

HOMWORK 4

Varun Kaundinya
9082410649

Instructions: Use this latex file as a template to develop your homework. Submit your homework on time as a single zip (containing the pdf and the code) file to Canvas. Please, check Piazza for updates about the homework. It is recommended that you use Python for your solutions. You are allowed to use the libraries numpy, matplotlib, pandas and pytorch.

1 Questions (50 pts)

- (10 Points) Suppose the world generates a single observation $x \sim \text{multinomial}(\theta)$, where the parameter vector $\theta = (\theta_1, \dots, \theta_k)$ with $\theta_i \geq 0$ and $\sum_i \theta_i = 1$. Note $x \in \{1, \dots, k\}$. You know θ and want to predict x . Call your prediction \hat{x} . What is your expected 0-1 loss:

$$\mathbb{E}[1\{\hat{x} \neq x\}]$$

using the following two prediction strategies respectively? Prove your answer.

- (5 pts) Strategy 1: $\hat{x} \in \text{argmax}_x \theta_x$, the outcome with the highest probability.
- (5 pts) Strategy 2: You mimic the world by generating a prediction $\hat{x} \sim \text{multinomial}(\theta)$. (Hint: your randomness and the world's randomness are independent)

The image shows handwritten mathematical derivations for the softmax function and cross-entropy loss. The top section is titled 'softmax' and shows the derivative of the softmax output with respect to the input. The bottom section is titled 'cross-entropy loss' and shows the derivative of the cross-entropy loss with respect to the softmax output.

softmax

$$\sigma_j = \frac{e^{x_j}}{\sum_i e^{x_i}}$$

$$\frac{\partial \sigma_j}{\partial x_k} = \frac{e^{x_j} \delta_{jk} - e^{x_j} e^{x_k}}{(\sum_i e^{x_i})^2} = \frac{e^{x_j} (\delta_{jk} - \sigma_j \sigma_k)}{(\sum_i e^{x_i})^2}$$

$$\frac{\partial \sigma_j}{\partial x_k} = \sigma_j (\delta_{jk} - \sigma_j)$$

$$\frac{\partial \sigma_j}{\partial x_k} = \sigma_j (\delta_{jk} - \sigma_j)$$

cross-entropy loss

$$L = -\sum_i y_i \log(\hat{y}_i) = -\sum_i y_i \log(\sigma_i)$$

$$\frac{\partial L}{\partial \sigma_i} = -\sum_j y_j \frac{\partial \log(\sigma_j)}{\partial \sigma_j} = -\sum_j y_j \frac{1}{\sigma_j} \frac{\partial \log(\sigma_j)}{\partial \sigma_j}$$

$$= -\sum_j y_j \frac{1}{\sigma_j} \frac{1}{\sigma_j} = -\sum_j y_j \frac{1}{\sigma_j^2}$$

$$= -\sum_j y_j \frac{1}{\sigma_j^2} = -\sum_j y_j \frac{1}{\sigma_j^2}$$

- (10 points) Like in the previous question, the world generates a single observation $x \sim \text{multinomial}(\theta)$. Let $c_{ij} \geq 0$ denote the loss you incur, if $x = i$ but you predict $\hat{x} = j$, for $i, j \in \{1, \dots, k\}$. $c_{ii} = 0$ for all i . This is a way to generalize different costs on false positives vs false negatives from binary classification to multi-class classification. You want to minimize your expected loss:

$$\mathbb{E}[c_{x\hat{x}}]$$

Derive your optimal prediction \hat{x} .

