

Problem 1

Show that Eqn. (8) Lecture 3 holds, i.e., the partial derivative of $J(\beta)$ with respect to β_j is

$$\frac{\partial}{\partial \beta_j} J(\beta) = \frac{1}{N} \sum_{i=1}^N \left(\sigma(\beta^T x_i) - y_i \right) x_{ij}$$

Log loss:

$$J(\beta) = -\frac{1}{N} \sum_{i=1}^N [y_i \log(p_i) + (1-y_i) \log(1-p_i)] \quad (7)$$

$$\sigma(t) = \frac{1}{1 + \exp(-t)}$$

$$\sigma(x_i^T \beta) = \frac{1}{1 + \exp(-x_i^T \beta)} = p_i$$

$$\therefore \frac{\partial p_i}{\partial \beta_j} = \frac{\partial}{\partial \beta_j} (\sigma(x_i^T \beta))$$

$$= \sigma'(x_i^T \beta) \frac{\partial}{\partial \beta_j} (x_i^T \beta)$$

$$= \frac{\exp(-x_i^T \beta) x_{ij}}{(1 + \exp(-x_i^T \beta))^2}$$

$$= \sigma(x_i^T \beta) (1 - \sigma(x_i^T \beta)) x_{ij}$$

$$\frac{\partial}{\partial \beta_j} J(\beta) = -\frac{1}{N} \sum_{i=1}^N \left[y_i \frac{1}{p_i} \frac{\partial p_i}{\partial \beta_j} + (1-y_i) \frac{1}{1-p_i} \left(-\frac{\partial p_i}{\partial \beta_j} \right) \right]$$

$$= -\frac{1}{N} \sum_{i=1}^N \left[\frac{y_i}{p_i} + \frac{1-y_i}{1-p_i} \right] \frac{\partial p_i}{\partial \beta_j}$$

now we can substitute $p_i \approx \frac{\partial p_i}{\partial \beta_j}$:

$$\frac{\partial}{\partial \beta_j} J(\beta) = -\frac{1}{N} \sum_{i=1}^N \left[\frac{y_i}{\sigma(x_i^T \beta)} + \frac{1-y_i}{1-\sigma(x_i^T \beta)} \right] \sigma(x_i^T \beta) (1 - \sigma(x_i^T \beta)) x_{ij}$$

$$\begin{aligned}
&= -\frac{1}{N} \sum_{i=1}^N (y_i - \sigma(x_i^T \beta)) x_{ij} \\
&= \frac{1}{N} \sum_{i=1}^N (\sigma(\beta^T x_i) - y_i) x_{ij}
\end{aligned}$$

Problem 2

Consider the softmax regression. The number of data instances is N , and the number of classes is K .

(a) What is the value of $\sum_{k=1}^K y_k^{(i)}$, where $y_k^{(i)}$ is defined in Lecture 4.

(b) Let $t_k = s_k(\mathbf{x}_i)$. So p_k can be rewritten as:

$$p_k^{(i)} = \frac{\exp(t_k)}{\sum_{j=1}^K \exp(t_j)}$$

Let \bar{k} , $1 \leq \bar{k} \leq K$, be another index. Show that

$$\frac{1}{p_k^{(i)}} \frac{\partial p_k^{(i)}}{\partial t_{\bar{k}}} = \begin{cases} -p_{\bar{k}}^{(i)} & \text{if } k \neq \bar{k} \\ 1 - p_k^{(i)} = 1 - p_{\bar{k}}^{(i)} & \text{if } k = \bar{k} \end{cases}$$

(c) Show Eqn. (9) Lecture 4 holds, i.e.,

$$\nabla_{\beta^{(\bar{k})}} J(B) = \frac{1}{N} \sum_{i=1}^N (p_{\bar{k}}^{(i)} - y_{\bar{k}}^{(i)}) \mathbf{x}_i \quad (1)$$

(a) Lecture 4: $y_k^{(i)}$ is the target probability that the i th instance belongs to class k ; either 0 or 1. Since an i th instance can only belong to 1 class, if $y_k^{(i)} = 1$ then $y_j^{(i)} = 0$ where $j \neq k$. This summation thus will always add up to 1, because each instance can belong only to one class. $\sum_{k=1}^K y_k^{(i)} = 1$

