

# CSE 6242 Assignment 3

Vincent La (Georgia Tech ID - vla6)

October 21, 2017

## Question 1: Data Preprocessing

## Question 2: Theory

**Part a: Write down the formula for the loss function used in Logistic Regression, the expression that you want to minimize:  $L(\theta)$**

Taken from the lecture “MLE and Iterative Optimization”:

$$\hat{\theta}_{MLE} = \operatorname{argmin}_{\theta} \sum_{i=1}^n \log(1 + e^{y^{(i)} * \langle \theta, x^{(i)} \rangle})$$

Thus, the loss function is

$$L(\theta) = \sum_{i=1}^n \log(1 + e^{y^{(i)} * \langle \theta, x^{(i)} \rangle})$$

where  $y^{(i)} = 1$  or  $y^{(i)} = -1$ .

**Part b: Derive the gradient of the loss function with respect to model parameters:  $\frac{dL(\theta)}{d\theta}$  or  $\frac{\partial L(\theta)}{\partial \theta_j}$ .**

$$\frac{\partial L(\theta)}{\partial \theta_j} = \frac{\partial \sum_{i=1}^n \log(1 + e^{y^{(i)} * \langle \theta, x^{(i)} \rangle})}{\partial \theta_j}$$

Furthermore, we know that  $\frac{\partial}{\partial x} \log(x) = \frac{1}{x}$ . Also, we can use the chain rule here to complete the derivative.

$$\begin{aligned} \frac{\partial L(\theta)}{\partial \theta_j} &= \sum_{i=1}^n \frac{1}{1 + e^{y^{(i)} * \langle \theta, x^{(i)} \rangle}} * \frac{\partial}{\partial \theta_j} e^{y^{(i)} * \langle \theta, x^{(i)} \rangle} \\ \frac{\partial L(\theta)}{\partial \theta_j} &= \sum_{i=1}^n \frac{e^{y^{(i)} * \langle \theta, x^{(i)} \rangle}}{1 + e^{y^{(i)} * \langle \theta, x^{(i)} \rangle}} * \frac{\partial}{\partial \theta_j} (y^{(i)} * \langle \theta, x^{(i)} \rangle) \\ \frac{\partial L(\theta)}{\partial \theta_j} &= \sum_{i=1}^n \frac{(e^{y^{(i)} * \langle \theta, x^{(i)} \rangle}) * y^{(i)} x_j^{(i)}}{1 + e^{y^{(i)} * \langle \theta, x^{(i)} \rangle}} \end{aligned}$$

**Part c: Based on this gradient, express the Stochastic Gradient Descent (SGD) update rule that uses a single sample  $\langle x^{(i)}, y^{(i)} \rangle$  at a time.**

Stochastic Gradient Descent (SGD) can be used when your data is very big. The steps are:

1. Initialize the dimensions of  $\theta$  vector to random values.
2. Pick one labeled data vector  $(x^{(i)}, y^{(i)})$  randomly, and update for each  $j = 1, \dots, d : \theta_j \leftarrow \theta_j - \alpha \frac{\partial \log(1 + \exp(y^{(i)} \langle \theta, x^{(i)} \rangle))}{\partial \theta_j}$
3. Repeat step (2) until the updates of the dimensions of  $\theta$  become too small.

So basically substituting in the expression for  $\frac{\partial}{\partial \theta_j} \log(1 + \exp(y^{(i)} \langle \theta, x^{(i)} \rangle))$  that we found previously, we get:

$$\frac{\partial}{\partial \theta_j} \log(1 + \exp(y^{(i)} \langle \theta, x^{(i)} \rangle)) = \frac{(e^{y^{(i)} \langle \theta, x^{(i)} \rangle}) * y^{(i)} x_j^{(i)}}{1 + e^{y^{(i)} \langle \theta, x^{(i)} \rangle}}$$

Thus, the update rule becomes,

for each  $j = 1, \dots, d : \theta_j \leftarrow \theta_j - \alpha * \frac{(e^{y^{(i)} \langle \theta, x^{(i)} \rangle}) * y^{(i)} x_j^{(i)}}{1 + e^{y^{(i)} \langle \theta, x^{(i)} \rangle}}$

**Part d: Write pseudocode for training a model using Logistic Regression and SGD.**

1. for (j from 1, ..., d), initialize  $\theta_j$  to random values. (Initialize the dimensions of  $\theta$  vector to random values.)
2. Pick one labeled data vector randomly, call it  $(x^{(i)}, y^{(i)})$
3. for (j from 1, ..., d) set  $\theta_j$  equal to  $\theta_j - \alpha * \frac{(e^{y^{(i)} \langle \theta, x^{(i)} \rangle}) * y^{(i)} x_j^{(i)}}{1 + e^{y^{(i)} \langle \theta, x^{(i)} \rangle}}$ , where  $\alpha$  is some step size, decaying as the gradient descent iterations increase.
4. Repeat step (3) until the updates of the dimensions of  $\theta$  become too small.

**Part e: Estimate the number of operations per epoch of SGD, where an epoch is one complete iteration through all the training samples. Express this in Big-O notation, in terms of the number of samples (n) and the dimensionality of each sample (d).**