

### Problem 1.

[2pts] Derive the forward and backward propagation equations for the cross-entropy loss function.

*Hint: The gradient update equations for the sum of squares loss can be found in p396 of Elements of Statistical Learning, which is used for regression tasks. Here we consider the classification task and replace the sum of squares loss by the cross-entropy loss. You should highlight in your solution what parts have changed.*

For cross-entropy deviance, we have

$$R(\theta) = - \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log(f_k(x_i)),$$

and the corresponding classifier is  $G(x) = \operatorname{argmax}_k f_k(x)$

Let  $z_{mi} = \sigma(\alpha_{0m} + \alpha_m^T x_i)$

$z_i = (z_{1i}, z_{2i}, \dots, z_{Mi})$ , then we have

$$\begin{aligned} R(\theta) &= \sum_{i=1}^N R_i \\ &= \sum_{i=1}^N \sum_{k=1}^K (-y_{ik} \log(f_k(x_i))), \end{aligned}$$

With derivatives

$$\left\{ \begin{array}{l} \frac{dR_i}{d\beta_{km}} = -\frac{y_{ik}}{f_k(x_i)} g'_k(\beta_k^T z_i) z_{mi}, \\ \frac{dR_i}{d\alpha_{m1}} = -\sum_{k=1}^K \frac{y_{ik}}{f_k(x_i)} g'_k(\beta_k^T z_i) \beta_{km} \sigma'(\alpha_{0m} + \alpha_m^T x_i) x_{i1} \end{array} \right. \quad (1)$$

Given these derivatives, a gradient descent update at the  $(r+1)$ st iteration has the form

$$\begin{cases} \beta_{km}^{(r+1)} = \beta_{km}^r - \gamma_r \sum_{i=1}^N \frac{d R_i}{d \beta_{km}^r} \\ \alpha_{ml}^{(r+1)} = \alpha_{ml}^{(r)} - \gamma_r \sum_{i=1}^N \frac{d R_i}{d \alpha_{ml}^{(r)}} \end{cases}$$

where  $\gamma_r$  is the learning rate, we write (1)

as

$$\begin{cases} \frac{d R_i}{d \beta_{km}} = \delta_{ki} \otimes_{mi} \\ \frac{d R_i}{d \alpha_{ml}} = s_{mi} x_{il} \end{cases}$$

From their definitions, we have

$$s_{mi} = \zeta'(\alpha_{om} + \alpha_m^\top x_i) \sum_{k=1}^K \beta_{km} \delta_{ki}$$

known as back-propagation equations.

## Problem 2.

[2pts] In a  $K$ -class classification problem, suppose that for given neural network parameter (weights)  $\theta$  and input  $x$ , the output of the network are functions  $f_k(x, \theta)$ ,  $k = 1, \dots$ , which are nonnegative numbers that sum to 1. Show that minimizing the cross-entropy loss is equivalent to maximum likelihood estimation, if  $f_k(x, \theta)$  is interpreted as the likelihood  $P(Y = k|x, \theta)$ .

In a  $K$ -class classification problem, we are given a neural-network that takes an input  $x$  and produces  $K$  non-negative numbers as outputs.

These outputs can be interpreted as probabilities of the input belonging to each of the  $K$  classes. Let  $f(x, \theta_0)$  denote the output of the network when its weights are set to 0, i.e. the network is not trained.

Now suppose we have a training set of inputs  $X_i$  and their corresponding true labels  $y_i$ , where  $y_i$  is an integer in the range 1 to  $K$  indicating the true classes of  $X_i$ .

$$\text{Let } p(y_i | X_i, \omega) = \frac{f(X_i, \omega)_{y_i}}{\sum_{k=1}^K f(X_i, \omega)_k}$$

where  $f(X_i, \omega)_k$  denotes the  $k$ th output of the network for input  $X_i$  when its weights are set to  $\omega$ .

The maximum-likelihood estimate of the weights  $w$  is obtained by maximizing the likelihood of the training set given the weights.

Assuming the training examples are independent and identically distributed, the likelihood is given by the product of the probabilities assigned by the network to the true labels:

$$L(w) = \prod_{i=1}^N p(y_i | x_i, w)$$

where  $N$  is the size of the training set.

