

HW2

Ibrahim Gabr

Problem 1: Softmax

Let us start by providing the definition of the softmax function for the binary case - that is $C = 2$, we

$$\text{know that } \hat{p}(y = c_1 | \mathbf{x}; \mathbf{W}) = \frac{\exp(\mathbf{w}_{c_1} \cdot \mathbf{x})}{\sum_{y=1}^2 \exp(\mathbf{w}_y \cdot \mathbf{x})} \text{ and } \hat{p}(y = c_2 | \mathbf{x}; \mathbf{W}) = \frac{\exp(\mathbf{w}_{c_2} \cdot \mathbf{x})}{\sum_{y=1}^2 \exp(\mathbf{w}_y \cdot \mathbf{x})}.$$

We can now compute the log-odds for these two classes as follows:

$$\begin{aligned} \log\left(\frac{\hat{p}(y = c_1 | \mathbf{x}; \mathbf{W})}{\hat{p}(y = c_2 | \mathbf{x}; \mathbf{W})}\right) &= \log(\hat{p}(y = c_1 | \mathbf{x}; \mathbf{W})) - \log(\hat{p}(y = c_2 | \mathbf{x}; \mathbf{W})) \\ &= \log\left(\frac{\exp(\mathbf{w}_{c_1} \cdot \mathbf{x})}{\sum_{y=1}^2 \exp(\mathbf{w}_y \cdot \mathbf{x})}\right) - \log\left(\frac{\exp(\mathbf{w}_{c_2} \cdot \mathbf{x})}{\sum_{y=1}^2 \exp(\mathbf{w}_y \cdot \mathbf{x})}\right) \\ &= \mathbf{w}_{c_1} \cdot \mathbf{x} - \log\left(\sum_{y=1}^2 \exp(\mathbf{w}_y \cdot \mathbf{x})\right) - \mathbf{w}_{c_2} \cdot \mathbf{x} + \log\left(\sum_{y=1}^2 \exp(\mathbf{w}_y \cdot \mathbf{x})\right) \\ &= \mathbf{w}_{c_1} \cdot \mathbf{x} - \mathbf{w}_{c_2} \cdot \mathbf{x} \\ &= \mathbf{w}_{c_1}^T \mathbf{x} - \mathbf{w}_{c_2}^T \mathbf{x} \\ &= (\mathbf{w}_{c_1}^T - \mathbf{w}_{c_2}^T) \mathbf{x} \\ &= \mathbf{v} \cdot \mathbf{x} \quad \text{where } \mathbf{v} = \mathbf{w}_{c_1}^T - \mathbf{w}_{c_2}^T \end{aligned}$$

Note: $\mathbf{w} \cdot \mathbf{x} = \mathbf{w}^T \mathbf{x}$.

The above shows that the log-odds expression for the $C = 2$ corresponds to a linear function. We can now show that we can take this linear model at arrive back at our definition of the softmax function.

$$\log \left(\frac{\hat{p}(y = c_1 | \mathbf{x}; \mathbf{W})}{\hat{p}(y = c_2 | \mathbf{x}; \mathbf{W})} \right) = \mathbf{w}_{c_1} \cdot \mathbf{x} - \mathbf{w}_{c_2} \cdot \mathbf{x}$$

$$\frac{\hat{p}(y = c_1 | \mathbf{x}; \mathbf{W})}{\hat{p}(y = c_2 | \mathbf{x}; \mathbf{W})} = \exp(\mathbf{w}_{c_1} \cdot \mathbf{x} - \mathbf{w}_{c_2} \cdot \mathbf{x})$$

$$\frac{\hat{p}(y = c_1 | \mathbf{x}; \mathbf{W})}{1 - \hat{p}(y = c_1 | \mathbf{x}; \mathbf{W})} = \exp(\mathbf{w}_{c_1} \cdot \mathbf{x} - \mathbf{w}_{c_2} \cdot \mathbf{x})$$

