Intro to Deep Learning          Yifan Li
HW 1                            yl506

Problem 1:          logistic regression (binary outcome)

a)

$$z = W_1 X_1 + W_2 X_2 + \cdots + W_m X_m + b$$

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

$$L(y, \sigma(z)) = -y \log \sigma(z) - (1-y) \log(1-\sigma(z))$$

$$\frac{\partial L}{\partial W_i} = \frac{\partial L}{\partial \sigma(z)} \frac{\partial \sigma(z)}{\partial z} \frac{\partial z}{\partial W_i} = \frac{\sigma(z)-y}{(1-\sigma(z))\sigma(z)} \cdot \sigma(z)(1-\sigma(z)) \cdot X_i = X_i(\sigma(z)-y)$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \sigma(z)} \frac{\partial \sigma(z)}{\partial z} \frac{\partial z}{\partial b} = \frac{\sigma(z)-y}{\sigma(z)(1-\sigma(z))} \cdot \sigma(z)(1-\sigma(z)) \cdot 1 = \sigma(z)-y$$

b)          MLP w/ a single hidden layer (binary outcome)

$$z_i = W_{1i} X_1 + W_{2i} X_2 + \cdots + W_{mi} X_m + b$$

$$\sigma(z_i) = \frac{1}{1+e^{-z_i}}$$

$$\zeta = V_1 \sigma(z_1) + V_2 \sigma(z_2) + \cdots + V_k \sigma(z_k) + c$$

$$\sigma(\zeta) = \frac{1}{1+e^{-\zeta}}$$

$$L(y, \sigma(\zeta)) = -y \log \sigma(\zeta) - (1-y) \log(1-\sigma(\zeta))$$

$$\frac{\partial L}{\partial V_i} = \frac{\partial L}{\partial \sigma(\zeta)} \frac{\partial \sigma(\zeta)}{\partial \zeta} \frac{\partial \zeta}{\partial V_i} = \frac{\sigma(\zeta)-y}{\sigma(\zeta)(1-\sigma(\zeta))} \cdot \sigma(\zeta)(1-\sigma(\zeta)) \cdot \sigma(z_i) = \sigma(z_i)(\sigma(\zeta)-y)$$
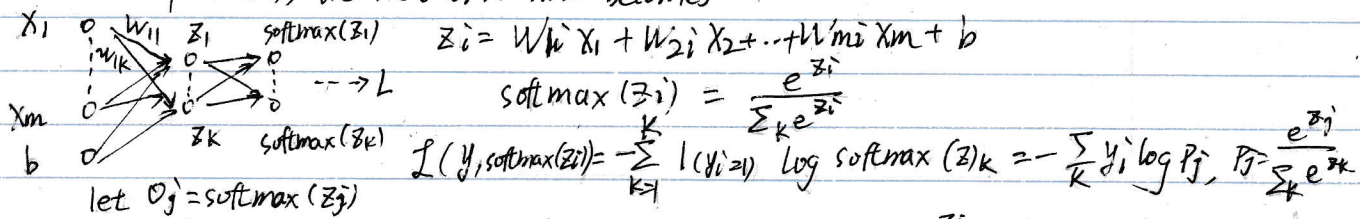
$$\frac{\partial L}{\partial c} = \frac{\partial L}{\partial \sigma(\zeta)} \frac{\partial \sigma(\zeta)}{\partial \zeta} \frac{\partial \zeta}{\partial c} = \sigma(\zeta)-y$$

$$\frac{\partial L}{\partial W_{ij}} = \frac{\partial L}{\partial \sigma(\zeta)} \frac{\partial \sigma(\zeta)}{\partial \zeta} \frac{\partial \zeta}{\partial \sigma(z_j)} \frac{\partial \sigma(z_j)}{\partial z_j} \frac{\partial z_j}{\partial W_{ij}} = \frac{\sigma(\zeta)-y}{\sigma(\zeta)(1-\sigma(\zeta))} \cdot \sigma(\zeta)(1-\sigma(\zeta)) \cdot V_j$$

$$\cdot \sigma(z_j)(1-\sigma(z_j)) \cdot Y_i = (\sigma(\zeta)-y) V_j \sigma(z_j)(1-\sigma(z_j)) X_i$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \sigma(\zeta)} \frac{\partial \sigma(\zeta)}{\partial \zeta} \frac{\partial \zeta}{\partial \sigma(z_j)} \frac{\partial \sigma(z_j)}{\partial z_j} \frac{\partial z_j}{\partial b} = (\sigma(\zeta)-y) V_j \sigma(z_j)(1-\sigma(z_j))$$

c)    When using a softmax (multi-class) setup
in part (a), the network now becomes

