

CS57800: Statistical Machine Learning

HOMEWORK 3

Ruoyu Wu
wu1377@purdue.edu

November 15, 2018

1 Open-ended Questions

1.1 Gradient of the loss function

1.1.1 Base Formula

$$z = w^T x \tag{1}$$

$$g(z) = \frac{1}{(1 + e^{-z})} \tag{2}$$

$$Err(w) = - \sum_i y^i \log(g(w, x_i)) + (1 - y^i) \log(1 - g(w, x_i)) \tag{3}$$

1.1.2 Derivative of Logistic Function

$$\frac{d}{dx} g(z) = \frac{d}{dx} \left(\frac{1}{1 + e^{-z}} \right) \tag{4}$$

$$= \frac{-\frac{d}{dx}(1 + e^{-x})}{(1 + e^{-x})^2} \tag{5}$$

$$= \frac{e^{-x}}{(1 + e^{-x})^2} \tag{6}$$

$$= \left(\frac{1}{1 + e^{-x}} \right) \left(\frac{e^{-x}}{1 + e^{-x}} \right) \tag{7}$$

$$= g(z)(1 - g(z)) \tag{8}$$

1.1.3 Derivative(Gradient) of Logistic Loss Function

$$\frac{\partial \text{Err}(w)}{\partial w_j} = - \sum_i [y^i \frac{\partial}{\partial w_j} \log(g(w, x_i)) + (1 - y^i) \frac{\partial}{\partial w_j} \log(1 - g(w, x_i))] \quad (9)$$

$$= - \sum_i [y^i \frac{\frac{\partial}{\partial w_j} (g(w, x_i))}{g(w, x_i)} + (1 - y^i) \frac{\frac{\partial}{\partial w_j} (1 - g(w, x_i))}{1 - g(w, x_i)}] \quad (10)$$

$$= - \sum_i [y^i \frac{\frac{\partial}{\partial w_j} g(w^T x_i)}{g(w, x_i)} + (1 - y^i) \frac{\frac{\partial}{\partial w_j} (1 - g(w^T x_i))}{1 - g(w, x_i)}] \quad (11)$$

$$= - \sum_i [y^i \frac{g(w^T x_i)(1 - g(w^T x_i)) \frac{\partial}{\partial w_j} (w^T x_i)}{g(w, x_i)} + (1 - y^i) \frac{g(w^T x_i)(1 - g(w^T x_i)) \frac{\partial}{\partial w_j} (w^T x_i)}{1 - g(w, x_i)}] \quad (12)$$

$$= - \sum_i [y^i (1 - g(w^T x_i)) x_{ij} + (1 - y^i) g(w^T x_i) x_{ij}] \quad (13)$$

$$= - \sum_i [g(w^T, x_i) - y^i] x_{ij} \quad (14)$$

1.2 Prove that the logistic loss function is convex

1.2.1 Prove that $-\log(g(w, x_i))$ is convex

Theorem: Function convex \iff its epigraph(i.e. the set of points lying above the graph of function) is convex.

$$A = \{(w, t) \mid -\log(g(w, x_i)) \leq t\} \quad (15)$$

$$= \{(w, t) \mid \log(1 + e^{-W^T x_i}) \leq t\} \quad (16)$$

$$= \{(w, t) \mid 1 + e^{-W^T x_i} \leq e^t\} \quad (17)$$

$$= \{(w, t) \mid e^{-t} + e^{-W^T x_i - t} \leq 1\} \quad (18)$$

Since $f(w, t) = e^{-t} + e^{-W^T x_i - t}$ is convex, A is a sublevel set of $f(w, t)$, then A is a convex set. Since A is the epigraph of $-\log(g(w, x_i))$, $-\log(g(w, x_i))$ is convex.

1.2.2 Prove that $-\log(1 - g(w, x_i))$ is convex

Similarly, we can easily prove that $\log(1 - g(w, x_i))$ is convex function of w .

1.2.3 Prove that $\text{Err}(w)$ is convex

The positive linear combination of convex function is still a convex function. And we have proven that for each x_i , $\log(g(w, x_i))$ and $\log(1 - g(w, x_i))$ are convex functions for w . Therefore, $\text{Err}(w)$ is a convex function for w .

