

## HW2. Back Propagation

Question 2.

$$b) E = -\sum_i t_i \log(t_i) ; x_i = \frac{e^{s_i}}{\sum_{c=1}^m e^{s_c}}$$

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial x_i} \cdot \frac{\partial x_i}{\partial s_i} \cdot \frac{\partial s_i}{\partial w_{ij}} \quad x_i = \frac{e^{s_i}}{e^{s_1} + e^{s_2} + \dots + e^{s_m}}$$

$$\frac{\partial E}{\partial x_c} = -\left(\frac{t_i}{x_c}\right) \cdot \frac{-e^{s_c} \sum e^{s_c} - e^{s_c} \cdot e^{s_c}}{\left(\sum_{c=1}^m e^{s_c}\right)^2} = x_i \cdot (1-x_c)$$

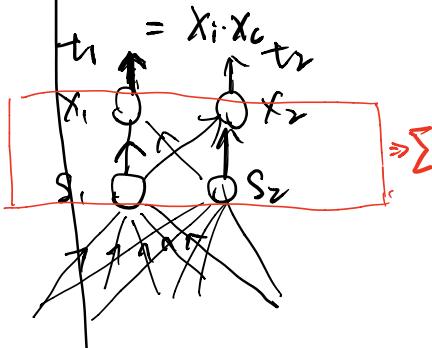
$$\frac{\partial x_c}{\partial s_i} = \begin{cases} \frac{-e^{s_c} \cdot e^{s_c}}{\left(\sum_{c=1}^m e^{s_c}\right)^2} & \text{if } i=c \\ \frac{-e^{s_i} \cdot e^{s_c}}{\left(\sum_{c=1}^m e^{s_c}\right)^2} & \text{if } i \neq c \end{cases}$$

$$\boxed{\frac{\partial E}{\partial x_c}} \cdot \frac{\partial x_c}{\partial s_i} \cdot \frac{\partial s_i}{\partial w_{ji}}$$

$$\frac{e^{s_c} \left( \sum_{c=1}^m e^{s_c} - e^{s_c} \right)}{\left( \sum_{c=1}^m e^{s_c} \right)^2}$$

$$\textcircled{1} \quad \frac{\partial E}{\partial x_c} = -\frac{t_i}{x_c}$$

$$\textcircled{2} \quad \frac{\partial x_c}{\partial s_i} = \begin{cases} \text{if } c=i : \frac{e^{s_c} \sum e^{s_c} - e^{s_c} \cdot e^{s_c}}{\left( \sum_{c=1}^m e^{s_c} \right)^2} = x_i \cdot (1-x_c) \\ \text{if } c \neq i : -\frac{e^{s_c} \cdot e^{s_i}}{\left( \sum_{c=1}^m e^{s_c} \right)^2} = -x_i x_c \end{cases}$$



$$\textcircled{3} \quad \frac{\partial s_i}{\partial w_{ij}} = y_j$$

$$\Rightarrow \frac{\partial E}{\partial w_{ij}} = -\frac{t_i}{x_c} \cdot \left[ x_c \cdot (1-x_c) + \sum_{i \neq c} (-x_i x_c) \right] \cdot y_j$$

$$= \left[ -t_i(1-x_c) + \sum_{i \neq c} t_i x_c \right] \cdot y_j$$

$$= \left[ -t_i + t_i x_c + \sum_{i \neq c} t_i x_c \right] \cdot y_j$$

$$= \left[ \sum_i t_i x_c - t_i \right] \cdot y_j$$

## ② Class 7:

### Logistic Regression

$$p(y_i|x_i) = f(x_i) \cdot (1-f(x_i))^{1-y_i}$$

$$\prod_{i=1}^N p(y_i|x_i) = \prod_{i=1}^N f(x_i)^{y_i} \cdot (1-f(x_i))^{1-y_i} \Rightarrow \text{LIKELIHOOD}$$

$$\log \prod_{i=1}^N p(y_i|x_i) = \sum_{i=1}^N \log [f(x_i)^{y_i} \cdot (1-f(x_i))^{1-y_i}]$$

$$L(\theta) = \sum_{j=1}^N y_j \log f(x_i) + (1-y_j) \log (1-f(x_i))$$

$$\text{Maximum Likelihood: } f(x_i) = \text{sigmoid}(z_i) = \frac{1}{1+e^{-z_i}} = \frac{1}{1+e^{-\theta^T x}} ; z = \theta^T x$$

