HOMEWORK 4

Hafeez Ali Anees Ali
908 460 3423

## 1. Questions (50 pts)

(10 Points) Suppose the world generates a single observation x ~ multinomial(θ), where the parameter vector θ = (θ1, . . . , θk) with θi ≥ 0 and Pki θi = 1. Note x ∈ {1, . . . , k}. You know θ and want to predict x. Call your prediction ˆx. What is your expected 0-1 loss:
E[1{ˆx≠ x}]
using the following two prediction strategies respectively? Prove your answer.
(a) (5 pts) Strategy 1: ˆx ∈ argmaxxθx, the outcome with the highest probability.
(b) (5 pts) Strategy 2: You mimic the world by generating a prediction ˆx ~ multinomial(θ). (Hint: your randomness and the world's randomness are independent)

1) 

$$x \sim \text{multinomial}(\theta)$$
$$\theta = (\theta_1, \ldots \theta_k) \quad \text{with} \quad \theta_i \geq 0 \quad \& \quad \sum_i^k \theta_i = 1$$
$$x \in \{1, \ldots k\}$$
$$E[1\{\hat{x} \neq x\}] = ?$$

a) Strategy 1: $\hat{x} = \text{argmax}_x \theta_x$

$$E1[\hat{x} \neq x] = \sum_{k=1}^{k} P(\hat{x} \neq x) = 1 - P(\hat{x} = x)$$
$$= 1 - \theta_{\hat{x}}$$

$$\boxed{E1[\hat{x} \neq x] = 1 - \theta_{\hat{x}}}$$

b) Strategy 2: $\hat{x} \sim \text{multinomial}(\theta)$

$$E1[\hat{x} \neq x] = \sum_{y=1}^{k} P(\hat{x} \neq y, \hat{x} = y)$$
$$= \sum_{y=1}^{k} P(\hat{x} \neq y) P(\hat{x} = y)$$

(Because it is given that the world's randomness & my randomness are independent of each other)

$$E1[\hat{x} \neq x] = \sum_{y=1}^{k} (1 - P(\hat{x} = y)) P(\hat{x} = y)$$
$$= \sum_{y=1}^{k} (1 - \theta_y) \theta_y$$
$$= \sum_{y=1}^{k} \theta_y - \sum_{y=1}^{k} \theta_y^2 = \boxed{1 - \sum_{y=1}^{k} \theta_y^2}$$

2. (10 points) Like in the previous question, the world generates a single observation x ~ multinomial(θ). Let cij ≥ 0 denote the loss you incur, if x = i but you predict ˆx = j, for i, j ∈ {1, . . . , k}. cii = 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:

E[cxˆx]

Derive your optimal prediction ˆx.

2)

$$x \sim \text{multinomial}(\theta)$$

$c_{ij} \geq 0$ is the loss if $x = i$ but prediction $\hat{x} = j$, $i, j \in \{1, \dots k\}$

$c_{ii} = 0 \quad \forall i$

$$\mathbb{E}[c_{x\hat{x}}] = \sum_{j=1}^{k} \sum_{\substack{i=1, i \neq j}}^{k} c_{ij} \, \theta_i \, \hat{\theta}_j$$

We need $\hat{x}$ such that $\sum_{j=1}^{k} \sum_{\substack{i=1, i \neq j}}^{k} c_{ij} \, \theta_i \, \hat{\theta}_j$ is minimum.

$$\hat{x} \ \underline{\underline{\triangleq}} \ \arg\min_j \left( \sum_{j=1}^{k} \sum_{i=1}^{k} c_{ij} \theta_i \right)$$

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 = winit.
• Given example x ∈ Rn, predict according to the linear threshold function w °§ x ≥ 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 ← w − (w °§ x)x.

So the update step differs from that of Perceptron shown in class in that (w°§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 °§ x ≥ 0 for which ∥ v ∥ 2 = 1. Suppose moreover that

• Each example vector x has ∥ x ∥ 2 = 1; The initial hypothesis vector winit satisfies ∥ winit ∥ 2 = 1 and winit °§ v ≥ γ for some fixed γ > 0;

• Each example vector x satisfies |w°§x| ∥ w ∥ 2 ≥ δ, where w is the current hypothesis vector when x is received. (Note that for a unit vector x, this quantity |w°§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 2 δ2 ln(1/γ) many

mistakes.