$z_i = W_{1i} x_1 + W_{2i} x_2 + \ldots + W_{mi} x_m + b$
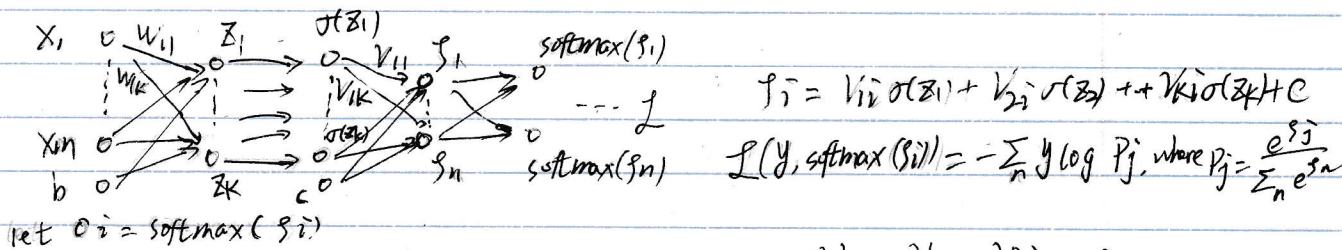
$\text{softmax}(z_i) = \dfrac{e^{z_i}}{\sum_k e^{z_i}}$

$L(y, \text{softmax}(z_i)) = -\sum_{k=1}^{K} I(y_i = 1) \log \text{softmax}(z)_k = -\sum_K y_i \log P_j, \quad P_j = \dfrac{e^{z_j}}{\sum_k e^{z_k}}$

let $O_j = \text{softmax}(z_j)$

$\cdot \ \dfrac{\partial L}{\partial W_{ij}} = \dfrac{\partial L}{\partial O_j} \dfrac{\partial O_j}{\partial z_j} \dfrac{\partial z_j}{\partial W_{ij}} = (P_j - y_j) x_i$, where $P_j = \dfrac{e^{z_j}}{\sum_k e^{z_k}}$

$\cdot \ \dfrac{\partial L}{\partial b} = P_j - y_j$     The loss function changes, so does $\dfrac{\partial L}{\partial W_{ij}}$ and $\dfrac{\partial L}{\partial b}$

in part (b), the network now becomes

$\hat{y}_i = V_{1i} \sigma(z_1) + V_{2i} \sigma(z_2) + \ldots + V_{ki} \sigma(z_k) + C$

$L(y, \text{softmax}(\hat{y}_i)) = -\sum_n y \log P_j$, where $P_j = \dfrac{e^{\hat{y}_j}}{\sum_n e^{\hat{y}_n}}$

let $O_i = \text{softmax}(\hat{y}_i)$

$\cdot \ \dfrac{\partial L}{\partial V_{ij}} = \dfrac{\partial L}{\partial O_j} \dfrac{\partial O_j}{\partial \hat{y}_j} \dfrac{\partial \hat{y}_j}{\partial V_{ij}} = (P_j - y_j) \sigma(z_i) \quad \cdot \ \dfrac{\partial L}{\partial C} = \dfrac{\partial L}{\partial O_j} \dfrac{\partial O_j}{\partial \hat{y}_j} \dfrac{\partial \hat{y}_j}{\partial C} = P_j - y_j$

$\cdot \ \dfrac{\partial L}{\partial W_{ij}} = \dfrac{\partial L}{\partial O_j} \dfrac{\partial O_j}{\partial \hat{y}_j} \dfrac{\partial \hat{y}_j}{\partial \sigma(z_i)} \dfrac{\partial \sigma(z_i)}{\partial z_i} \dfrac{\partial z_i}{\partial W_{ij}} = (P_j - y_j) V_{ij} \sigma(z_i)(1-\sigma(z_i)) x_i$

$\cdot \ \dfrac{\partial L}{\partial b} = \dfrac{\partial L}{\partial O_j} \dfrac{\partial O_j}{\partial \hat{y}_j} \dfrac{\partial \hat{y}_j}{\partial \sigma(z_i)} \dfrac{\partial \sigma(z_i)}{\partial z_i} \dfrac{\partial z_i}{\partial b} = (P_j - y_j) V_{ij} \sigma(z_i)(1-\sigma(z_i))$

The difference is that since the softmax is used for multi-class outcome,
a few terms of the original gradient change as well.

Problem 5.

a)    Yes.

b)    About 12 hours.

c)    I adhered to the Duke Community Standard in the completion of
        this assignment.