(b) case 1: $k = \bar{k}$

$$\frac{\partial p_k^{(i)}}{\partial t_k} = \frac{\exp(t_k) \sum_{j=1}^K \exp(t_j) - \exp(t_k) \exp(t_k)}{\sum_{j=1}^K \exp(t_j)^2}$$

$$\partial t_k \quad \sum_{j=1}^K \exp(t_j)$$

$$\therefore \frac{\partial p_k^{(i)}}{\partial t_k} = p_k^{(i)} (1 - p_k^{(i)})$$

$$\text{Thus } \frac{1}{p_k^{(i)}} \frac{\partial p_k^{(i)}}{\partial t_k} = 1 - p_k^{(i)}$$

case 2 : $k \neq \bar{k}$

$$\frac{\partial p_k^{(i)}}{\partial t_{\bar{k}}} = - \frac{\exp(t_k) \exp(t_{\bar{k}})}{\left(\sum_{j=1}^K \exp(t_j)\right)^2} = -p_k^{(i)} p_{\bar{k}}^{(i)}$$

$$\therefore \frac{1}{p_k^{(i)}} \frac{\partial p_k^{(i)}}{\partial t_{\bar{k}}} = -p_{\bar{k}}^{(i)}$$

$$(c) \quad \text{given } \mathcal{J}(\beta) = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_k^{(i)} \log(p_k^{(i)})$$

$$\} \quad p_k^{(i)} = \frac{\exp(x_i^T \beta^{(k)})}{\sum_{j=1}^K \exp(x_i^T \beta^{(j)})}$$

$$\frac{\partial \mathcal{J}(\beta)}{\partial \beta^{(\bar{k})}} = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_k^{(i)} \frac{\partial}{\partial \beta^{(\bar{k})}} \log(p_k^{(i)})$$

$$\frac{\partial}{\partial \beta^{(\bar{k})}} \log(p_k^{(i)}) = \frac{1}{p_k^{(i)}} \frac{\partial p_k^{(i)}}{\partial \beta^{(\bar{k})}}$$

from softmax properties:

$$\frac{\partial p_k^{(i)}}{\partial t_{\bar{k}}} = \begin{cases} p_k^{(i)} (1 - p_k^{(i)}) & \text{if } k = \bar{k} \\ -p_k^{(i)} p_{\bar{k}}^{(i)} & k \neq \bar{k} \end{cases}$$

$$\frac{\partial t_{\bar{k}}}{\partial \beta^{(\bar{k})}} = x_i$$

$$\therefore \frac{\partial p_k^{(i)}}{\partial \beta^{(\bar{k})}} = \frac{\partial p_k^{(i)}}{\partial t_{\bar{k}}} \cdot \frac{\partial t_{\bar{k}}}{\partial \beta^{(\bar{k})}} = \frac{\partial p_k^{(i)}}{\partial t_{\bar{k}}} \cdot x_i$$

$$\frac{\partial \mathcal{J}(\beta)}{\partial \beta^{(k)}} = -\frac{1}{N} \sum_{i=1}^N (y_i^{(k)} - p_i^{(k)}) x_i$$

now we substitute

$$\frac{\partial \mathcal{J}(\beta)}{\partial \beta^{(k)}} = -\frac{1}{N} \sum_{i=1}^N (y_i^{(k)} - p_i^{(k)}) x_i$$

$$\therefore \nabla_{\beta^{(k)}} \mathcal{J}(\beta) = -\frac{1}{N} \sum_{i=1}^N (p_i^{(k)} - y_i^{(k)}) x_i$$

Problem 3 (MATH 5388 Only)

Show the cost function of softmax regression (Eqn. 8 Lecture 4) is equivalent to that of the logistic regression (Eqn. 7 lecture 3), if there are only 2 classes: 1 and 0. Hint: suppose the coefficient for the first class 1 is $\beta^{(1)}$ and the coefficient for the second class 0 is $\beta^{(2)}$. Consider consolidating the coefficients to a single parameter $\beta = \beta^{(1)} - \beta^{(2)}$.

softmax regression for k classes:

$$\mathcal{J}(\beta) = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_i^{(k)} \log(p_i^{(k)})$$

we are told only 2 classes

$$\therefore K = 2$$

$$\mathcal{J}(\beta) = -\frac{1}{N} \sum_{i=1}^N [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$$

$$p_i^{(1)} = \frac{\exp(\beta^{(1)T} x_i)}{\exp(\beta^{(1)T} x_i) + \exp(\beta^{(2)T} x_i)}$$


$$p_i^{(2)} = \frac{\exp(\beta^{(2)T} x_i)}{\exp(\beta^{(1)T} x_i) + \exp(\beta^{(2)T} x_i)} = 1 - p_i^{(1)}$$

$$\text{Hint: } \beta = \beta^{(1)} - \beta^{(2)}$$

$$\therefore p_1^{(i)} = \frac{\exp(\beta^T x_i)}{\exp(\beta^T x_i) + 1}$$

$$p_2^{(i)} = \frac{1}{\exp(\beta^T x_i) + 1} = 1 - p_1^{(i)}$$

$$\begin{aligned} J(\beta) &= -\frac{1}{N} \sum_{i=1}^N [y_1^{(i)} \log(p_1^{(i)}) + y_2^{(i)} \log(p_2^{(i)})] \\ &= -\frac{1}{N} \sum_{i=1}^N \left[y_i \log\left(\frac{\exp(\beta^T x_i)}{\exp(\beta^T x_i) + 1}\right) + (1 - y_i) \log\left(\frac{1}{\exp(\beta^T x_i) + 1}\right) \right] \end{aligned}$$


 same as that of
 logistic regression