---

3) Perceptron Convergence

It is given that the initial hypothesis vector $w^{init}$ satisfies
$\|w^{init}\|_2 = 1$ & $w^{init} \cdot v \geq \gamma$ for some fixed $\gamma > 0$

$$w^{init} \cdot v \geq \gamma \qquad —— ①$$

Let's say a mistake occurred for some data point $x$, then the new hypothesis vector is given by

$$w_1 = w^{init} - (w^{init} \cdot x)\, x$$

$$w_1 \cdot v = w^{init} \cdot v - (w^{init} \cdot x)\, x \cdot v$$

According to the linear threshold function, in the case of an error $w \cdot x < 0$
$$\Rightarrow w^{init} \cdot x < 0$$
$$\text{So} \quad -(w^{init} \cdot x) > 0$$
$$-(w^{init} \cdot x)\, x \cdot v > 0 \qquad —— ②$$

Using ① & ②, we get:

$$w^{init} \cdot v - (w^{init} \cdot x)\, x \cdot v \geq \gamma$$
$$\Rightarrow w_1 \cdot v \geq \gamma$$
$$\Rightarrow w_t \cdot v \geq \gamma \qquad \forall t$$

Since $W_1 = w^{init} - (w^{init} \cdot x) x$

$$\|W_1\|^2 = \|w^{init}\|^2 - 2\|w^{init}\|^2 \|x\|^2 |\cos\theta|$$
$$+ \|w^{init}\|^2 \cos^2\theta \|x\|^4$$

Sine $\|x\| = 1$ is given for all $x$

& taking $\theta$ to be angle between $w^{init}$ & $x$

we have

$$\|W_1\|^2 = \|w^{init}\|^2 \left[ 1 - 2|\cos\theta| + \cos^2\theta \right] \quad -\text{③}$$

It is also given that $x$ satisfies $\dfrac{|w \cdot x|}{\|w\|_2} \geq \delta$

But $\dfrac{|w \cdot x|}{\|w\|_2} = |\cos\theta| \Rightarrow$

$$|\cos\theta| \geq \delta$$
$$-|\cos\theta| \leq -\delta$$
$$1 - |\cos\theta| \leq 1 - \delta \quad -\text{④}$$

Continuing with ③ we have,

$$\|W_1\|^2 = \|w^{init}\|^2 (1 - \cos\theta)^2$$

But $\|w^{init}\|_2 = 1$ is also given

So, $\|W_1\|^2 = (1 - \cos\theta)^2$

$\|W_1\| = (1 - |\cos\theta|)$

From ④, we have $\|W_1\| = (1 - |\cos\theta|) \leq 1 - \delta$

$$-\text{⑤}$$

Extending this, we have $\|W_t\| \leq (1-\delta)^t$ for $t$ mistakes.

The geometric interpretation of a misclassification update is it shrinks the hypothesis vector afer every update.

$$\Rightarrow \quad 1 \geq (1-\delta) \geq 0$$

$$(1-\delta) \leq 1$$
$$0 \geq -\delta$$
$$\delta \geq 0$$

Also, $\quad \delta \leq |\cos\theta| \leq 1$

$$\Rightarrow \quad 0 \leq \delta \leq 1$$

Additionally, $w_t \cdot v \geq \gamma \Rightarrow w_t \cdot v \geq 0 \quad$ since $\gamma > 0$
$$\Rightarrow \quad |w_t \cdot v| \geq \gamma$$

$$\Rightarrow \quad \|w_t\| \|v\| \geq |w_t \cdot v| \geq \gamma$$

Now, $\|v\| = 1$

$$\Rightarrow \quad \|w_t\| \geq \gamma \qquad -\textcircled{6}$$

From $\textcircled{5}$ & $\textcircled{6}$ $\quad \gamma \leq \|w_t\| \leq (1-\delta)^t \qquad -\textcircled{7}$

Since $\quad 0 \leq \delta \leq 1$, we have $\delta \geq \delta^2$
$$\Rightarrow \quad 1 - \delta \leq 1 - \delta^2$$

So $\textcircled{7}$ can be rewritten as $\quad \gamma \leq \|w_t\| \leq (1-\delta^2)^t$