Given $x = \text{multinomial}(0)$
 $C_{ij} > 0 \iff \pi_i = i, \text{ but } \hat{\pi}_i = j$
 $C_{ij} = 0 \iff \pi_i = i, \hat{\pi}_i = j$ (independence)
 $P[C, \hat{\pi}] = \sum_i \sum_j C_{ij} P(\pi_i = i, \hat{\pi}_i = j)$
 $= \sum_i \sum_j C_{ij} P(\pi_i = i) P(\hat{\pi}_i = j)$
 $= \sum_i P(\hat{\pi}_i = j) \sum_j C_{ij} P(\pi_i = i)$
 substituting $P(\pi_i = i) = \theta_i$
 $= \sum_j P(\hat{\pi}_i = j) \sum_i C_{ij} \theta_i$
 $\hat{\pi}_i = N$, assume $P(\hat{\pi}_i = j) = 1$, otherwise 0
 then $\sum_j \sum_i C_{ij} \theta_i$ as rest will be 0
 $16[C_{ij}] = \sum_{j=N} P(\hat{\pi}_i = j = N) \sum_i C_{ij} \theta_i$
 $= \arg\max_{\theta} P(\hat{\pi}_i = N) \left(\sum_i C_{ij} \theta_i \right)$
 substitute $P(\hat{\pi}_i = N) = \frac{1}{N}$ per assumption
 $\hat{\pi}_i = \arg\max_{\theta} \frac{1}{N} \sum_i C_{ij} \theta_i$
 $\hat{\pi}_i = \arg\max_{\theta} \sum_i C_{ij} \theta_i$
 $\hat{\pi}_i = \arg\max_{\theta} \sum_i C_{ij} \theta_i$

3. (30 Points) The Perceptron Convergence Theorem shows that the Perceptron algorithm will not make too many mistakes as long as every example is “far” from the separating hyperplane of the target halfspace. In this problem, you will explore a *variant* of the Perceptron algorithm and show that it performs well (given a little help in the form of a good initial hypothesis) as long as every example is “far” (in terms of angle) from the separating hyperplane of the *current hypothesis*.

Consider the following variant of Perceptron:

- Start with an initial hypothesis vector $w = w^{\text{init}}$.
- Given example $x \in \mathbb{R}^n$, predict according to the linear threshold function $w \cdot x \geq 0$.
- Given the true label of x , update the hypothesis vector w as follows:
 - If the prediction is correct, leave w unchanged.
 - If the prediction is incorrect, set $w \leftarrow w - (w \cdot x)x$.

So the update step differs from that of Perceptron shown in class in that $(w \cdot x)x$ (rather than x) is added or subtracted to w . (Note that if $\|x\|_2 = 1$, then this update causes vector w to become orthogonal to x , i.e., we add or subtract the multiple of x that shrinks w as much as possible.)

Suppose that we run this algorithm on a sequence of examples that are labeled according to some linear threshold function $v \cdot x \geq 0$ for which $\|v\|_2 = 1$. Suppose moreover that

- Each example vector x has $\|x\|_2 = 1$;
- The initial hypothesis vector w^{init} satisfies $\|w^{\text{init}}\|_2 = 1$ and $w^{\text{init}} \cdot v \geq \gamma$ for some fixed $\gamma > 0$;
- Each example vector x satisfies $\frac{|w \cdot x|}{\|w\|_2} \geq \delta$, where w is the current hypothesis vector when x is received. (Note that for a unit vector x , this quantity $\frac{|w \cdot x|}{\|w\|_2}$ is the cosine of the angle between vectors w and x .)

Show that under these assumptions, the algorithm described above will make at most $\frac{2}{\delta^2} \ln(1/\gamma)$ many mistakes.

Upper Bound :

Given the update step, $w \leftarrow w - (w \cdot x)x$.

Taking a square on either side,

$$\|w'\|^2 = \|w\|^2 - 2\|w(w \cdot x)x\| + \|(w \cdot x)x\|^2$$

On Simplifying, with $\|x\|^2 = 1$

$$\begin{aligned}\|\mathbf{w}'\|^2 &= \|\mathbf{w}\|^2 - \|\mathbf{w}x\|^2, \text{ Dividing both sides by } \|\mathbf{w}\|^2 \text{ and substituting } \delta \\ \frac{\|\mathbf{w}'\|^2}{\|\mathbf{w}\|^2} &= 1 - \frac{\|\mathbf{w}x\|^2}{\|\mathbf{w}\|^2} \\ \frac{\|\mathbf{w}'\|^2}{\|\mathbf{w}\|^2} &\leq 1 - \delta^2\end{aligned}$$