$$\text{Hence: } L(\theta) = \sum_{i=1}^N y_i \log \frac{1}{1+e^{\theta^T x}} + (1-y_i) \log \left( \frac{e^{\theta^T x}}{1+e^{\theta^T x}} \right)$$

$$\frac{\partial f(x_i)}{\partial \theta} = \frac{+e^{z_i}}{(1+e^{z_i})^2} \cdot \frac{\partial z}{\partial \theta} = f(x_i) \cdot (1-f(x_i)) \cdot (x_i)$$

$$\frac{\partial L(\theta)}{\partial \theta} = \sum_{i=1}^N y_i \frac{f(x_i) \cdot (1-f(x_i)) \cdot (x_i)}{f(x_i)} + (1-y_i) \frac{-1}{1-f(x_i)} \cdot f(x_i) \cdot (1-f(x_i)) \cdot (x_i)$$

$$\leftarrow \sum_{i=1}^N y_i \cdot (1-f(x_i))^{(1-x_i)} + (1-y_i) \cdot (f(x_i))^{(1-x_i)} \cdot (x_i)$$

$$\begin{aligned} \text{MAX-LIKELIHOOD} &= \sum_{i=1}^N \left[ y_i \cdot (1-f(x_i))^{(1-x_i)} + (1-y_i) \cdot (f(x_i))^{(1-x_i)} \right] (x_i) \\ &= \sum_{i=1}^N \left[ y_i \cdot y_i \cdot f(x_i) + y_i \cdot f(x_i) - f(x_i) \cdot (1-x_i) \right] (x_i) \end{aligned}$$

$$\boxed{\sum_{i=1}^N (y_i - f(x_i)) \cdot (x_i) = 0} \quad \boxed{f(x_i) = \frac{1}{1+e^{-\theta^T x}}}$$

Gradient Descent

Likelihood

$$\text{③ posterior} \Rightarrow p(\theta|x) = \frac{p(x|\theta) p(\theta)}{p(x)}$$



## Midterm Exam

- Total duration: 90 minutes.
- You **can** use one page as a cheat sheet.
- You **cannot** consult your notes, textbook, your neighbor, or Google.
- Maximum points: 60.

- 
1. **(5 points)** Please write down the time at the *start* and *end* of your exam. The difference should not exceed 90 minutes. Please also include your *signature* on the first page; by doing so, you are affirming the NYU Tandon School of Engineering student code of conduct.

2. (10 points) This is a slight variant of a homework problem. Let  $\{x_1, x_2, \dots, x_n\}$  be a set of points in  $d$ -dimensional space, and let  $\{\beta_1, \beta_2, \dots, \beta_n\}$  be positive numbers that represent weights. Suppose we wish to produce a single point estimate  $\mu \in \mathbb{R}^d$  that minimizes the *weighted* squared-error:

$$\beta_1 \|x_1 - \mu\|_2^2 + \beta_2 \|x_2 - \mu\|_2^2 + \dots + \beta_n \|x_n - \mu\|_2^2$$

Find a closed form expression for  $\mu$  and prove that your answer is correct.

2.

$$\begin{aligned} f(\mu) &= \beta_1 \|x_1 - \mu\|_2^2 + \beta_2 \|x_2 - \mu\|_2^2 + \dots + \beta_n \|x_n - \mu\|_2^2 \\ f(\mu) &= \sum_{i=1}^n \beta_i \|x_i - \mu\|^2 = \sum_{i=1}^n \beta_i \|x_i^T x_i - 2x_i^T \mu + \mu^T \mu\| \\ \frac{\partial f(\mu)}{\partial \mu} &= \sum_{i=1}^n \beta_i (-2x_i^T + 2\mu) = 0 \\ \sum_{i=1}^n \beta_i (-2x_i^T) &= \sum_{i=1}^n \beta_i (2\mu) \\ \mu &= \frac{\sum_{i=1}^n \beta_i x_i^T}{\sum_{i=1}^n \beta_i} \end{aligned}$$

3. **(10 points)** Assuming a classification application with  $n$  data points in  $d$  dimensions where  $n = 10d = 10 * 10^8$ , rank the following algorithms in decreasing order of (i) training times; (ii) testing times. Include a few sentences justifying your reasoning.

- a. Perceptron.
- b. Nearest neighbors.
- c. Kernel perceptron with a polynomial kernel of order 4.
- d. Kernel perceptron with gaussian kernel.
- e. Linear support vector machines.

