

2. Softmax regression gradient calculation

Consider a simple Softmax regression model,

$$\hat{y} = \sigma(Wx + b)$$

→ Multiclass classification problem

$$x \in \mathbb{R}^d, W \in \mathbb{R}^{k \times d}, b \in \mathbb{R}^k$$

where d is the input dimension, k is the number of classes, σ is the softmax function:

$$\sigma(a)_i = \frac{\exp(a_i)}{\sum_j \exp(a_j)}$$

(a) Given the cross-entropy loss → Between y true label and \hat{y} softmax

$$l(y, \hat{y}) = - \sum_i y_i \log \hat{y}_i$$

y is the one-hot vector representing true labels $([0, 0, \dots, 1, 0, 0, \dots]^T$ with the 1 corresponding to the true label), derive $\frac{\partial l}{\partial W_{i,j}}$. (You can use your results from Assignment 1)

(b) What happens to the loss function and the gradients when $y_{c_1} = 1, \hat{y}_{c_2} = 1, c_1 \neq c_2$? Why there is no need to worry about this situation?

Remember Sigmoid is used for binary classification methods and Softmax is used on multiclass problems.

Softmax is an extension of the Sigmoid function.

$$\text{Sigmoid Function: } \frac{1}{1 + e^{-x}}$$

$$\text{Softmax Function: } \frac{e^{x_i}}{\sum_{j=1}^K e^{x_j}}$$

2(a) Find $\frac{\partial l}{\partial W_{i,j}}$, derivative of scalar by matrix

$$Wx + b = \sum_{j=1}^d (w_{i,j} x_j) + b$$

2(b) When y_{c_1} and $\hat{y}_{c_2} = 1$ and $c_1 \neq c_2$ it means that Softmax model is classifying c_1 true class like c_2 class

$$l(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_i y_i \log \hat{y}_i$$

Class

There is no need to worry because $\log(\hat{y}_{c_2}) = 0$
 $y_i \cdot 0 = 0$ and for that