

# Machine Learning Methodologies and Applications (AI6012)

## Individual Assignment

Leon Sun (G2204908A)

30 September 2023

### 1 Question 1 (10 marks)

Multi-class classification, or Multinomial Logistic Regression, can be approached using softmax regression. The softmax function is generally defined as:

$$P(y = c | x) = \frac{\exp(w^{(c)}x)}{\sum_{i=1}^C \exp(w^{(i)}x)} \text{ or } \frac{1}{1 + \sum_{i \neq C} \exp(w^{(i)}x)} \quad (1)$$

Since the sum of all the conditional probabilities for the softmax is 1, we can summarise the probabilities for all classes to:

$$\sum_{c=0}^C P(y = c | x) = 1 \quad (2)$$

By introducing the set of logits into we can arrived at the following parametric equations for multinomial logistic regression. Suppose there are C classes, 0, 1, ..., C-1:

$$\text{For } c > 0 : P(y = c | x) = \frac{\exp(-w^{(c)^T}x)}{1 + \sum_{c=1}^{C-1} \exp(-w^{(c)^T}x)} = \hat{y}_c \quad (3)$$

$$\text{For } c > 0 : P(y = 0 | x) = \frac{1}{1 + \sum_{c=1}^{C-1} \exp(-w^{(c)^T}x)} = \hat{y}_0 \quad (4)$$

Given a set of N training input-output pairs like  $x_i, y_i, i = 1, \dots, N$ , which are i.i.d, we can define the likelihood as the product of likelihoods of each individual pairs.

$$\mathcal{L}(w_c) = \prod_{i=1}^N l(w_c | \{x_i, y_i\}) = \prod_{i=1}^N P(y_i | x_i; w_c) \quad (5)$$

Hence the maximum likelihood estimation can be represented in the following ln function which converts the product into a sum.

$$\hat{w}_c = \operatorname{argmax}_{w_c} \prod_{i=1}^N P(y_i | x_i; w_c) = \operatorname{argmax}_{w_c} \sum_{i=1}^N \sum_{c=0}^{C-1} (y_i \ln(g(x_i; w_c))) \quad (6)$$

$$\ln \hat{w} = -\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^{C-1} y_i \ln(P(y_i | x_i; w_c)) \quad (7)$$

We will now derive the learning procedure for the multinomial logistic classification using Gradient Descent optimisation method.

$$w_{t+1} = w_t - \rho \frac{\partial E(w)}{\partial w} \quad (8)$$

$$\begin{aligned} \frac{\partial E(w)}{\partial w} &= \frac{\partial \left( -\sum_{i=1}^N \sum_{c=1}^C y_i \cdot \ln(P(y_i | x_i; w_c)) \right)}{\partial w_c} \\ &= -\sum_{i=1}^N \sum_{c=1}^C \frac{\partial (y_i \cdot \ln(p(y_i | x_i; w_c)))}{\partial w_c} \end{aligned}$$

Let  $P(y_i | x_i; w_c)$  be  $f(z)$ . Using chain rule:

$$\begin{aligned}\frac{\partial \ln f(z)}{\partial z} &= \frac{\partial \ln f(z)}{\partial f(z)} \cdot \frac{\partial f(z)}{\partial z} \\ &= \frac{1}{f(z)} \cdot \frac{df(z)}{\partial z}\end{aligned}\tag{9}$$

We can differentiate using:

$$\frac{\partial (y_i \cdot \ln(P(y_i | x_i; w_c)))}{\partial w_c} = y_i \cdot \frac{1}{P(y_i | x_i; w_c)} \cdot \frac{\partial (P(y_i | x_i; w_c))}{\partial w_c}\tag{10}$$

If  $y = c$ , subbing equations (3) into (10) gives us:

$$\frac{\partial (P(y_i | x_i; w_c))}{\partial w_c} = \frac{\partial}{\partial w_c} \left( \frac{\exp(-w_c \cdot x_i)}{1 + \sum_{c=1}^{c-1} \exp(-w_c \cdot x_i)} \right)$$

Using quotient rule:

$$\begin{aligned}& \frac{\partial}{\partial w_c} \left( \frac{\exp(-w_c \cdot x_i)}{1 + \sum_{c=1}^{c-1} \exp(-w_c \cdot x_i)} \right) \\ &= \frac{\left( \frac{\partial}{\partial w_c} \exp(-w_c \cdot x_i) \right) \left( 1 + \sum_{c=1}^{c-1} \exp(-w_c \cdot x_i) \right) + \left( \frac{\partial}{\partial w_c} \left( 1 + \sum_{c=1}^{c-1} \exp(-w_c \cdot x_i) \right) \right) (\exp(-w_c \cdot x_i))}{\left( 1 + \sum_{c=1}^{c-1} \exp(-w_c \cdot x_i) \right)^2} \\ &= \frac{(x_i \cdot \exp(-w_c \cdot x_i)) \left( 1 + \sum_{i=1}^{c-1} \exp(-w_c \cdot x_i) \right) + x_i \cdot \exp(-w_c \cdot x_i) \cdot (\exp(-w_c \cdot x_i))}{\left( 1 + \sum_{c=1}^{c-1} \exp(-w_c \cdot x_i) \right)^2} \\ &= \frac{x_i \exp(-w_c \cdot x_i) \left( 1 + \sum_{c=1}^{c-1} \exp(-w_c \cdot x_i) - \exp(-w_c \cdot x_i) \right)}{\left( 1 + \sum_{c=1}^{c-1} \exp(-w_c \cdot x_i) \right)^2} \\ &= \frac{x_i \cdot \exp(-w_c \cdot x_i)}{1 + \sum_{c=1}^{c-1} \exp(-w_c \cdot x_i)} \cdot \frac{1 + \sum_{c=1}^{c-1} \exp(-w_c \cdot x_i) - \exp(-w_c \cdot x_i)}{1 + \sum_{c=1}^{c-1} \exp(-w_c \cdot x_i)} \\ &= \frac{x_i \cdot \exp(-w_c \cdot x_i)}{1 + \sum_{c=1}^{c-1} \exp(-w_c \cdot x_i)} \cdot \left( \frac{1 + \sum_{c=1}^{c-1} \exp(-w_c \cdot x_i)}{1 + \sum_{c=1}^{c-1} \exp(-w_c \cdot x_i)} - \frac{\exp(-w_c \cdot x_i)}{1 + \sum_{c=1}^{c-1} \exp(-w_c \cdot x_i)} \right) \\ &= x_i \cdot \hat{y}_c (1 - \hat{y}_c)\end{aligned}\tag{11}$$

Subbing (11) back into (10):

$$\begin{aligned}& \frac{\partial (P(y_i | x_i; w_c))}{\partial w_c} \\ &= y_i \cdot \frac{1}{P(y_i | x_i; w_c)} \cdot \frac{\partial (P(y_i | x_i; w_c))}{\partial w_c} \\ &= y_i \cdot \frac{1}{\hat{y}_c} \cdot (x_i \cdot \hat{y}_c (1 - \hat{y}_c)) \\ &= x_i (y_i - y_i \cdot \hat{y}_c)\end{aligned}\tag{12}$$

Putting them back together, we will sub (12) into the gradient descent rule (8)

$$\begin{aligned}
\mathbf{w}_{t+1} &= \mathbf{w}_t - \rho \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} \\
&= \mathbf{w}_t - \rho \left( - \sum_{i=1}^N \sum_{c=1}^C \frac{\partial (y_i \cdot \ln(p(y_i | x_i; w_c)))}{\partial w_c} - \lambda \mathbf{w} \right) \\
&= \mathbf{w}_t + \rho \left( \sum_{i=1}^N (y_i - \hat{y}_i) x_i - \lambda \mathbf{w} \right)
\end{aligned} \tag{13}$$

## 2 Question 2 (5 marks)

### 2.2. Answer:

C=0.01	C=0.05	C=0.1	C=0.5	C=1
0.84958	0.85038	0.85038	0.85050	0.85032

Table 1: Classification accuracy on running linear kernel SVM on 3-fold cross-validation using training set with different values of the parameter C in {0.01, 0.05, 0.1, 0.5, 1}

### 2.3. Answer:

	g=0.01	g=0.05	g=0.1	g=0.5	g=1
C=0.01	0.76377	0.76433	0.77096	0.76377	0.76377
C=0.05	0.78871	0.83293	0.83029	0.76961	0.76377
C=0.1	0.83226	0.83754	0.83717	0.79092	0.76377
C=0.5	0.84411	0.84546	0.84559	0.82507	0.77772
C=1	0.84682	0.84589	0.84651	0.82956	0.78668

Table 2: Classification accuracy on running rbf kernel SVM on 3-fold cross-validation using training set with parameter gamma in {0.01, 0.05, 0.1, 0.5, 1} and different values of the parameter C in {0.01, 0.05, 0.1, 0.5, 1}

### 2.4. Answer:

	kernal=Linear, C=1
Accuracy of SVMs	0.84733

Table 3: Classification accuracy on running linear kernel SVM on 3-fold cross-validation using test set with different values with C=1

## 3 Question 3 (5 marks)

## 4 Question 4 (5 marks)

	<b>kernal=Linear, C=1</b>
<b>Accuracy of SVMs</b>	0.84733

Table 4: Classification accuracy on running linear kernel SVM on 3-fold cross-validation using test set with different values with C=1