After One mistake, $\|\mathbf{w}'\|^2 \leq 1 - \delta^2$

After m mistakes, $\|\mathbf{w}\|_m \leq \sqrt{(1 - \delta^2)^m}$

Lower Bound :

$$w^{init} = -(wx)x$$

$$w_1 = w^{init} - (wx)x$$

$$w_1 - w^{init} = -(wx)x$$

After m mistakes,

$$w_m - w_{init} = -\sum_1^m (wx)x$$

Multiplying by v on either side,

$$w_m \cdot v - w_{init} \cdot v = -\sum_1^m (wx)x \cdot v$$

$$w_m \cdot v = -\sum_1^m (wx)x \cdot v + w_{init} \cdot v$$

Ignoring the negative term and substituting γ ,

$$w_m \cdot v \geq \gamma, \text{ Given } \|v\| = 1$$

$$\|\mathbf{w}_m\| \geq \gamma$$

Combining both the bounds

$$\gamma \leq \|\mathbf{w}\|_m \leq \sqrt{(1 - \delta^2)^m}$$

$$\gamma \leq \sqrt{(1 - \delta^2)^m}$$

Taking Log (ln) on either sides

$$\frac{m}{2} \ln(1 - \delta^2) \geq \ln(\gamma)$$

By applying Taylor's expansion on $\ln(1 - \delta^2)$ we get $\ln(1 - \delta^2)$ approx equal to $-\delta^2$

$$-\frac{m}{2} \delta^2 \geq \ln(\gamma)$$

$$m \leq \frac{2}{\delta^2} \ln(1/\gamma)$$

Proved.

2 Programming (60 pts)

In this exercise, you will derive, implement back-propagation for a simple neural network, and compare your output with some standard library's output. Consider the following 3-layer neural network.

$$\mathbf{y} = \mathbf{f}(\mathbf{x}) = \mathbf{g}(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}))$$

Suppose $\mathbf{x} \in \mathbb{R}^d$, $\mathbf{W}_1 \in \mathbb{R}^{d_1 \times d}$ and $\mathbf{W}_2 \in \mathbb{R}^{k \times d_1}$ i.e., $f: \mathbb{R}^d \mapsto \mathbb{R}^k$, Let $\sigma(\mathbf{z}) = [\sigma(z_1), \dots, \sigma(z_n)]$ for any $\mathbf{z} \in \mathbb{R}^n$ where $\sigma(t) = \frac{1}{1 + \exp(-t)}$ is the sigmoid (logistic) activation function and $\mathbf{g}(\mathbf{z})_i = \frac{\exp(z_i)}{\sum_{i=1}^k \exp(z_i)}$ is the softmax function. Suppose that the true pair is (\mathbf{x}, \mathbf{y}) where $\mathbf{y} \in \{0, 1\}^k$ with exactly one of the entries equal to 1 and you are working with the cross-entropy loss function given below,

$$L(\mathbf{x}, \mathbf{y}) = - \sum_{i=1}^k \mathbf{y} \log(\hat{\mathbf{y}}) .$$

1. Derive backpropagation updates for the above neural network. (10 pts)

softmax activation: $\hat{y}_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$ Softmax in the output layer

derivative of softmax has 2 scenarios

for $i = j$: $\frac{\partial \hat{y}_i}{\partial z_i} = \frac{e^{z_i} - e^{z_i} e^{z_i}}{(e^{z_i})^2} = \frac{e^{z_i} - e^{2z_i}}{(e^{z_i})^2} = \frac{e^{z_i}(1 - e^{z_i})}{(e^{z_i})^2} = \frac{1 - e^{z_i}}{e^{z_i}} = 1 - \hat{y}_i$

for $i \neq j$: $\frac{\partial \hat{y}_i}{\partial z_j} = \frac{-e^{z_i} e^{z_j}}{(e^{z_i})^2} = -\frac{e^{z_i} e^{z_j}}{(e^{z_i})^2} = -\frac{e^{z_j}}{e^{z_i}} = -\frac{\hat{y}_j}{\hat{y}_i}$

cross entropy loss: $\mathcal{L} = -\sum_i y_i \log(\hat{y}_i)$

derivative of cross entropy loss

for $i = j$: $\frac{\partial \mathcal{L}}{\partial z_i} = -\sum_i y_i \frac{\partial \log(\hat{y}_i)}{\partial z_i} = -\sum_i y_i \frac{1}{\hat{y}_i} \frac{\partial \hat{y}_i}{\partial z_i} = -\sum_i y_i \frac{1}{\hat{y}_i} (1 - \hat{y}_i) = -\sum_i y_i \frac{1 - \hat{y}_i}{\hat{y}_i} = -\sum_i \frac{y_i - \hat{y}_i}{\hat{y}_i} = -\sum_i \frac{y_i}{\hat{y}_i} + \sum_i \frac{\hat{y}_i}{\hat{y}_i} = -\sum_i \frac{y_i}{\hat{y}_i} + \sum_i 1 = -\sum_i \frac{y_i}{\hat{y}_i} + 1$

