

# AI6102 Assignment 1

Sourabh Vyas

16/09/2025

## Question 1 (10 marks)

Consider a multi-class classification problem of  $C$  classes. Based on the parametric forms of the conditional probabilities of each class introduced on the 39th Page (“Extension to Multiple Classes”) of the lecture notes of L4, derive the learning procedure of regularized logistic regression for multi-class classification problems.

Hint: define a loss function by borrowing an idea from binary classification, and derive the gradient descent rules to update  $\{\mathbf{w}^{(c)}\}$  for  $c = 1, \dots, C - 1$ .

### Answer

**1. Conditional Probabilities** We treat class 0 as the baseline and model the others via a “softmax” over negative scores. Then for  $c = 0, \dots, C - 1$ :

$$P(y = 0 \mid \mathbf{x}) = \frac{1}{1 + \sum_{k=1}^{C-1} \exp(-\mathbf{w}^{(k)\top} \mathbf{x})}, \quad P(y = c \mid \mathbf{x}) = \frac{\exp(-\mathbf{w}^{(c)\top} \mathbf{x})}{1 + \sum_{k=1}^{C-1} \exp(-\mathbf{w}^{(k)\top} \mathbf{x})}.$$

This ensures all  $C$  probabilities sum to 1.

**2. Multiclass Negative Log Likelihood** We assume each example  $(x_i, y_i)$  is drawn independently, so the likelihood of all labels given all inputs is

$$L(\{\mathbf{w}^{(c)}\}) = \prod_{i=1}^N P(y_i \mid \mathbf{x}_i)$$

Taking the logarithm turns the product into a sum:

$$\ell(\{\mathbf{w}^{(c)}\}) = \sum_{i=1}^N \log P(y_i \mid \mathbf{x}_i)$$

Two cases for each  $i$ :

If  $y = c \mid c = 1, \dots, C - 1$ :

$$P(y_i = c \mid \mathbf{x}_i) = \frac{\exp(-\mathbf{w}^{(c)\top} \mathbf{x}_i)}{D_i}, \quad D_i = 1 + \sum_{k=1}^{C-1} \exp(-\mathbf{w}^{(k)\top} \mathbf{x}_i)$$

Taking its log

$$\log P(y_i = c \mid \mathbf{x}_i) = -\mathbf{w}^{(c)\top} \mathbf{x}_i - \log D_i$$

If  $y_i = 0$ :

$$P(y_i = 0 \mid \mathbf{x}_i) = \frac{1}{D_i}, \quad \log P(y_i = 0 \mid \mathbf{x}_i) = -\log D_i$$

Introduce  $\mathbb{I}(y_i = c)$  which is 1 when  $y_i = c$ . Then for each  $i$

$$-\log P(y_i | \mathbf{x}_i) = \log D_i + \sum_{c=1}^{C-1} \mathbb{I}(y_i = c) \mathbf{w}^{(c)\top} \mathbf{x}_i$$

Summing the negative log-probabilities over  $i = 1, \dots, N$  yields the total unregularized loss:

$$\mathcal{L} = \sum_{i=1}^N \left[ \log \left( 1 + \sum_{k=1}^{C-1} e^{-\mathbf{w}^{(k)\top} \mathbf{x}_i} \right) + \sum_{c=1}^{C-1} \mathbb{I}(y_i = c) \mathbf{w}^{(c)\top} \mathbf{x}_i \right]$$

**3. Regularized Negative Log-Likelihood** The unregularized negative log-likelihood is

$$\mathcal{L} = \sum_{i=1}^N \left[ \log \left( 1 + \sum_{k=1}^{C-1} e^{-\mathbf{w}^{(k)\top} \mathbf{x}_i} \right) + \sum_{c=1}^{C-1} \mathbb{I}(y_i = c) \mathbf{w}^{(c)\top} \mathbf{x}_i \right]$$

Add an  $\ell_2$  penalty:

$$\mathcal{L}_{\text{reg}} = \mathcal{L} + \frac{\lambda}{2} \sum_{c=1}^{C-1} \|\mathbf{w}^{(c)}\|^2.$$

**4. Gradient Computation** Regularized negative log-likelihood for class  $c$ :

$$\mathcal{L}_{\text{reg}} = \sum_{i=1}^N \left[ \log \left( 1 + \sum_{k=1}^{C-1} e^{-\mathbf{w}^{(k)\top} \mathbf{x}_i} \right) + \sum_{k=1}^{C-1} \mathbb{I}(y_i = k) \mathbf{w}^{(k)\top} \mathbf{x}_i \right] + \frac{\lambda}{2} \sum_{k=1}^{C-1} \|\mathbf{w}^{(k)}\|^2$$

Consider the loss for a single example  $i$ :

$$\ell_i = \log D_i + \sum_{k=1}^{C-1} \mathbb{I}(y_i = k) \mathbf{w}^{(k)\top} \mathbf{x}_i, \quad D_i = 1 + \sum_{k=1}^{C-1} e^{-\mathbf{w}^{(k)\top} \mathbf{x}_i}$$

Taking the gradient w.r.t  $\mathbf{w}^{(c)}$  splits into two: Derivative of  $\log D_i$ :

$$\frac{\partial}{\partial \mathbf{w}^{(c)}} \log D_i = \frac{1}{D_i} \frac{\partial}{\partial \mathbf{w}^{(c)}} \left( 1 + \sum_k e^{-\mathbf{w}^{(k)\top} \mathbf{x}_i} \right) = -\frac{e^{-\mathbf{w}^{(c)\top} \mathbf{x}_i}}{D_i} \mathbf{x}_i = -P(y = c | \mathbf{x}_i) \mathbf{x}_i$$

Derivative of indicator term:

$$\frac{\partial}{\partial \mathbf{w}^{(c)}} [\mathbb{I}(y_i = c) \mathbf{w}^{(c)\top} \mathbf{x}_i] = \mathbb{I}(y_i = c) \mathbf{x}_i$$

Combining these two for an  $i$ :

$$\frac{\partial \ell_i}{\partial \mathbf{w}^{(c)}} = (\mathbb{I}(y_i = c) - P(y = c | \mathbf{x}_i)) \mathbf{x}_i$$

Summing the per-example gradients over  $i = 1, \dots, N$  yields

$$\nabla_{\mathbf{w}^{(c)}} \mathcal{L} = \sum_{i=1}^N (\mathbb{I}(y_i = c) - P(y = c | \mathbf{x}_i)) \mathbf{x}_i$$

Gradient for regularizer

$$\frac{\partial (\frac{\lambda}{2} \|\mathbf{w}^{(c)}\|^2)}{\partial \mathbf{w}^{(c)}} = \lambda \mathbf{w}^{(c)}$$

Adding the  $\ell_2$  Regularizer

$$\nabla_{\mathbf{w}^{(c)}} \mathcal{L}_{\text{reg}} = \sum_{i=1}^N (\mathbb{I}(y_i = c) - P(y = c | \mathbf{x}_i)) \mathbf{x}_i + \lambda \mathbf{w}^{(c)}$$

**5. Gradient Descent Update** With learning rate  $\eta > 0$ , update each class-weight:

$$\mathbf{w}^{(c)} \leftarrow \mathbf{w}^{(c)} - \eta \left[ \nabla_{\mathbf{w}^{(c)}} \mathcal{L}_{\text{reg}} \right]$$

**5. Prediction Rule** For a new input  $\mathbf{x}^*$ ,

$$y^* = \arg \max_{c \in \{0, \dots, C-1\}} P(y = c \mid \mathbf{x}^*).$$