$$\frac{1 - \hat{p}(y = c_1 | \mathbf{x}; \mathbf{W})}{\hat{p}(y = c_1 | \mathbf{x}; \mathbf{W})} = \frac{1}{\exp(\mathbf{w}_{c_1} \cdot \mathbf{x} - \mathbf{w}_{c_2} \cdot \mathbf{x})}$$

$$\frac{1}{\hat{p}(y = c_1 | \mathbf{x}; \mathbf{W})} = 1 + \frac{1}{\exp(\mathbf{w}_{c_1} \cdot \mathbf{x} - \mathbf{w}_{c_2} \cdot \mathbf{x})}$$

$$\frac{1}{\hat{p}(y = c_1 | \mathbf{x}; \mathbf{W})} = \frac{1 + \exp(\mathbf{w}_{c_1} \cdot \mathbf{x} - \mathbf{w}_{c_2} \cdot \mathbf{x})}{\exp(\mathbf{w}_{c_1} \cdot \mathbf{x} - \mathbf{w}_{c_2} \cdot \mathbf{x})}$$

$$\hat{p}(y = c_1 | \mathbf{x}; \mathbf{W}) = \frac{\exp(\mathbf{w}_{c_1} \cdot \mathbf{x} - \mathbf{w}_{c_2} \cdot \mathbf{x})}{1 + \exp(\mathbf{w}_{c_1} \cdot \mathbf{x} - \mathbf{w}_{c_2} \cdot \mathbf{x})}$$

$$\hat{p}(y = c_1 | \mathbf{x}; \mathbf{W}) = \frac{\exp(\mathbf{w}_{c_1} \cdot \mathbf{x}) \exp(-\mathbf{w}_{c_2} \cdot \mathbf{x})}{1 + \exp(\mathbf{w}_{c_1} \cdot \mathbf{x}) \exp(-\mathbf{w}_{c_2} \cdot \mathbf{x})}$$

$$\hat{p}(y = c_1 | \mathbf{x}; \mathbf{W}) = \frac{\exp(\mathbf{w}_{c_1} \cdot \mathbf{x})}{\exp(\mathbf{w}_{c_1} \cdot \mathbf{x}) + \exp(\mathbf{w}_{c_2} \cdot \mathbf{x})}$$

Problem 2

To show that the softmax model as stated in (1) is *overparametrized* we can write the probabilities as follows:

$$p(y=1|\mathbf{x}; \mathbf{W}) = \frac{\exp(\mathbf{w}'_1)}{1 + \sum_{y=1}^{C-1} \exp(\mathbf{w}'_y \cdot \mathbf{x})}$$

where $\exp(\mathbf{w}'_i) = \exp(\mathbf{w}_i)\exp(-\mathbf{w}_C)$ $\forall i \neq C$ and

$$\exp(\mathbf{w}'_C) = \exp(\mathbf{w}_C)\exp(-\mathbf{w}_C) = 1$$

$$p(y=2|\mathbf{x}; \mathbf{W}) = \frac{\exp(\mathbf{w}'_2)}{1 + \sum_{y=1}^{C-1} \exp(\mathbf{w}'_y \cdot \mathbf{x})}$$

⋮

$$p(y=C-1|\mathbf{x}; \mathbf{W}) = \frac{\exp(\mathbf{w}'_{C-1})}{1 + \sum_{y=1}^{C-1} \exp(\mathbf{w}'_y \cdot \mathbf{x})}$$

Using the fact that the probabilities from 1 to C must sum to one, we get:

$$\begin{aligned} p(y=C|\mathbf{x}; \mathbf{W}) &= 1 - \sum_{c=1}^{C-1} p(y=c|\mathbf{x}; \mathbf{W}) \\ &= 1 - \frac{\sum_{c=1}^{C-1} \exp(\mathbf{w}'_c \cdot \mathbf{x})}{1 + \sum_{y=1}^{C-1} \exp(\mathbf{w}'_y \cdot \mathbf{x})} \\ &= \frac{1}{1 + \sum_{c=1}^{C-1} \exp(\mathbf{w}'_c \cdot \mathbf{x})} \end{aligned}$$

This expression shows that we can write $p(y=c|\mathbf{x}; \mathbf{W})$ for any c by using $C-1$ parameter vectors, where we interpret each vector parameter as a difference with respect to our base category. We usually select the maximum parameter value to subtract in order to assist with the computation.