for $i \neq j$: $\frac{\partial \mathcal{L}}{\partial z_j} = -\sum_i y_i \frac{\partial \log(\hat{y}_i)}{\partial z_j} = -\sum_i y_i \frac{1}{\hat{y}_i} \frac{\partial \hat{y}_i}{\partial z_j} = -\sum_i y_i \frac{1}{\hat{y}_i} \left(-\frac{\hat{y}_j}{\hat{y}_i}\right) = \sum_i \frac{y_i \hat{y}_j}{\hat{y}_i^2} = \sum_i \frac{y_i}{\hat{y}_i} \hat{y}_j = \sum_i \frac{y_i}{\hat{y}_i} \hat{y}_j$

gradient w.r.t. weights \rightarrow cross entropy gradient

$\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial \hat{y}_n} \frac{\partial \hat{y}_n}{\partial z_n} \frac{\partial z_n}{\partial w_2}$

$= (1 - \hat{y}_n) \frac{\partial (w_2 x_1)}{\partial w_2}$

$= (1 - \hat{y}_n) x_1$

gradient w.r.t. activations of layer \rightarrow hidden

$\frac{\partial \mathcal{L}}{\partial a_1} = \frac{\partial \mathcal{L}}{\partial \hat{y}_n} \frac{\partial \hat{y}_n}{\partial z_n} \frac{\partial z_n}{\partial a_1} = (1 - \hat{y}_n) \frac{\partial (w_2 x_1)}{\partial a_1}$

$= (1 - \hat{y}_n) w_2$

gradient w.r.t. input of hidden layer \rightarrow sigmoid

$\frac{\partial \mathcal{L}}{\partial z_1} = \frac{\partial \mathcal{L}}{\partial a_1} \frac{\partial a_1}{\partial z_1} = (1 - \hat{y}_n) \sigma(z_1) (1 - \sigma(z_1))$

gradient w.r.t. w_1

$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial z_1} \frac{\partial z_1}{\partial w_1} = (1 - \hat{y}_n) \sigma(z_1) (1 - \sigma(z_1)) x$

2. Implement it in numpy or pytorch using basic linear algebra operations. (e.g. You are not allowed to use auto-grad, built-in optimizer, model, etc. in this step. You can use library functions for data loading, processing, etc.). Evaluate your implementation on MNIST dataset, report test error, and learning curve. (25 pts)

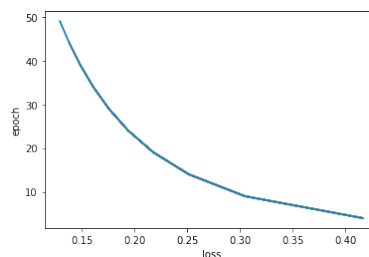


Figure 1: Batch Size = 32, Learning Rate = 0.05, Accuracy = 96.25%

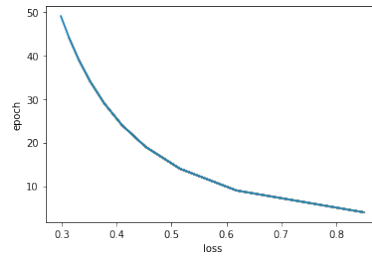


Figure 2: Batch Size = 32, Learning Rate = 0.01, Accuracy = 94.23%

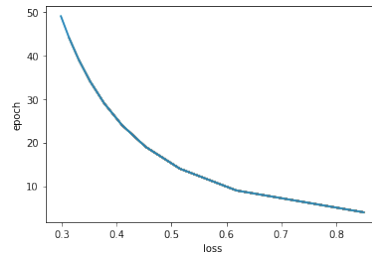


Figure 3: Batch Size = 64, Learning Rate = 0.01, Accuracy = 92.51%

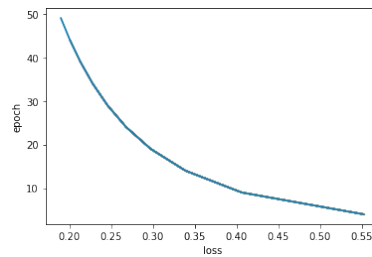


Figure 4: Batch Size = 64, Learning Rate = 0.05, Accuracy = 95.48%

3. Implement the same network in `pytorch` (or any other framework). You can use all the features of the framework e.g. auto-grad etc. Evaluate it on MNIST dataset, report test error, and learning curve. (20 pts)