4. (15 points) The following represents python code for an algorithm that attempts to perform linear regression. (a) Identify the algorithm. (b) Explain why this algorithm may not converge as implemented below, and identify the line in the algorithm that makes this happen. (c) Suggest a way to fix this algorithm.

```

def optim_alg(init, steps, grad):
    xs = [init]
    for step in steps:
        xs.append(xs[-1] - step * grad(xs[-1]))
    return xs

def linear_reg_grad(X, y, w):
    row_id = numpy.random.randint(X.shape[1])
    x = X[row_id, :]
    return x.T.dot(x.dot(w) - y)

input_to_optim_alg = lambda w: linear_reg_grad(X, y, w)
learning_rates = [0.001]*300
ws = optim_alg(w0, learning_rates, input_to_optim_alg)

```

a)

$$\begin{aligned}
 & \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) \\
 &= \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (y_i - f(x_i))^2 \\
 &= \frac{\partial f}{\partial \theta_0} = \frac{1}{N} \sum_{i=1}^N (y_i - \theta_0 x_i) (-x_i) = 0 \\
 & \frac{\partial f}{\partial \theta_0} = \frac{1}{N} \sum_{i=1}^N (y_i - \theta_0 x_i - \theta_0) \cdot (-1) = 0 \\
 & \sum_{i=1}^N \theta_0 x_i = \sum_{i=1}^N y_i x_i - \theta_0 \sum_{i=1}^N x_i \\
 & \sum_{i=1}^N \theta_0 x_i^2 = \sum_{i=1}^N y_i x_i - \sum_{i=1}^N x_i \left[ \frac{1}{N} \sum_{i=1}^N y_i - \frac{1}{N} \theta_0 \sum_{i=1}^N x_i \right] \\
 & \sum_{i=1}^N \theta_0 x_i^2 = \sum_{i=1}^N y_i x_i - \sum_{i=1}^N x_i \cdot \frac{1}{N} \cdot \sum_{i=1}^N y_i + \frac{1}{N} \theta_0 \sum_{i=1}^N x_i \cdot \sum_{i=1}^N x_i \\
 & \theta_0 = \frac{\sum_{i=1}^N y_i x_i - \frac{1}{N} \sum_{i=1}^N x_i \sum_{i=1}^N y_i}{\sum_{i=1}^N x_i^2 - \frac{1}{N} \sum_{i=1}^N x_i \sum_{i=1}^N x_i}
 \end{aligned}$$

(2-Dimension-case:  $y_i = \theta_0 x_i + \theta_1$ )

b) Gradient Descent may not converge in 300 steps.  
It needs more steps to converge.  
2. learning rate can be smaller to make sure it converges.

c) tolerance > 0.0001  
Set..

General Case  $\Rightarrow$  Stochastic GD.

$$\begin{aligned}
 J(\theta) &= \frac{1}{2} \sum_{i=1}^N (h_\theta(x_i) - y_i)^2 \\
 \frac{\partial J(\theta)}{\partial \theta_j} &= \frac{1}{2} \cdot 2(h_\theta(x_i) - y_i) \cdot \frac{\partial h_\theta(x_i)}{\partial \theta_j} \\
 &= \frac{1}{2} \cdot 2(h_\theta(x_i) - y_i) \cdot \frac{\partial}{\partial \theta_j} h_\theta(x_i) \\
 &= (h_\theta(x_i) - y_i) \cdot x_i
 \end{aligned}$$

5. (10 points) To combat the COVID-19 pandemic, an enterprising NYU Tandon graduate student decides to build a logistic regression model to predict the conditional likelihood of a person being one of three states – susceptible, infected, cured – based on forehead temperature measurements over the last 30 days. Fortunately, such measurements for a population of 10,000 persons is available.

a), b)

K=3

$$d=30, n=10000 \\ nx^30$$

$$X \in \mathbb{R}^{n \times 30}, Y \in \mathbb{R}^{n \times 1}, y_i \in \{-1, 0, 1\}$$

- a. Identify the parameters of the problem (number of samples  $n$ , data dimension  $d$ , number of classes  $k$ .)  
 b. If  $X$  and  $y$  denote the arrays that encode the training data points and labels, what are the sizes of  $X$  and  $y$ ?  
 c. Starting from the definition of conditional likelihood, derive the loss function used to train the model.