Problem 3

Let us begin by looking at the solution for MLE in the case of logistic regression. This will be critical in proving the hessian matrix \mathbf{H} is positive definite. That is, all entries on the diagonal are strictly positive.

Let us also define the following quantity: $\gamma_i \equiv \frac{1}{1 + \exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)} = \sigma(w_0 + \mathbf{w} \cdot \mathbf{x}_i)$.

$$\operatorname{argmin}_{\mathbf{w}} \log p(\mathbf{y} | \mathbf{X}; \mathbf{w}) = - \sum_{i=1}^d y_i \log \sigma(w_0 + \mathbf{w} \cdot \mathbf{x}_i) + (1 - y_i) \log(1 - \sigma(w_0 + \mathbf{w} \cdot \mathbf{x}_i))$$

$$\text{where } \sigma(w_0 + \mathbf{w} \cdot \mathbf{x}_i) = \frac{1}{1 + \exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)}$$

Let us now compute the derivatives with respect to w_0 and w_j :

$$\begin{aligned} \frac{\delta}{\delta w_0} \log \sigma(w_0 + \mathbf{w} \cdot \mathbf{x}_i) &= \frac{1}{\sigma(w_0 + \mathbf{w} \cdot \mathbf{x}_i)} \frac{\delta}{\delta w_0} \sigma(w_0 + \mathbf{w} \cdot \mathbf{x}_i) \\ &= \frac{(1 + \exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i))^{-2} \exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)}{(1 + \exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i))^{-1}} \\ &= \frac{\exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)}{1 + \exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)} \\ \frac{\delta}{\delta w_j} \log \sigma(w_0 + \mathbf{w} \cdot \mathbf{x}_i) &= \frac{1}{\sigma(w_0 + \mathbf{w} \cdot \mathbf{x}_i)} \frac{\delta}{\delta w_j} \sigma(w_0 + \mathbf{w} \cdot \mathbf{x}_i) \\ &= \frac{\exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i) x_{ij}}{1 + \exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)} \end{aligned}$$

And again for the remaining part of the expression:

$$\begin{aligned} \frac{\delta}{\delta w_0} \log(1 - \sigma(w_0 + \mathbf{w} \cdot \mathbf{x}_i)) &= \frac{\delta}{\delta w_0} \log \left(\frac{\exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)}{1 + \exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)} \right) \\ &= \frac{\delta}{\delta w_0} \left[-w_0 - \mathbf{w} \cdot \mathbf{x}_i - \log(1 + \exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)) \right] \\ &= -1 - \frac{\exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)}{1 + \exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)} \\ &= -\frac{1}{1 + \exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)} \\ &= -\gamma_i \end{aligned}$$

$$\begin{aligned}
\frac{\delta}{\delta w_j} \log(1 - \sigma(w_0 + \mathbf{w} \cdot \mathbf{x}_i)) &= \frac{\delta}{\delta w_j} \log \left(\frac{\exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)}{1 + \exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)} \right) \\
&= \frac{\delta}{\delta w_j} \left[-w_0 - \mathbf{w} \cdot \mathbf{x}_i - \log(1 + \exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)) \right] \\
&= -x_{ij} - \frac{\exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)}{1 + \exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)} x_{ij} \\
&= -\frac{x_{ij}}{1 + \exp(-w_0 - \mathbf{w} \cdot \mathbf{x}_i)} \\
&= -\gamma_i x_{ij}
\end{aligned}$$