Maximizing the likelihood is equivalent to minimizing the negative log-likelihood:

$$J(w) = - \sum_{i=1}^N (\log(p(y_i | x_i, w)))$$

The negative log-likelihood is equivalent to the cross entropy loss, which is defined as follows:

$$J(\omega) = - \sum_{i=1}^n (y_i \times \log(p(y_i | x_i, \omega)))$$

where  $y_i$  is a one-hot vector representing the true label of  $x_i$ , i.e.  $y_i$  has a 1 in the position corresponding to the true class of  $x_i$  and 0s elsewhere.

### Problem 3.

Train a one hidden layer neural network with scalar input and output  $(X_i, Y_i)_{i=1}^n$ . Consider a toy problem setting where the activation function degenerates to a linear (identity) function<sup>[1]</sup>, so that for any input  $x$ , the network output is calculated as

$$z_j := \alpha_j x; \quad (1)$$

$$f(x) := \frac{1}{\sqrt{p}} \sum_{j=1}^p \beta_j z_j \quad (2)$$

where  $z_j$ 's are neurons in the hidden layer and  $(\alpha_j)_{j=1}^p$  and  $(\beta_j)_{j=1}^p$  are the weights. Define the loss function  $E := \sum_{i=1}^n (Y_i - f(X_i))^2$ . Suppose that the gradient descent step size is

small, so that we focus on the continuous-time version of backpropagation:

$$\frac{d\alpha_j}{dt} = -\frac{\partial E}{\partial \alpha_j}; \quad (3)$$

$$\frac{d\beta_j}{dt} = -\frac{\partial E}{\partial \beta_j}. \quad (4)$$

1. [2pt] Express the right sides of (3) and (4) in terms of  $X_i$ ,  $\alpha_j$ ,  $\beta_j$ , and the residue  $R_i := Y_i - f(X_i)$  ( $i = 1, \dots, n$ ).
2. [1pt] Show that when training completes (reaches a stationary point), the neural net function  $\hat{f}(x) = \frac{1}{\sqrt{p}} \sum_{j=1}^p \hat{\alpha}_j \hat{\beta}_j x$  coincides with the least squares estimator. Hint: set the right sides of (3) and (4) to zero.
3. [1pt bonus] Show that for each  $j$ ,  $\alpha_j^2 - \beta_j^2$  remains a constant during training. Hint: take the ratio of (3) to (4) and use part 1).

$$\begin{aligned}
 1. \frac{d\mathcal{E}}{d\alpha_j} &= \frac{d}{d\alpha_j} \sum_{i=1}^n (Y_i - f(x_i))^2 = -2 \sum_{i=1}^n (Y_i - f(x_i)) \frac{df(x_i)}{d\alpha_j} \\
 &= -\frac{2}{Jp} \sum_{i=1}^n (Y_i - f(x_i)) \beta_j x_i = -\frac{2}{Jp} \sum_{i=1}^n R_i \beta_j x_i
 \end{aligned}$$

Similarity,

$$\begin{aligned}
 \frac{d\mathcal{E}}{d\beta_j} &= \frac{d}{d\beta_j} \sum_{i=1}^n (Y_i - f(x_i))^2 = -2 \sum_{i=1}^n (Y_i - f(x_i)) \frac{df(x_i)}{d\beta_j} \\
 &= -\frac{2}{Jp} \sum_{i=1}^n (Y_i - f(x_i)) z_j = -\frac{2}{Jp} \sum_{i=1}^n (Y_i - f(x_i)) \alpha_j x_i \\
 &= -\frac{2}{Jp} \sum_{i=1}^n R_i \alpha_j x_i
 \end{aligned}$$

Summary :

$$\left\{ \begin{array}{l} \frac{d\alpha_j}{d\mathcal{E}} = -\frac{d\mathcal{E}}{d\alpha_j} = \frac{2}{Jp} \sum_{i=1}^n R_i \beta_j x_i \\ \frac{d\beta_j}{d\mathcal{E}} = -\frac{d\mathcal{E}}{d\beta_j} = \frac{2}{Jp} \sum_{i=1}^n R_i \alpha_j x_i \end{array} \right. \times$$

2. To show that when the training completes, the neural net function

$$\hat{f}(x) = \frac{1}{Jp} \sum_{j=1}^J \hat{\alpha}_j \hat{\beta}_j x$$

coincides with the least square estimators, we need to show that, at the stationary point, the weight  $\hat{\alpha}_j$  and  $\hat{B}_j$  have converged to the values that minimize the least square loss function.

