

Machine Learning Homework 2-Saket Vishwasrao

Written Problems

1. The gradient for the given objective function will be

$$\nabla_{\mathbf{w}_c} L = \lambda \mathbf{w}_c + \sum_{i=1}^n \mathbf{x}_i \left(\frac{\exp(\mathbf{w}_c^\top \mathbf{x}_i)}{\sum_{c'} \exp(\mathbf{w}_{c'}^\top \mathbf{x}_i)} - I(y_i = c) \right) \quad (1)$$

(2)

At the optimum the gradient is zero. Hence we have

$$\lambda \mathbf{w}_c = - \sum_{i=1}^n \mathbf{x}_i \left(\frac{\exp(\mathbf{w}_c^\top \mathbf{x}_i)}{\sum_{c'} \exp(\mathbf{w}_{c'}^\top \mathbf{x}_i)} - I(y_i = c) \right) \quad (3)$$

(4)

Summing over all classes in each dimension we get,

$$\lambda \sum_{c=1}^C w_c[j] = - \sum_{c=1}^C \sum_{i=1}^n x_i[j] \left(\frac{\exp(\mathbf{w}_c^\top \mathbf{x}_i)}{\sum_{c'} \exp(\mathbf{w}_{c'}^\top \mathbf{x}_i)} - I(y_i = c) \right) \quad (5)$$

(6)

where $j=1,2,3,\dots,d$

Rearranging the summation

$$\lambda \sum_{c=1}^C w_c[j] = - \sum_{i=1}^n x_i[j] \left(\sum_{c=1}^C \frac{\exp(\mathbf{w}_c^\top \mathbf{x}_i)}{\sum_{c'} \exp(\mathbf{w}_{c'}^\top \mathbf{x}_i)} - \sum_{c=1}^C I(y_i = c) \right) \quad (7)$$

(8)

Now summing over all classes, the term $\sum_{c=1}^C \frac{\exp(\mathbf{w}_c^\top \mathbf{x}_i)}{\sum_{c'} \exp(\mathbf{w}_{c'}^\top \mathbf{x}_i)} = 1$, because it is the sum of all $p(y = c | \mathbf{x}; W)$ over all the classes.

$\sum_{c=1}^C I(y_i = c)$ is a $C \times 1$ matrix that has only one out of its C elements equal to one corresponding to the predicted class for that particular sample and rest all are zero. Hence $\sum_{c=1}^C I(y_i = c) = 1$

Thus we have

$$\lambda \sum_{c=1}^C w_c[j] = 0 \quad (9)$$

which implies

$$\sum_{c=1}^C w_c[j] = 0 \quad (10)$$

2. (a) From equations 5 and 6 in the problem statement, we have

$$\mathbf{w} = \mathbf{w} - \eta(y'_i - y_i)\mathbf{x} \quad (11)$$

where y'_i is the predicted label and y_i is the given label. η is a constant. Hence we have

$$\nabla J = \eta(y'_i - y_i)\mathbf{x} \quad (12)$$

Therefore

$$J = \int \eta(y'_i - y_i) \mathbf{x} d\mathbf{w} \quad (13)$$

$$\implies J = \int \eta(\text{sgn}(\mathbf{w}^T \mathbf{x}) - y_i) \mathbf{x} d\mathbf{w} \quad (14)$$

$$\implies J = \int \eta(\text{sgn}(\mathbf{w}^T \mathbf{x}) - y_i) \mathbf{x} d\mathbf{w} \quad (15)$$

$$\implies J = \eta(|\mathbf{w}^T \mathbf{x}| - y_i \mathbf{w}^T \mathbf{x}) \quad (16)$$

Since we have $y'_i = \frac{|\mathbf{w}^T \mathbf{x}|}{\mathbf{w}^T \mathbf{x}}$

$$J = \eta(y'_i \mathbf{w}^T \mathbf{x} - y_i \mathbf{w}^T \mathbf{x}) \quad (17)$$

$$\implies J = \eta(y'_i - y_i) \mathbf{w}^T \mathbf{x} \quad (18)$$

(b) Using the result from equation 18, we have a batch training function as

$$F = \arg \min_w \sum_{i=1}^N |\eta(y'_i - y_i) \mathbf{w}^T \mathbf{x}| \quad (19)$$

- (c) For the correct classification of all data points, there should be no loss. Hence the value of objective F (eq 19) is zero for correct classification. There will always be a trivial solution to the eq(19) with $\mathbf{w} = 0$ besides the optimum non-trivial \mathbf{w} that minimizes the loss function.
- (d) For a non-linear data set, again we have a trivial $\mathbf{w} = 0$ as the solution. Since we cannot have a non-trivial \mathbf{w} that exactly classifies data, optimizing the loss function will give $\mathbf{w} = 0$ which is not what we want. We want a \mathbf{w} that correctly classifies the most number of samples i.e. increases the accuracy. This is a degeneracy as in practice most datasets are non-linear.
- (e) The perceptron update still works because in perceptron we do not try to find the optimum \mathbf{w} over the entire batch of data but we update \mathbf{w} for each sample. Thus we can stop the optimisation after we have evaluated a specific number of iterations which may not arrive at the optimum \mathbf{w} , but gives sufficient accuracy i.e. we do not approach the trivial solution. Hence perceptron update works well despite the degeneracies.