$$\begin{aligned} y_i &= \theta^T x \\ f(y_i) &= \text{sigmoid}(y_i) = \frac{1}{1 + e^{-\theta^T x}} \\ L(\theta) &= \sum_{i=1}^N f(y_i)^{y_i} f(1-y_i)^{1-y_i} \\ \log(L(\theta)) &= \sum_{i=1}^N (\log f(y_i)^{y_i} f(1-y_i)^{1-y_i}) \\ &= \sum_{i=1}^N y_i \log f(y_i) + (1-y_i) \log f(1-y_i) \\ \frac{\partial \log(L(\theta))}{\partial \theta} &= \frac{\partial f(y_i)}{\partial \theta} \cdot \frac{\partial z_i}{\partial x_i} = \frac{1}{1 + e^{-\theta^T x}} \cdot \frac{\partial z_i}{\partial \theta} \\ &= \frac{e^{-\theta^T x}}{(1 + e^{-\theta^T x})^2} \cdot x_i \\ &= f(y_i) \cdot (1 - f(y_i)) \cdot x_i \\ \frac{\partial \log(L(\theta))}{\partial \theta} &= \sum_{i=1}^N \left[ y_i \cdot \frac{f(y_i) \cdot (1-f(y_i))}{f(1-y_i)} - (1-y_i) \cdot \frac{f(1-y_i) \cdot (1-f(1-y_i))}{f(y_i)} \right] \cdot x_i \\ &= \sum_{i=1}^N [y_i \cdot (1-f(y_i)) - (1-y_i) \cdot f(y_i)] \cdot x_i \\ &= \sum_{i=1}^N [y_i - y_i f(y_i) - f(y_i) + f(y_i) f(y_i)] \cdot x_i \\ &= \sum_{i=1}^N [y_i - f(y_i)] \cdot x_i \end{aligned}$$

Bernoulli  
Case.

c) Method 1.

Binary Case (Bernoulli Case) =

Step 1: Problem 1:  $X \in \mathbb{R}^{n \times 30}, Y \in \mathbb{R}^{n \times 1}, y_i \in \{-1, 0, 1\}$

Transform label 0 into label 1  
Temporary for step 1.

One-vs-rest

$$y_i \in \{-1, 1\}$$

Step 2: Problem 2:  $X \in \mathbb{R}^{10000 \times 30}, Y \in \mathbb{R}^{10000 \times 1}, y_i \in \{-1, 0, 1\}$   
Transform label 0 into label -1 Temporary.

Method 2.

Multinomial Logistic Regression.

Softmax:

$$h_\theta(x) = \frac{e^{\theta^T x}_i}{\sum_{j=1}^k e^{\theta^T x}_j}$$

$$L(\theta) = \prod_{i=1}^n \left( \frac{e^{\theta^T x_i}}{\sum_{j=1}^k e^{\theta^T x_j}} \right)^{I(y_i=j)} = \prod_{i=1}^n \left( \frac{e^{\theta_i^T x_i}}{e^{\theta_1^T x_i} + e^{\theta_2^T x_i} + e^{\theta_3^T x_i}} \right)^{I(y_i=j)}$$

$$\left( \frac{e^{\theta_i^T x_i}}{e^{\theta_1^T x_i} + e^{\theta_2^T x_i} + e^{\theta_3^T x_i}} \right)^{I(y_i=2)} \left( \frac{e^{\theta_i^T x_i}}{e^{\theta_1^T x_i} + e^{\theta_2^T x_i} + e^{\theta_3^T x_i}} \right)^{I(y_i=3)}$$

6. (10 points) Recall that the kernel trick is basically used to compute inner products between data points that are implicitly transformed into some higher dimensional space via a mapping  $\phi$ :

$$K(x, z) = \langle \phi(x), \phi(z) \rangle.$$

Let's assume that data points are 2-dimensional. So if  $K(x, y) = (1 + x^T z)^2$  is the quadratic kernel, write down the explicit form of the mapping  $\phi$  (i.e., write down what  $\phi(x)$  and  $\phi(z)$  are for this kernel).

$$\begin{aligned} K(x, z) &= (1 + x^T z)^2 \\ &= (1 + x^T z) \cdot (1 + x^T z) \\ &= (1 + z^T x) \cdot (1 + x^T z) \\ &= \phi(x) \cdot \phi(z) \end{aligned}$$

$$\phi(x) = (1 + z^T x)$$