For $0 \leq \delta \leq 1$,

$-1 \leq -\delta \leq 0$

$0 \leq 1 - \delta \leq 1$

$\Rightarrow (1-\delta)^2 \leq (1-\delta)$

but $(1-\delta) \leq 1 - \delta^2$

$\Rightarrow (1-\delta)^2 \leq 1 - \delta^2$

From ⑨, $\gamma \leq \|w_t\| \leq (1-\delta)^{2(t/2)} \leq (1-\delta^2)^{t/2}$

$\Rightarrow \gamma \leq \|w_t\| \leq (1-\delta^2)^{t/2}$

$\Rightarrow \gamma^2 \leq (1-\delta^2)^t \quad (\gamma > 0 \ \& \ (1-\delta^2)^{t/2} > 0)$

$\ln \gamma^2 \leq \ln(1-\delta^2)^t \quad \text{(applying ln on both sides)}$

$2 \ln \gamma \leq t \ln(1-\delta^2) \quad - \ ⑧$

Now, $e^{-\gamma} = 1 - x + \dfrac{x^2}{2!} - \dfrac{x^3}{3!} \ldots$

So $e^{-\delta^2} \geq 1 - \delta^2$

$\ln e^{-\delta^2} \geq \ln(1-\delta^2)$

From ⑧ $\Rightarrow$ $2 \ln \gamma \leq t \ln(1-\delta^2) \leq t \ln e^{-\delta^2}$

$2 \ln \gamma \leq t \ln e^{-\delta^2}$

$2 \ln \gamma \leq -\delta^2 t$

$-2 \ln \gamma \geq \delta^2 t$

$\boxed{t \leq \dfrac{2}{\delta^2} \ln\left(1/\gamma\right)}$

$\Rightarrow$ number of mistakes t is bounded by this inequality

## 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. $y = f(x) = g(W2\sigma(W1x))$

Suppose $x \in Rd, W1 \in Rd1°\o d$ and $W2 \in Rk°\o d1$ i.e., $f : Rd \xrightarrow{7} Rk$, Let $\sigma(z) = [\sigma(z1), \ldots, \sigma(zn)]$ for any $z \in Rn$ where $\sigma(t) = 1 / 1+\exp(-t)$ is the sigmoid (logistic) activation function and $g(z)i = Pexp(zi) k \ i=1 \ exp(zi)$ is the softmax function. Suppose that the true pair is $(x, y)$ where $y \in \{0, 1\}k$ with exactly one of the entries equal to 1 and you are working with the cross-entropy loss function given below,

$L(x, y) = -Xki=1y \ log(\hat{y})$ .

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

2] **Programming:**

1)
$$y = f(x) = g(W_2\, \sigma(W_1 x))$$
$$x \in \mathbb{R}^d, \quad W_1 \in \mathbb{R}^{d_1 \times d}, \quad W_2 \in \mathbb{R}^{k \times d_1}.$$
$$f: \mathbb{R}^d \to \mathbb{R}^k$$
$$\sigma(z) = [\sigma(z_1), \dots, \sigma(z_n)] \quad \text{where } z \in \mathbb{R}^n$$
$$\sigma(t) = \frac{1}{1 + e^{-t}}$$
$$g(z)_i = \frac{e^{z_i}}{\sum_{i=1}^{k} e^{z_i}}$$

Cross entropy loss, $L(x, y) = -\sum_{i=1}^{k} y \log(\hat{y}) \qquad y \in \{0, 1\}^k$
$$= -\sum_{i=1}^{k} y \log(g(z)_i)$$



$\longrightarrow$ CE Loss

* Derivative of softmax function:

Let $p(z)_i = e^{z_i}$
$$q(z)_i = \sum_{\ell=1}^{k} e^{z_\ell}$$

$$\frac{dg(z)_i}{dz_j} = \frac{d}{dz_j}\left(\frac{e^{z_i}}{\sum_{\ell=1}^{k} e^{z_\ell}}\right) = \frac{d}{dz_j}\left(\frac{p(z)_i}{q(z)_i}\right)$$

$$p'(z)_i = \frac{d(p(z)_i)}{dz_j} = \frac{d(e^{z_i})}{dz_j} = \frac{d\, e^{z_i}}{dz_i} \cdot \frac{dz_i}{dz_j} \qquad \text{⊠}$$

$$= e^{z_i} \frac{dz_i}{dz_j} = \begin{cases} e^{z_i} & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