1.3 What is regularization and why is it used

Regularization is a process(term) when learning a model. It is normally used to prevent overfitting.

1.4 The gradient with L2 regularization

$$Err(w) = - \sum_i [y^i \log(g(w, x_i)) + (1 - y^i) \log(1 - g(w, x_i))] + \frac{1}{2} \lambda ||w^2|| \quad (19)$$

$$\frac{\partial Err(w)}{\partial w_j} = - \sum_i [g(w^T, x_i) - y^i] x_{ij} + \lambda w_j \quad (20)$$

1.5 Stopping criteria for GD and SGD

For GD, the algorithm will stop when the accuracy of learned model on the validation set has not improved for 5 consecutive steps, or the max_iteration is reached. For SGD, the algorithm will stop when the accuracy of learned model on the validation set has not improved for 5 consecutive steps, or the max_iteration is reached.

1.6 Effect of bias term in GD/SGD

Bias allows us to shift the activation function to the left or right. If we have: $z = w^T x + b$, then we have $g(z) = \frac{1}{(1+e^{w^T x + b})}$. I think bias term is useful, especially in the case of one-layer neural network.

2 Batch Gradient Descent with Logistic Function

2.1

Algorithm 1 BatchGradientDescentLogistic

```

1: procedure GBDTRAINING(DATA, W,  $\alpha$ ,  $\lambda$ )
2:   if It is the first step then
3:     size_data  $\leftarrow$  sizeof(D)
4:     size_feature  $\leftarrow$  sizeof(W) + 1
5:      $w_j \leftarrow 0$ , for all  $j = 1 \dots \text{size\_feature}$ .
6:      $\text{gradient}_j \leftarrow 0$ , for  $j = 1 \dots \text{size\_feature}$ 
7:   end if
8:   for (x, y) in D do
9:      $\text{gradient}_j \leftarrow \text{gradient}_j + [g(w^T, x) - y]x_j$ , for all  $j = 1 \dots \text{size\_feature}$ 
10:  end for
11:   $\text{gradient}_j \leftarrow \frac{\text{gradient}_j}{\text{size\_data}} + \lambda w_j$ , for all  $j = 1 \dots \text{size\_feature}$ 
12:   $W_j \leftarrow W_j - \alpha \text{gradient}_j$ , for all  $j = 1 \dots \text{size\_feature}$ 
13:  return W
14: end procedure
15: procedure TRAIN()
16:   devide DATA into D(for training) and V(for validation)
17:   while True do
18:     if Accuracy on V has not improved for consecutive 5 steps then
19:       break
20:     end if
21:     for Each clf in classifiers do
22:       clf.train(D)
23:     end for
24:   end while
25: end procedure

```

2.2 Graph

Blue lines indicate training curve and green lines indicate testing curve.