Putting this all together, we get:

$$\begin{aligned}
\frac{\delta}{\delta w_0} \log p(\mathbf{y} | \mathbf{X}; \mathbf{w}) &= - \sum_{i=1}^d y_i \frac{\delta}{\delta w_0} \log \sigma(w_0 + \mathbf{w} \cdot \mathbf{x}_i) + (1 - y_i) \frac{\delta}{\delta w_0} \log(1 - \sigma(w_0 + \mathbf{w} \cdot \mathbf{x}_i)) = 0 \\
&= - \sum_{i=1}^d y_i (1 - \gamma_i) + (1 - y_i)(-\gamma_i) = 0 \\
&= - \sum_{i=1}^d y_i - \gamma_i = 0 \\
&= \sum_{i=1}^d \gamma_i - y_i = 0 \\
\\
\frac{\delta}{\delta w_j} \log p(\mathbf{y} | \mathbf{X}; \mathbf{w}) &= - \sum_{i=1}^d y_i \frac{\delta}{\delta w_j} \log \sigma(w_0 + \mathbf{w} \cdot \mathbf{x}_i) + (1 - y_i) \frac{\delta}{\delta w_j} \log(1 - \sigma(w_0 + \mathbf{w} \cdot \mathbf{x}_i)) = 0 \\
&= - \sum_{i=1}^d (y_i - \gamma_i) x_{ij} = 0 \\
&= \sum_{i=1}^d (\gamma_i - y_i) x_{ij} = 0
\end{aligned}$$

In matrix form, we can express the above as the following:

$$\boldsymbol{\gamma} = \begin{pmatrix} \gamma_1 \\ \gamma_2 \\ \gamma_3 \\ \vdots \\ \gamma_d \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_d \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ x_{21} & x_{22} & \cdots & x_{2d} \\ x_{31} & x_{32} & \cdots & x_{3d} \\ \vdots & \ddots & \ddots & \vdots \\ x_{d1} & x_{d2} & \cdots & x_{dd} \end{pmatrix}$$

Therefore, the gradient is $\nabla L = \frac{\delta}{\delta \mathbf{w}} \log p(\mathbf{y}|\mathbf{X}; \mathbf{w}) = (\boldsymbol{\gamma} - \mathbf{y})^T \mathbf{X}$.

With some manipulation, we know that this is equal to $\nabla L = \mathbf{X}^T (\boldsymbol{\gamma} - \mathbf{y})$.

In order to create \mathbf{H} , we must now calculate second order derivatives where $w_k \neq w_j$. From the above calculations, it can be shown that:

$$\frac{\delta}{\delta w_k} \gamma_i = \gamma_i(1 - \gamma_i)x_{ik}$$

Thus:

$$\begin{aligned} \frac{\delta^2}{\delta w_j \delta w_k} \log p(\mathbf{y}|\mathbf{X}; \mathbf{w}) &= \sum_{i=1}^d \frac{\delta}{\delta w_k} \gamma_i x_{ij} \\ &= \sum_{i=1}^d \gamma_i(1 - \gamma_i)x_{ik}x_{ij} \\ &= \mathbf{a}_k^T \mathbf{R} \mathbf{a}_j > 0 \end{aligned}$$

where $\mathbf{a}_k = [x_{1k}, x_{2k}, \dots, x_{dk}]^T$,

$$\mathbf{R} = \begin{pmatrix} \gamma_1(1 - \gamma_1) & 0 & 0 & \cdots & 0 \\ 0 & \gamma_2(1 - \gamma_2) & 0 & \cdots & 0 \\ 0 & 0 & \gamma_3(1 - \gamma_3) & \cdots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \gamma_d(1 - \gamma_d) \end{pmatrix}$$

Here, \mathbf{a}_k and \mathbf{a}_j represent columns of \mathbf{X} . In particular, the expression $\mathbf{a}_k^T \mathbf{R} \mathbf{a}_j$ gives us the derivative for the kth-jth entry.

As such, we can express the hessian as $\mathbf{H} = \mathbf{X}^T \mathbf{R} \mathbf{X}$.