$$q'(x) = \frac{d}{dz_j}\left(\sum_{\ell=1}^{k} e^{z_\ell}\right) = \frac{d}{dz_j}\left(\sum_{\ell=1,\ell\neq j}^{k} e^{z_\ell} + e^{z_j}\right)$$

$$= \frac{d}{dz_j}\left(\sum_{\ell=1,\ell\neq j}^{k} e^{z_\ell}\right) + \frac{d}{dz_j}\left(e^{z_j}\right)$$

$$= 0 + e^{z_j}$$

$$= e^{z_j}$$

So, 
$$\frac{dg(z)_i}{dz_j} = e^{z_i}\frac{\sum_{\ell=1}^{k} e^{z_\ell} - e^{z_j}e^{z_i}}{\left(\sum_{\ell=1}^{k} e^{z_\ell}\right)^2} \quad (\text{when } i = j)$$

$$= e^{z_i}\frac{\left(\sum_{\ell=1}^{k} e^{z_\ell} - e^{z_j}\right)}{\left(\sum_{\ell=1}^{k} e^{z_\ell}\right)^2}$$

$$= \frac{e^{z_i}}{\sum_{\ell=1}^{k} e^{z_\ell}} \cdot \frac{\sum_{\ell=1}^{k} e^{z_\ell} - e^{z_j}}{\sum_{\ell=1}^{k} e^{z_\ell}}$$

$$= g(z)_i \left(1 - g(z)_j\right)$$

$$= \boxed{g(z)_i \left(1 - g(z)_i\right)} \quad (i = j)$$

$$\frac{dg(z)_i}{dz_j} = \frac{0 \times \sum_{\ell=1}^{k} e^{z_\ell} - e^{z_j}e^{z_i}}{\left(\sum_{\ell=1}^{k} e^{z_\ell}\right)^2} \quad (\text{when } i \neq j)$$

$$= \frac{-e^{z_j}e^{z_i}}{\left(\sum_{\ell=1}^{k} e^{z_\ell}\right)^2} = \boxed{-g(z)_i\, g(z)_j}$$

$$(i \neq j)$$

Derivation of cross-entropy loss with softmax function:

$$L = -\sum_{c=1}^{k} y_c \log(g(z)_c)$$

$$\frac{dL}{dz_i} = \frac{d}{dz_i}\left[-\sum_{c=1}^{k} y_c \log(g(z)_c)\right]$$

$$= -\sum_{c=1}^{k} y_c \frac{d(\log(g(z)_c))}{dz_i}$$

$$= -\sum_{c=1}^{k} y_c \frac{d(\log(g(z)_c))}{dg(z)_c} \cdot \frac{dg(z)_c}{dz_i}$$

$$= -\sum_{c=1}^{k} \frac{y_c}{g(z)_c} \frac{dg(z)_c}{dz_i}$$

$$= -\left[\frac{y_i}{g(z)_i} \cdot \frac{dg(z)_i}{dz_i} + \sum_{c=1, c\neq i}^{k} \frac{y_c}{g(z)_c} \frac{dg(z)_c}{dz_i}\right]$$

$$= -\frac{y_i}{g(z)_i} \cdot g(z)_i(1-g(z)_i) - \sum_{c=1, c\neq i}^{k} \frac{y_c}{g(z)_c} \cdot (-g(z)_c g(z)_i)$$

$$= -y_i + y_i g(z)_i + \sum_{c=1, c\neq i}^{k} y_c g(z)_i$$

$$= g(z)_i\left(y_i + \sum_{c=1, c\neq i}^{k} y_c\right) - y_i$$

$$= g(z)_i\left(\sum_{c=1}^{k} y_c\right) - y_i$$

$$= g(z)_i \cdot 1 - y_i \qquad \left(\text{because } \sum_{c=1}^{k} y_c = 1\right)$$

$$= \boxed{g(z)_i - y_i}$$

## Back propagation:

Using matrix notation, we can write

$$\frac{dL}{dz} = \hat{y} - y$$

$$\frac{dL}{dw^2} = \frac{dL}{dz^2} \frac{dz^2}{dw^2}$$

$$= (\hat{y} - y) \frac{d}{dw^2} \left( w^2 (A') \right)$$

$A' = $ output of first layer

$$= \boxed{(\hat{y} - y) \, A'}$$

$$\frac{dL}{dw'} = \frac{dL}{dz^2} \frac{dz^2}{dA'} \frac{dA'}{dz'} \frac{dz'}{dw'}$$