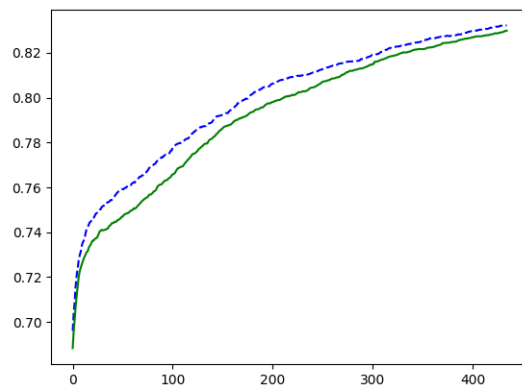


Figure 1: With regularization, type 1

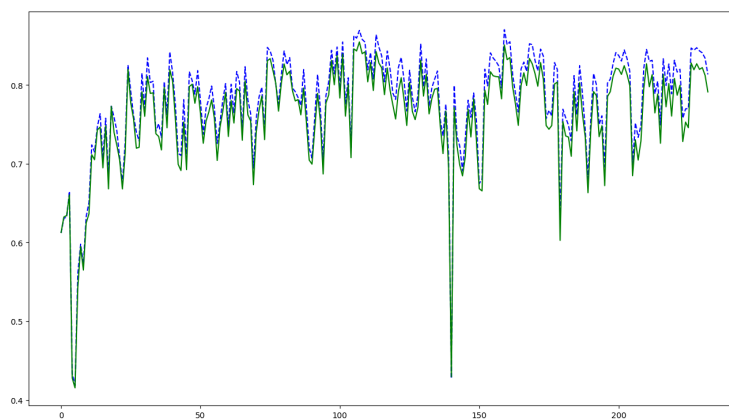


Figure 2: With regularization, type 2

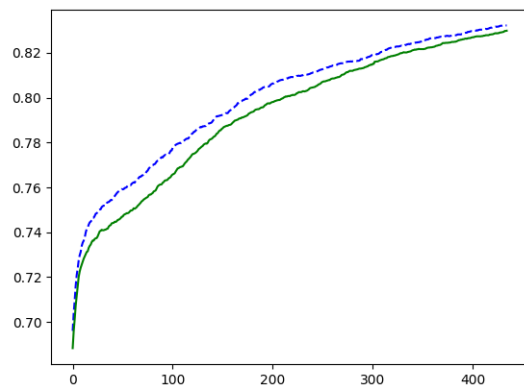


Figure 3: Without regularization, type 1

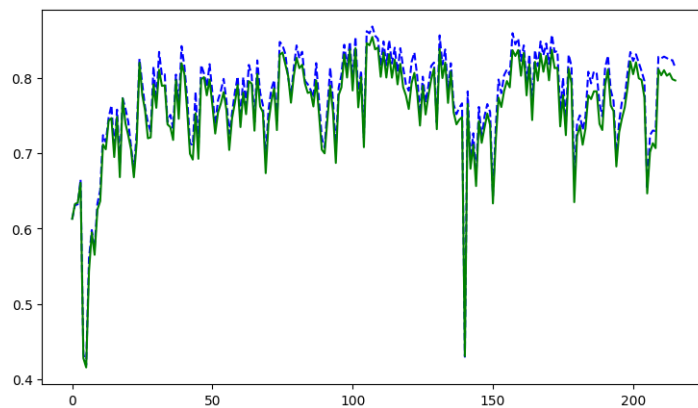


Figure 4: Without regularization, type 2

2.3 Analyze and compare

In terms of convergence speed, regularization term does not have huge impact. Type-2 converges much faster than type-1.

In terms of accuracy, regularization term does not have huge impact. Type-2.

It is because I set the factor of regularization term very small, so it does not make big difference. And type-2 can better represent the image because it filters some noises out and has much smaller feature space, so it achieves better converge speed and accuracy.

3 Stochastic Gradient Descent with Logistic Function

3.1

Algorithm 2 StochasticGradientDescentLogistic

```

1: procedure TRAINING( $D, W, \alpha, \lambda, \text{SIZE\_BATCH}, \text{MAX\_ITERS}$ )
2:   if It is the first step then
3:      $\text{size\_data} \leftarrow \text{sizeof}(D)$ 
4:      $\text{size\_feature} \leftarrow \text{sizeof}(W) + 1$ 
5:      $w_j \leftarrow 0$ , for all  $j = 1 \dots \text{size\_feature}$ .
6:      $\text{gradient}_j \leftarrow 0$ , for  $j = 1 \dots \text{size\_feature}$ 
7:   end if
8:    $S \leftarrow$  sample uniformly  $\text{size\_batch}$  of data from  $D$ 
9:   for  $(x, y)$  in  $S$  do
10:     $\text{gradient}_j \leftarrow \text{gradient}_j + [g(w^T, x) - y]x_j$ , for all  $j = 1 \dots \text{size\_feature}$ 
11:  end for
12:   $\text{gradient}_j \leftarrow \frac{\text{gradient}_j}{\text{size\_batch}} + \lambda w_j$ , for all  $j = 1 \dots \text{size\_feature}$ 
13:   $W_j \leftarrow W_j - \alpha \text{gradient}_j$ , for all  $j = 1 \dots \text{size\_feature}$ 
14:  return  $W$ 
15: end procedure
16: procedure TRAIN()
17:   devide DATA into  $D$ (for training) and  $V$ (for validation)
18:   while True do
19:     if Accuracy on  $V$  has not improved for consecutive 5 steps then
20:       break
21:     end if
22:     for Each clf in classifiers do
23:       clf.train( $D$ )
24:     end for
25:   end while

```

3.2 Graph

Blue lines indicate training curve and green lines indicate testing curve.