While the question states that the Hessian should be a $d \times d$ matrix, it is very simple to incorporate the constant:

$$\mathbf{H}_{(d+1) \times (d+1)} = \mathbf{X}_{(d+1) \times d}^T \times \mathbf{R}_{d \times d} \times \mathbf{X}_{d \times (d+1)}$$

Theorems from Linear Algebra inform us that a matrix \mathbf{H} is positive definite if it is non-singular ($\det(\mathbf{H}) \neq 0$), can be expressed as $\mathbf{H} = \mathbf{A}^T \mathbf{A}$ in addition to the columns of \mathbf{A} being linearly independent.

Let us first look at the bounds of the entries in R.

From our definition of γ_i above, we can infer that:

- $0 \leq \gamma_i \leq 1$

- $0 \leq (1 - \gamma_i) \leq 1$

As a consequence, we can infer that:

- $0 \leq \gamma_i(1 - \gamma_i) \leq 1$

Note that I have written $\mathbf{0} \leq$. We know that a positive definite matrix must have eigenvalues $> \mathbf{0}$. The fact that \mathbf{R} has linearly independent columns will ensure this. As such, all that is left is to show that $\mathbf{R}^T \mathbf{R} = \mathbf{H}$.

Since \mathbf{R} has linearly independent columns, we can go ahead and take the square root of all the elements in \mathbf{R} . Since \mathbf{R} is a diagonal matrix, it is important to note that $\mathbf{R}^{1/2} = (\mathbf{R}^{1/2})^T$.

We can now write \mathbf{H} as:

$$\begin{aligned}\mathbf{H} &= \mathbf{X}^T \mathbf{R}^{1/2} \mathbf{R}^{1/2} \mathbf{X} \quad \text{since } \mathbf{R}^{1/2} \text{ is a diagonal matrix it is true that } \mathbf{R}^{1/2} = (\mathbf{R}^{1/2})^T \\ &= (\mathbf{R}^{1/2} \mathbf{X})^T (\mathbf{R}^{1/2} \mathbf{X})\end{aligned}$$

$\mathbf{H} = \mathbf{X}^T \mathbf{R}^{1/2} \mathbf{R}^{1/2} \mathbf{X} = (\mathbf{R}^{1/2} \mathbf{X})^T (\mathbf{R}^{1/2} \mathbf{X})$. If we were to replace $\mathbf{R}^{1/2} \mathbf{X}$ with \mathbf{A} , we see that $\mathbf{H} = \mathbf{A}^T \mathbf{A}$.

Thus, \mathbf{H} is a positive definite matrix.

Using the Newton-Raphson method, we can show convexity in \mathbf{w} and thus a unique minimum.

$$\begin{aligned}\mathbf{w}_{t+1} &= \mathbf{w}_t + \mathbf{H}^{-1} \frac{\delta}{\delta \mathbf{w}} \log p(\mathbf{X}; \mathbf{w}) \\ &= \mathbf{w}_t + (\mathbf{X}^T \mathbf{R} \mathbf{X})^{-1} \mathbf{X}^T (\boldsymbol{\gamma} - \mathbf{y}) \\ &= (\mathbf{X}^T \mathbf{R} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{R} (\mathbf{X} \mathbf{w}_t + \mathbf{R}^{-1} (\boldsymbol{\gamma} - \mathbf{y})) \\ &= (\mathbf{X}^T \mathbf{R} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{R} \mathbf{v}_t \\ \text{where } \mathbf{v}_t &= (\mathbf{X} \mathbf{w}_t + \mathbf{R}^{-1} (\boldsymbol{\gamma} - \mathbf{y}))\end{aligned}$$

This looks almost identical to the solution to least squares. In our case, the optimal solution is given by:

$$\underset{\mathbf{w}}{\operatorname{argmin}} = \sum_{i=1}^d \gamma_i (1 - \gamma_i) (v_i - \mathbf{w} \cdot \mathbf{x}_i)^2$$

It can be shown that for each iteration, we are calculating the global minimum since the loss function is convex.