$$= (\hat{y} - y) \frac{d}{dA'} \left( A' w^2 \right) * \frac{d}{dz'} \left( \sigma(z') \right)$$

$$\frac{d}{dw'} \left( w'x \right)$$

$$= (\hat{y} - y) \, w^2 \, \sigma'(z') \, x$$

$$\frac{dL}{dw'} = \boxed{(\hat{y} - y) \, w^2 \, \sigma'(z') \, x}$$

$$\frac{dL}{dw'} = \boxed{(\hat{y} - y) \, w^2 \, \sigma'(w'x) \, x}$$

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)

I.

```
# batch size == 60000 (full dataset)
# and learning rate = 0.1
# and epochs = 1000
```

```
predict(W1, W2, X_test, Y_test)
```

```
89.99000000000001
```

```
test_error = 100 - predict(W1, W2, X_test, Y_test)
```

```
test_error
```

```
10.00999999999991
```

```python
def plot_cost():
    plt.figure()
    plt.plot(epochs, costs)
    plt.xlabel("epochs")
    plt.ylabel("cost")
    plt.show()
```

```
plot_cost()
```



```python
def plot_accuracy():
    plt.figure()
    plt.plot(epochs, accuracies)
    plt.xlabel("epochs")
    plt.ylabel("accuracy")
    plt.show()
```

```
plot_accuracy()
```

II.

```
predict(W1, W2, X_train, Y_train)
```

87.63

```
predict(W1, W2, X_test, Y_test)
```

88.29

```python
# batch_size = 64
# n_iterations = 100000
# learning rate = 0.5
```

```python
test_error = 100 - predict(W1, W2, X_test, Y_test)
test_error
```