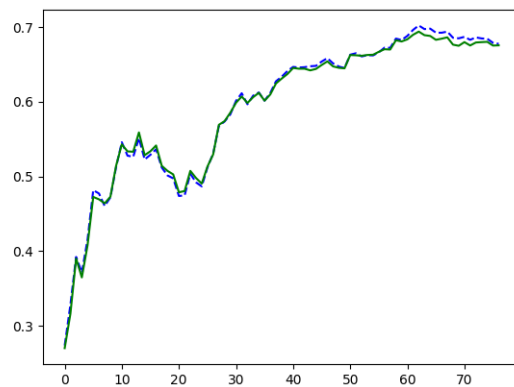


Figure 5: With regularization, type 1

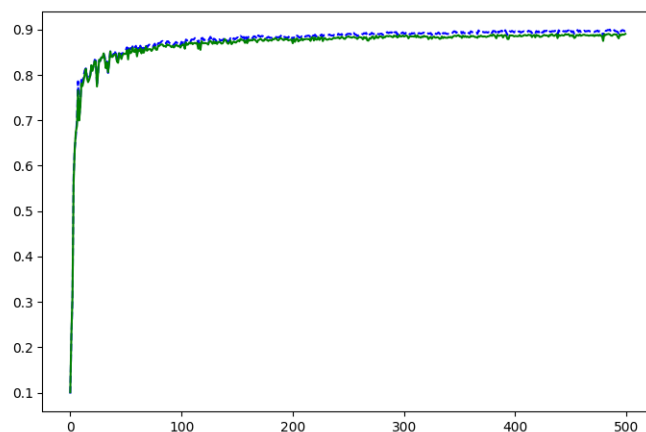


Figure 6: With regularization, type 2

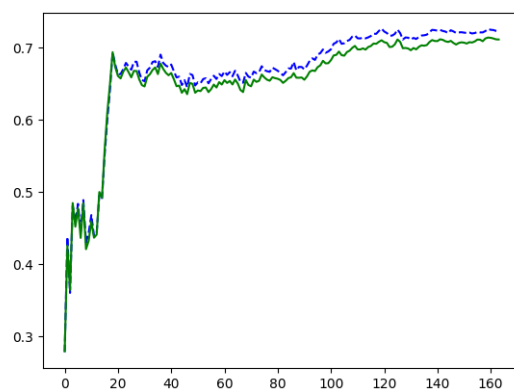


Figure 7: Without regularization, type 1

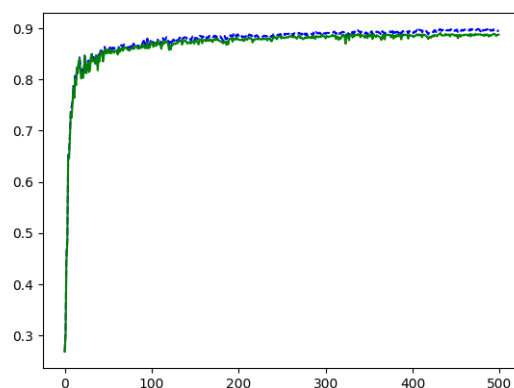


Figure 8: Without regularization, type 2

3.3 Analyze and compare

In terms of convergence speed, It is clear that my convergence condition is not good for the type-2 – they actually converge very fast(around 100 steps), and my convergence condition stop them at *max_iter*. Generally speaking, type-1 and type-2 converges roughly at the same speed. And with regularization term converges faster than without the term.

In terms of the accuracy, type-2 achieves better accuracy than type-1. Regularization term seems has no huge effect on accuracy in this setting.

It is because I set the factor of regularization term very small, so it does not make big difference. And type-2 can better represent the image because it filters some noises out and has much smaller feature space, so it achieves better converge speed and accuracy.