Problem 4

NOTE: This response was a joint effort between all contributors.

In this problem, we will represent the label for \mathbf{y}_i by an indicator vector \mathbf{t}_i . For example, if \mathbf{t}_1 refers to $c = 1$, then $\mathbf{t}_1^T = [1, 0, \dots, 0]$.

Therefore, \mathbf{t}_i represents a basis vector for a class c , containing a 1 at the j th position and 0 elsewhere. These vectors will be useful later to compute the derivatives.

Before we continue, let's define $\hat{\mathbf{p}}_i$ as a vector of probabilities for the i th row as $\hat{\mathbf{p}}_i = \hat{\mathbf{p}}(\mathbf{y}_i | \mathbf{x}_i; \mathbf{W})$ and $\mathbf{a}_i = \mathbf{w}_i \cdot \mathbf{x}_i$. This last definition allows us to rewrite the softmax model as follows:

$$\hat{\mathbf{p}}(\mathbf{y}_i = c | \mathbf{x}_i; \mathbf{W}) = \frac{\exp(a_c)}{\sum_{y=1}^C \exp(a_y)}$$

We have to rewrite now our cost function in terms of the log-loss as follows

$$\begin{aligned} \mathbf{W}^* = \underset{\mathbf{W}}{\operatorname{argmin}} J(\mathbf{x}_i, \mathbf{y}_i, \mathbf{W}) &= -\frac{1}{N} \sum_{i=1}^N t_i \log \hat{\mathbf{p}}(\mathbf{y}_i | \mathbf{x}_i; \mathbf{W}) + \lambda \sum_i \sum_j \mathbf{w}_{ij}^2 \\ &= -\frac{1}{N} \sum_{i=1}^N t_i \log \hat{\mathbf{p}}_i + \lambda \sum_i \sum_j \mathbf{w}_{ij}^2 \end{aligned}$$

First, let's take the derivative of the log-loss function J with respect to $\hat{\mathbf{p}}_i$.

$$(1) \quad \frac{\delta J}{\delta \hat{\mathbf{p}}_i} = -\frac{1}{N} \frac{t_i}{\hat{\mathbf{p}}_i}$$

$$(2) \quad \frac{\delta \hat{\mathbf{p}}_i}{\delta a_k} = \begin{cases} \frac{\exp(a_i)}{\sum_{y=1}^C \exp(a_y)} - \left(\frac{\exp(a_i)}{\sum_{y=1}^C \exp(a_y)} \right)^2 & \text{if } i = k \\ -\frac{\exp(a_i) \exp(a_k)}{\left(\sum_{y=1}^C \exp(a_y) \right)^2} & \text{if } i \neq k \end{cases}$$

$$= \begin{cases} \hat{\mathbf{p}}_i (1 - \hat{\mathbf{p}}_i) & \text{if } i = k \\ \hat{\mathbf{p}}_i \hat{\mathbf{p}}_k & \text{if } i \neq k \end{cases}$$

$$(3) \quad \frac{\delta J}{\delta a_i} = \sum_{k=1}^C \frac{\delta J}{\delta \hat{\mathbf{p}}_k} \frac{\delta \hat{\mathbf{p}}_k}{\delta a_i}$$

$$= \frac{\delta J}{\delta \hat{\mathbf{p}}_i} \frac{\delta \hat{\mathbf{p}}_i}{\delta a_i} - \sum_{k \neq i} \frac{\delta J}{\delta \hat{\mathbf{p}}_k} \frac{\delta \hat{\mathbf{p}}_k}{\delta a_i}$$

$$= -\frac{1}{N} \frac{t_i}{\hat{\mathbf{p}}_i} \hat{\mathbf{p}}_i (1 - \hat{\mathbf{p}}_i) + \frac{1}{N} \sum_{k \neq i} \frac{t_k}{\hat{\mathbf{p}}_k} \hat{\mathbf{p}}_k \hat{\mathbf{p}}_i$$

