element). Note that $f(x)$ is not differentiable at $x = 0$, *but* it is subdifferentiable.

- (a) The rectified linear unit activation function (ReLU) is defined as: [2]

$$\text{ReLU}(a) = \max(a, 0) = a\mathbb{I}(a > 0) \quad (3)$$

Sketch/plot $\text{ReLU}(a)$ for $a \in [-5, 5]$.

- (b) Write the piecewise subdifferential function $\partial\varphi(a)$, where $\varphi \equiv \text{ReLU}$. (Refer to **PMLI** 8.1.4.1 for an example of how to do this). [3]
- (c) Sketch/plot the subdifferential function $\partial\varphi(a)$, where $\varphi \equiv \text{ReLU}$ (label both axes). [2]

Problem 3 *Optimizers (6 pts)*

PMLI Section 8.4.5 describes approaches to reduce the variance in SGD.

- (a) Read the section
- (b) List three of the SGD variants described in the section and write their weight update equations. [6]

Problem 4 *Feed forward neural network (15 pts)*

The neural network shown in Figure 1 is to be trained for a binary classification problem, $y_n \in \{0, 1\}$. Logistic sigmoid activations are used in the hidden layer neurons.

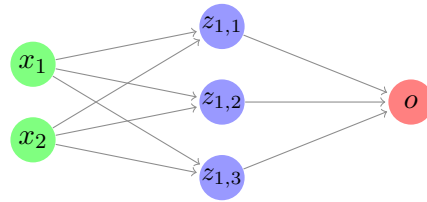


Figure 1

- [1] (a) What is the dimensionality of the input vector \mathbf{x} (i.e. what is D)?
- [2] (b) How many parameters (weights/biases) are required to train this network? (Show your calculations explicitly.)
- (c) We can write the hidden units (activations) at layer 1 as:

$$\mathbf{z}_1 = \begin{bmatrix} z_{1,1} \\ z_{1,2} \\ z_{1,3} \end{bmatrix} \quad (4)$$

Thus,

$$\mathbf{z}_1 = \sigma(\mathbf{b}_1 + \mathbf{W}_1 \mathbf{x}) \quad (5)$$

where $\mathbf{x} = [x_1, x_2]^\top$.

- [1] i. What is the length of \mathbf{b}_1 ?
- [1] ii. What are the dimensions of \mathbf{W}_1 ?
- [2] iii. Write an equation for $z_{1,3}$ in terms of x_1, x_2 , the relevant elements of \mathbf{W}_1 (e.g. $w_{3,1}$, etc.) and the relevant element(s) of \mathbf{b}_1 . The equation should be written in scalar form and should include the exponential function explicitly.
- (d) The weights in the final (output) layer are given by $\boldsymbol{\theta}_2 = (\mathbf{w}_2, b_2)$.
- [1] i. What is the length of \mathbf{w}_2 ?
- [1] ii. Given that this is a binary classification problem and only a single neuron is specified in the final (output) layer, what activation function should be used in the output layer?
- [2] iii. Write an explicit equation for o in terms of $\mathbf{w}_2, \mathbf{z}_1$ and b_2 .
- [1] iv. What does o represent?
- [3] (e) Write/derive the binary cross-entropy loss function \mathcal{L}_n for an observation n in this problem in terms of o_n and y_n .

Problem 5 *Backpropagation for MLP (20 pts)*

(From PMLI Exercise 13.1) Consider the following classification MLP with one hidden layer:

$$\mathbf{x} = \text{input vector} \in \mathbb{R}^D \quad (6)$$

$$\mathbf{z} = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b}_1) \in \mathbb{R}^K \quad (7)$$

$$\mathbf{h} = \text{ReLU}(\mathbf{z}) \in \mathbb{R}^K \quad (8)$$

$$\mathbf{a} = \mathbf{V}\mathbf{h} + \mathbf{b}_2 \in \mathbb{R}^C \quad (9)$$

$$\mathcal{L} = \text{CrossEntropy}(\mathbf{y}, \text{softmax}(\mathbf{a})) \in \mathbb{R} \quad (10)$$

where $\mathbf{x} \in \mathbb{R}^D$, $\mathbf{b}_1 \in \mathbb{R}^K$, $\mathbf{W}_1 \in \mathbb{R}^{K \times D}$, $\mathbf{V} \in \mathbb{R}^{C \times K}$, where D is the input dimension, K is the number of hidden units, and C is the number of classes. Show that the gradients for the parameters and input are given by:

(a)

$$\nabla_{\mathbf{V}} \mathcal{L} = \left[\frac{\partial \mathcal{L}}{\partial \mathbf{V}} \right]_{1,:} = \mathbf{u}_2 \mathbf{h}^\top \in \mathbb{R}^{C \times K} \quad (11)$$

(b)

$$\nabla_{\mathbf{b}_2} \mathcal{L} = \left(\frac{\partial \mathcal{L}}{\partial \mathbf{b}_2} \right)^\top = \mathbf{u}_2 \in \mathbb{R}^C \quad (12)$$

(c)

$$\nabla_{\mathbf{W}} \mathcal{L} = \left[\frac{\partial \mathcal{L}}{\partial \mathbf{W}} \right]_{1,:} = \mathbf{u}_1 \mathbf{x}^\top \in \mathbb{R}^{K \times D} \quad (13)$$

(d)

$$\nabla_{\mathbf{b}_1} \mathcal{L} = \left(\frac{\partial \mathcal{L}}{\partial \mathbf{b}_1} \right)^\top = \mathbf{u}_1 \in \mathbb{R}^K \quad (14)$$

(e)

$$\nabla_{\mathbf{x}} \mathcal{L} = \left(\frac{\partial \mathcal{L}}{\partial \mathbf{x}} \right)^\top = \mathbf{W}^\top \mathbf{u}_1 \in \mathbb{R}^D \quad (15)$$

where the gradients of the loss with respect to the logit layer and hidden layer are given by:

$$\mathbf{u}_2 = \left(\frac{\partial \mathcal{L}}{\partial \mathbf{a}} \right)^\top = \text{softmax}(\mathbf{a}) - \mathbf{y} \in \mathbb{R}^C \quad (16)$$

$$\mathbf{u}_1 = \left(\frac{\partial \mathcal{L}}{\partial \mathbf{z}} \right)^\top = \mathbf{V}^\top \mathbf{u}_2 \odot H(\mathbf{z}) \in \mathbb{R}^K \quad (17)$$

and where \odot denotes the element-wise (Hadamard) product and $H(\mathbf{z})$ is the Heaviside step function defined as:

$$H(z_i) = \begin{cases} 1 & z_i > 0 \\ 0 & \text{otherwise} \end{cases} \quad (18)$$