11.709999999999994

```python
plot_avg_accuracy()
```

```
plot_avg_cost()
```



III.

```
: # batch_size = 32
  # n_iterations = 500000
  # learning rate = 0.3
```

```
: predict(W1, W2, X_train, Y_train)
```
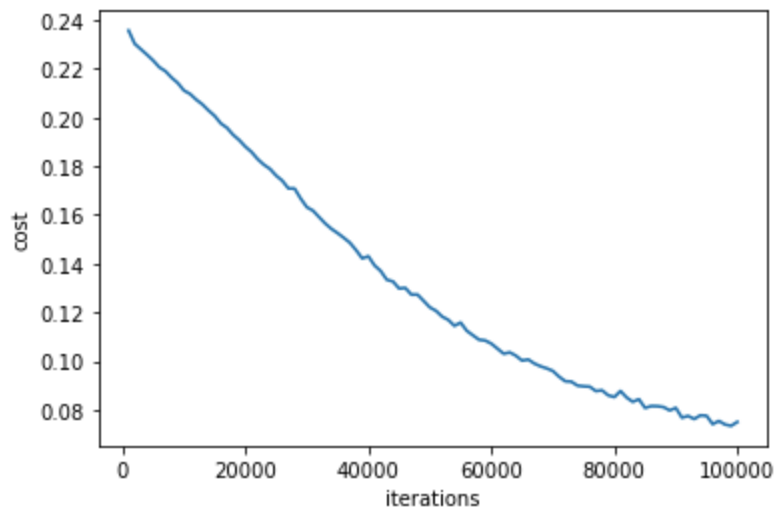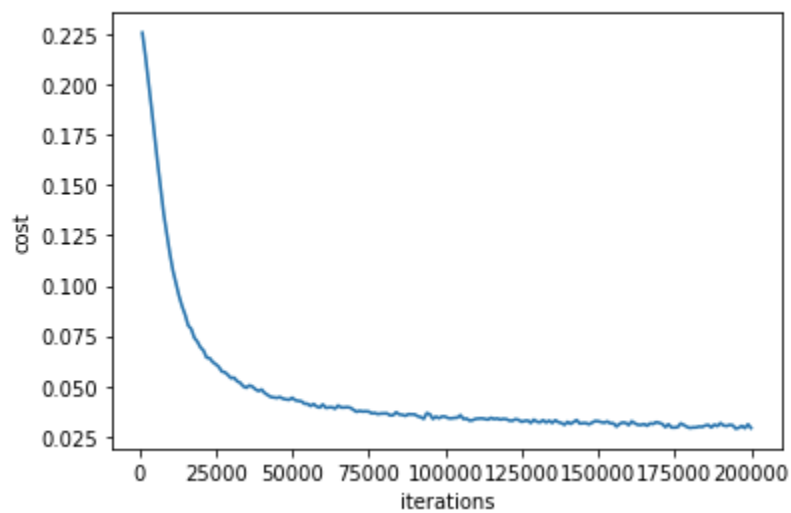
```
: 84.40666666666667
```

```
: predict(W1, W2, X_test, Y_test)
```

```
: 85.35000000000001
```

```
: test_error = 100 - predict(W1, W2, X_test, Y_test)
  test_error
```
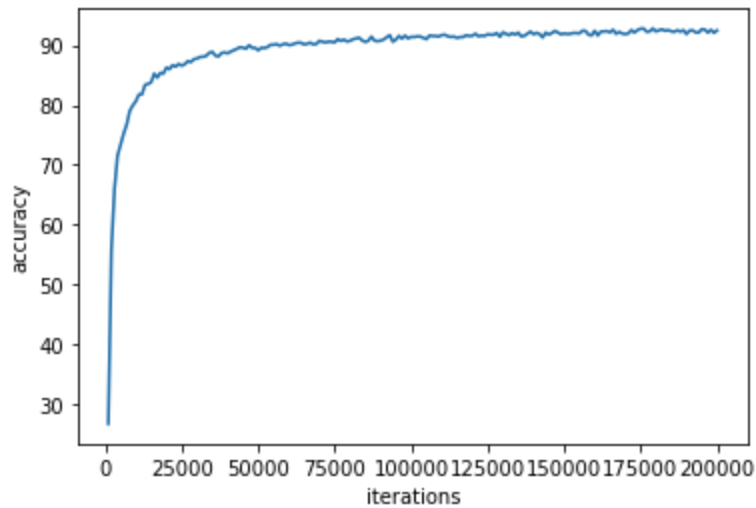
```
: 14.649999999999991
```

```
plot_avg_cost()
```



```
: plot_avg_accuracy()
```
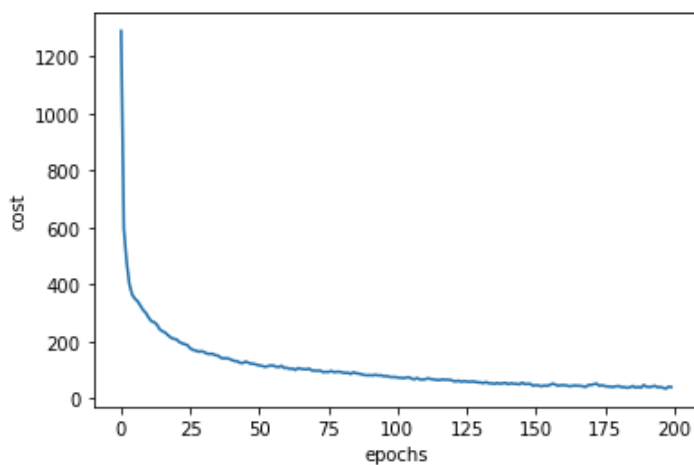


IV.

```
# batch_size = 128
# n_iterations = 200000
# learning rate = 0.7
```

```
predict(W1, W2, X_train, Y_train)
```

91.58166666666666

```
predict(W1, W2, X_test, Y_test)
```

91.79

```
test_error = 100 - predict(W1, W2, X_test, Y_test)
test_error
```

8.209999999999994

```
plot_avg_cost()
```

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)