$$\begin{aligned}
&= -\frac{1}{N} t_i (1 - \hat{p}_i) + \frac{1}{N} \sum_{k \neq i} t_k \hat{p}_i \\
&= -\frac{1}{N} t_i (1 - \hat{p}_i) + \frac{1}{N} \hat{p}_i \sum_{k \neq i} t_k \\
&= \frac{1}{N} \left[\hat{p}_i \sum_{k \neq i} t_k - t_i (1 - \hat{p}_i) \right] \\
&= \frac{1}{N} \left[\hat{p}_i \left(\sum_{k \neq i} t_k + t_i \right) - t_i \right] \\
&= \frac{1}{N} \left[\hat{p}_i \left(\sum_k t_k \right) - t_i \right] \quad \text{where } \sum_k t_k = 1 \\
&= \frac{1}{N} \left[\hat{p}_i - t_i \right]
\end{aligned}$$

$$(4) \quad \frac{\delta J}{\delta w_{ij}} = \sum_{i=1}^N \frac{\delta J}{\delta a_i} \frac{\delta a_i}{\delta w_{ij}} + \lambda \sum_i \sum_j \frac{\delta}{\delta w_{ij}} w_{ij}^2$$

$$\begin{aligned}
&= \frac{1}{N} \sum_{i=1}^N \left[\hat{p}_i - t_i \right] \frac{\delta a_i}{\delta w_{ij}} + 2\lambda w_{ij} \\
&= \frac{1}{N} \sum_{i=1}^N \left[\hat{p}_i - t_i \right] x_{ij} + 2\lambda w_{ij} \\
&= \frac{1}{N} \sum_{i=1}^N \left[\hat{p}_i - t_i \right] x_{ij} + 2\lambda \sum_{i=1}^N w_{ij} \\
&= \frac{1}{N} (\mathbf{X}^T (\hat{\mathbf{p}} - \mathbf{t})) + 2\lambda \mathbf{W} \mathbf{t}_i \\
&= \frac{1}{N} (\mathbf{X}^T (\hat{\mathbf{p}} - \mathbf{t})) + 2\lambda \mathbf{w}_i
\end{aligned}$$

$$\hat{\mathbf{p}} = \begin{pmatrix} \hat{p}_{11} & \hat{p}_{12} & \cdots & \hat{p}_{1C} \\ \hat{p}_{21} & \hat{p}_{22} & \cdots & \hat{p}_{2C} \\ \hat{p}_{31} & \hat{p}_{32} & \cdots & \hat{p}_{3C} \\ \vdots & \ddots & \ddots & \vdots \\ \hat{p}_{N1} & \hat{p}_{N2} & \cdots & \hat{p}_{NC} \end{pmatrix}, \quad \mathbf{t} = \begin{pmatrix} t_{11} & t_{12} & \cdots & t_{1C} \\ t_{21} & t_{22} & \cdots & t_{2C} \\ t_{31} & t_{32} & \cdots & t_{3C} \\ \vdots & \ddots & \ddots & \vdots \\ t_{N1} & t_{N2} & \cdots & t_{NC} \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1N} \\ x_{21} & x_{22} & \cdots & x_{2N} \\ x_{31} & x_{32} & \cdots & x_{3N} \\ \vdots & \ddots & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{NN} \end{pmatrix}$$

$$\mathbf{W} = \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1N} \\ w_{21} & w_{22} & \cdots & w_{2N} \\ w_{31} & w_{32} & \cdots & w_{3N} \\ \vdots & \ddots & \ddots & \vdots \\ w_{N1} & w_{N2} & \cdots & w_{NN} \end{pmatrix}$$

Therefore, we can write the equation for the stochastic gradient descent as follows

$$\mathbf{w}_{t+1} = \eta_t \frac{1}{N} (\mathbf{X}^T (\hat{\mathbf{p}} - \mathbf{t})) + 2\lambda \mathbf{w}_t$$