

Homework 5, Fall 2024

Mohammad Alshurbaji

11/27/2024

Problem 1: Vanilla Logistic Regression for multi-classification:

- k classes. - The input $x \in \mathbb{R}^d$

- Prob. to each class is:

$$P(Y=k | X=x) = \frac{\exp(w_k^T x)}{1 + \sum_{l=1}^{k-1} \exp(w_l^T x)}, \text{ for } k=1, 2, \dots, k$$

$$P(Y=k | X=x) = \frac{1}{1 + \sum_{l=1}^{k-1} \exp(w_l^T x)}$$

If we define $w_k = 0$, Then:

$$P(Y=k | X=x) = \frac{\exp(w_k^T x)}{1 + \sum_{l=1}^{k-1} \exp(w_l^T x)}, k=1, \dots, k$$

Question 1: What and how many para. are there to be opt.?

1. Weight vectors (w_k)

→ For each class k , there is $w_k \in \mathbb{R}^d$.

→ assuming $w_k = 0$.

→ means: The # of weight para.:

$$d \times (k-1)$$

2. Bias parameters: b_k

→ one bias for each class: $(k-1)$

∴ Total Parameters: $d \times (k-1) + (k-1)$

$$= (d+1)(k-1) \#$$

Question 2 $L(w_1, \dots, w_{k-1}) = \sum_{i=1}^n \ln P(Y=y_i | X=x_i)$

$$= \sum_{i=1}^n \ln \left(\frac{\exp(w_k^T x)}{1 + \sum_{l=1}^{k-1} \exp(w_l^T x)} \right)$$

$$= \sum_{i=1}^n \left(\ln(\exp(w_k^T x)) - \ln \left(1 + \sum_{l=1}^{k-1} \exp(w_l^T x) \right) \right)$$

$$= \underbrace{\sum_{i=1}^n w_k^T x}_I - \underbrace{\sum_{i=1}^n \ln \left(1 + \sum_{l=1}^{k-1} \exp(w_l^T x) \right)}_{II}$$

Question 3 The gradient of L w.r.t. w_k :

$$\frac{\partial L}{\partial w_k} = \sum_{i=1}^n$$

I: $\frac{\partial I}{\partial w_k} = x$

II: $\frac{\partial II}{\partial w_k} = \frac{\exp(w_k^T x) x}{1 + \sum_{l=1}^{k-1} \exp(w_l^T x)}$

I + II: $\frac{\partial L}{\partial w_k} = \sum_{i=1}^n \left(x - \frac{\exp(w_k^T x) x}{1 + \sum_{l=1}^{k-1} \exp(w_l^T x)} \right)$

$$= \sum_{i=1}^n \left(x - P(Y=k | X=x_i) x \right)$$

$$= \sum_{i=1}^n \left(\left(\frac{1}{Y=k} P(Y=k | X=x_i) \right) x \right) \quad \# \text{ (In terms of Probability)}$$

Question 4: With adding the regularization term:

$$f(w_1, \dots, w_{k-1}) = \underbrace{L(w_1, \dots, w_{k-1})}_I - \underbrace{\frac{\lambda}{2} \sum_{l=1}^{k-1} \|w_l\|_2^2}_II$$

$$\frac{\partial f}{\partial w_k} =$$

$$I: \frac{\partial I}{\partial w_k} = \sum_{i=1}^n \left((1 - P(Y=k | X=x_i)) X \right)$$

$$II: \frac{\partial II}{\partial w_k} = -\lambda w_k$$

$$\therefore I + II: \frac{\partial f}{\partial w_k} = \sum_{i=1}^n \left((1 - P(Y=k | X=x_i)) X \right) - \lambda w_k$$

#

Question 5: USPS Handwritten Recognition digit dataset.

→ Image size: 16 x 16

→ For each digit (i.e. 0, 1, ..., 9) there are 600 training samples. 500 testing ones.

(a)

→ Continue in the next pages.