I.



```
# Batch size = 64
# epochs = 200
# lr = 10

Starting testing...
Test Accuracy: 95.09 %
Test Error: 4.91 %
```
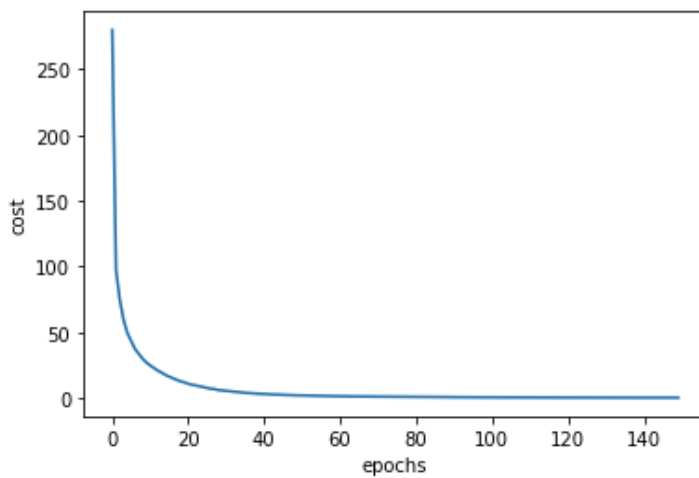
II.



```
|:  # Batch size = 32
    # epochs = 400
    # lr = 0.5

    Starting testing...
    Train Accuracy: 98.17 %
    Test Error: 1.83 %|
```

III.



```
:  # Batch size = 128
   # epochs = 150
   # lr = 1

   Starting testing...
   Test Accuracy: 98.15 %
   Test Error: 1.85 %
```

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)

On Numpy:

    a.  All weights initialized to 0

```
# batch_size = 128
# n_iterations = 200000
# learning rate = 0.7
```

```
predict(W1, W2, X_train, Y_train)
```
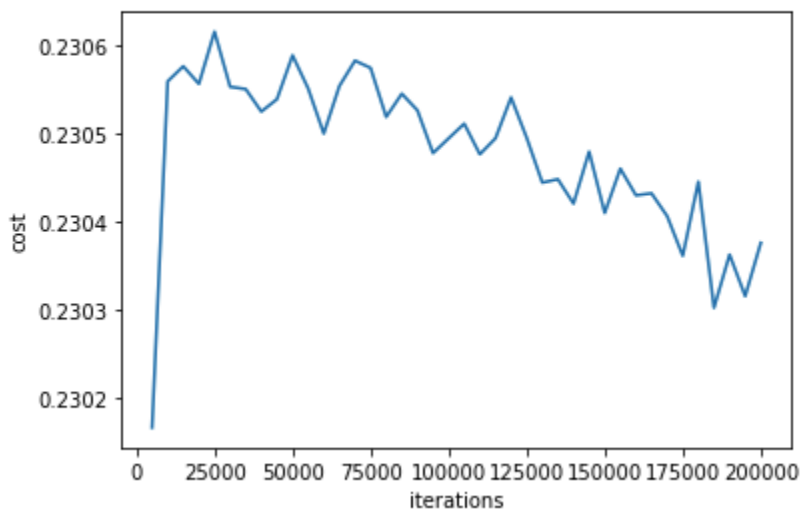
11.236666666666666
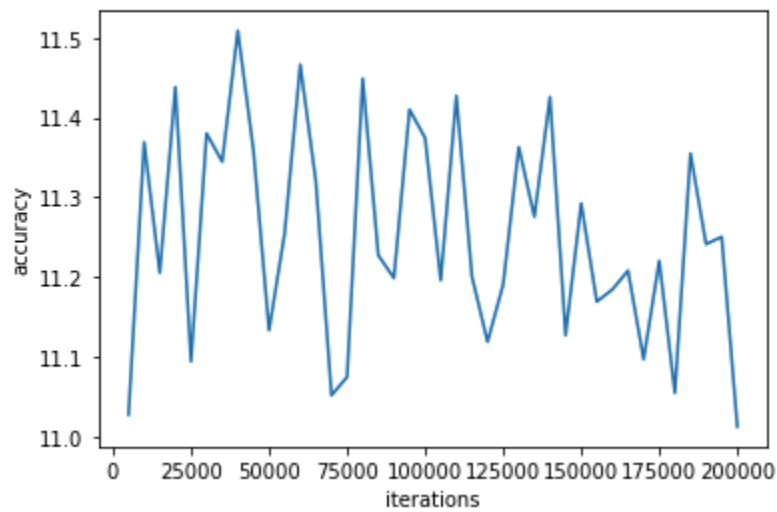
```
predict(W1, W2, X_test, Y_test)
```

11.35

```
test_error = 100 - predict(W1, W2, X_test, Y_test)
test_error
```

88.65

```
plot_avg_cost()
```

```
plot_avg_accuracy()
```



b. Weights initialized randomly between -1 and 1

```
# batch_size = 128
# n_iterations = 200000
# learning rate = 0.7
```

```
predict(W1, W2, X_train, Y_train)
```