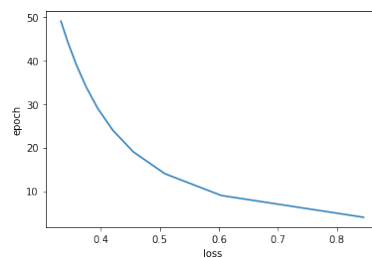


Figure 5: Batch Size = 32, Learning Rate = 0.01, Accuracy = 93.76%

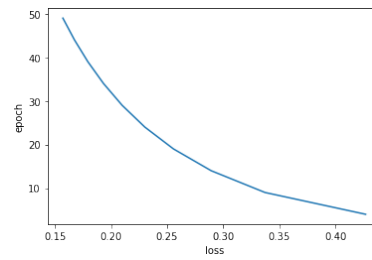


Figure 6: Batch Size = 32, Learning Rate = 0.01, Accuracy = 97.42%

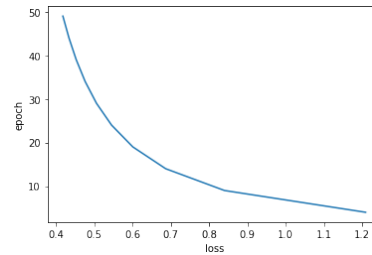


Figure 7: Batch Size = 64, Learning Rate = 0.01, Accuracy = 94.23%

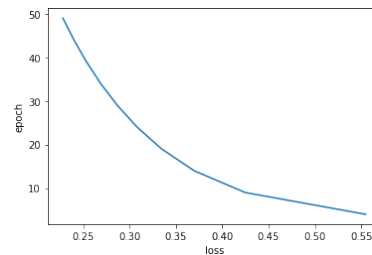


Figure 8: Batch Size = 64, Learning Rate = 0.05, Accuracy = 96.28%

4. Try different weight initializations a) all weights initialized to 0, and b) Initialize the weights randomly between -1 and 1. Report test error and learning curves for both. (You can use either of the implementations) (5 pts)

We see that when all weights are initialized to zero the convergence is very slow, The loss barely decreases over epochs, to compare the models below are the graph run on numpy implementation.

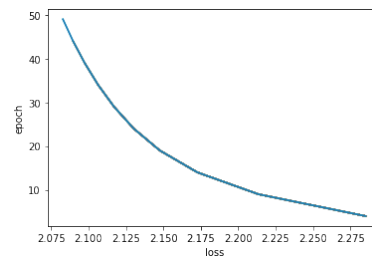


Figure 9: Zero Weights Batch Size = 32, Learning Rate = 0.05

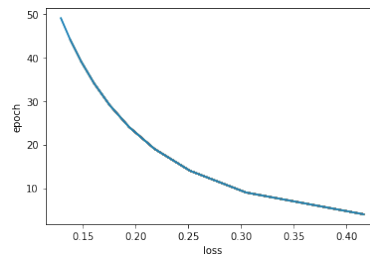


Figure 10: Batch Size = 32, Learning Rate = 0.05

You should play with different hyper-parameters like learning rate, batch size, etc. For Questions 2-3; you should report the answers for at least 3 different sets of hyperparameters. You should mention the values of those along with the results. Use $d_1 = 300$, $d_2 = 200$. For optimization use SGD (Stochastic gradient descent) without momentum, with some batch-size say 32, 64 etc. MNIST can be obtained from here (<https://pytorch.org/vision/stable/datasets.html>)

Check Piazza for a sample PyTorch project

To load the dataset:

```
import torch
import torch.nn.functional as F
from torch.utils.data import TensorDataset, DataLoader
import torchvision

mnist_data_train = torchvision.datasets.MNIST('path/to/download',
                                             train=True, download=True)
data_loader = torch.utils.data.DataLoader(mnist_data,
                                           batch_size=,
                                           shuffle=True)

mnist_data_test = torchvision.datasets.MNIST('path/to/download',
                                             train=False, download=True)
```