92.015

```
predict(W1, W2, X_test, Y_test)
```

91.61

```
test_error = 100 - predict(W1, W2, X_test, Y_test)
test_error
```
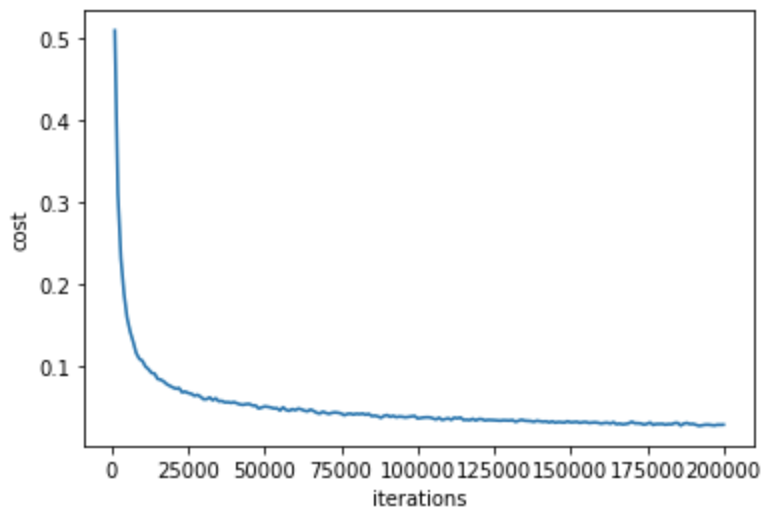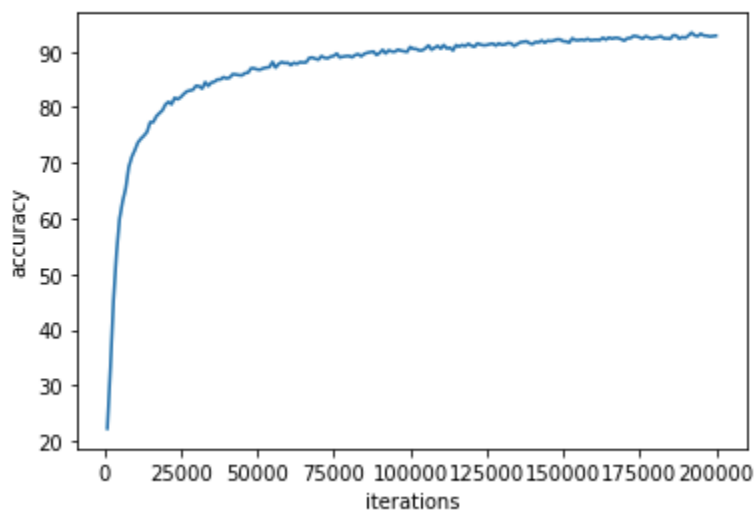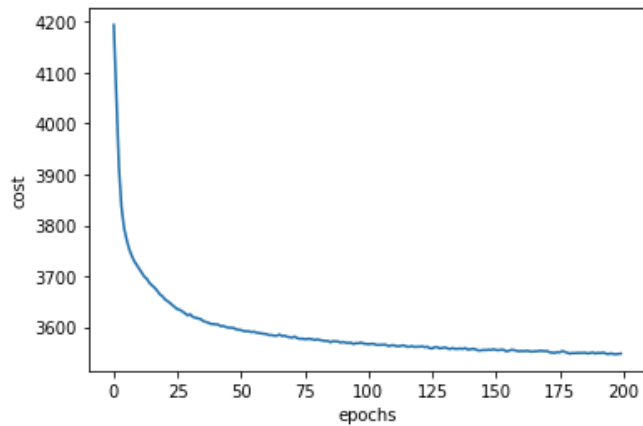
8.39

plot_avg_cost()



plot_avg_accuracy()



On Pytorch:

a.  All weights initialized to 0

```
# Batch size = 32
# epochs = 200
# lr = 0.1

Starting testing...
Train Accuracy: 11.35 %
Test Error: 88.65 %
```

```
plot_cost()
```



b. Weights initialized randomly between -1 and 1

```
# Batch size = 64
# epochs = 200
# lr = 10

Starting testing...
Test Accuracy: 96.81 %
Test Error: 3.19 %
```

```
plot_cost()
```