At the stationary point, the gradient descent ③ and ④ reduce to

$$\begin{cases} \frac{d\hat{\alpha}_j}{dt} = -\frac{dZ}{d\alpha_j} = \frac{2}{J_p} \sum_{i=1}^n \hat{B}_j X_i R_i \\ \frac{d\hat{B}_j}{dt} = -\frac{dZ}{dB_j} = \frac{2}{J_p} \sum_{i=1}^n \hat{\alpha}_j X_i R_i \end{cases}$$

Setting these two equations to zero, we get

$$\begin{cases} \sum_{i=1}^n X_i R_i \hat{B}_j = 0 \\ \sum_{i=1}^n X_i R_i \hat{\alpha}_j = 0 \end{cases} \Rightarrow \begin{cases} \frac{dZ}{d\alpha_j} = 0 \\ \frac{dZ}{dB_j} = 0 \end{cases} - \Phi$$

These equations are equivalent to the equations for the least square problem. To see this, let's rewrite  $f(x)$

$$\hat{f}(x) = \frac{1}{J_p} \sum_{j=1}^p \hat{\alpha}_j \hat{B}_j x$$

As for least square method,

$$\left\{ \begin{array}{l} \frac{d}{d\alpha_j} \sum_{i=1}^n (Y_i - f(x_i)) = \frac{dE}{d\alpha_j} = 0 \\ \frac{d}{d\beta_j} \sum_{i=1}^n (Y_i - f(x_i)) = \frac{dE}{d\beta_j} = 0 \end{array} \right. \quad - \textcircled{2}$$

$\textcircled{1} = \textcircled{2}$ . Hence, the solution  $\hat{\alpha}_j$ ,  $\hat{\beta}_j$  for the least square estimator are the same with neural net function. Therefore, the neural net function

$$\hat{f}(x) = \frac{1}{J_p} \sum_{j=1}^p \hat{\alpha}_j \hat{\beta}_j x \text{ coincides with the least square estimator.}$$

where  $\hat{\alpha}_j^{(0)}$  and  $\hat{\beta}_j^{(0)}$  are the initial values of  $\alpha_j$  and  $\beta_j$ , simplify the LHS, we get

$$\ln \left( \frac{x_i}{\hat{\alpha}_j^{(0)}} \right) = \frac{1}{J_p} \sum_i \frac{\hat{\beta}_j^{(0)} x_{ij} R_i}{\hat{\beta}_j} - \frac{1}{J_p} \sum_i \frac{\hat{\beta}_j x_{ij} R_i}{\hat{\beta}_j^{(0)}}$$

Expand  $R_i$  and  $\hat{\beta}_j^{(0)}$  in terms of  $\alpha_j$  and  $\hat{\alpha}_j^{(0)}$  using the expressions derived in prob 1, we obtain

$$\ln \left( \frac{x_i}{\hat{\alpha}_j^{(0)}} \right) = \frac{1}{J_p} \sum_i \frac{\alpha_j^{(0)} x_{ij} R_i}{\alpha_j} - \frac{1}{J_p} \sum_i \frac{\alpha_j x_{ij} R_i}{\alpha_j^{(0)}}$$

Simplify further, we get

$$\ln \left( \frac{\hat{x}_j}{x_j^{(0)}} \right) = \frac{2}{J_p} \sum_i \frac{\alpha_j^{(0)} x_{ij} R_i}{\hat{x}_j} - \frac{2}{J_p} \sum_i \frac{\alpha_j x_{ij} R_i}{x_j^{(0)}}$$

Using the fact that  $\hat{x}_j - \beta_j$  is a constant during training, we can write

$$\hat{x}_j - \beta_j = \frac{4}{p} \sum_i \alpha_j^{(0)} x_{ij} x_{ij} R_i - \frac{4}{p} \sum_i \sqrt{\alpha_j^{(0)} x_{ij}} x_{ij} R_i$$

Simplify further, we get

$$(\hat{x}_j - \beta_j) - (\hat{x}_j^{(0)} - \beta_j^{(0)}) = \frac{4}{p} \sum_i (\alpha_j^{(0)} x_{ij} - \beta_j^{(0)} \hat{x}_j) x_{ij} R_i$$

Since  $(\hat{x}_j^{(0)} - \beta_j^{(0)})$  is a constant, we conclude that  $(\hat{x}_j - \beta_j)$  is as well constant during training.

Using the equation  $\hat{x}_j = \sqrt{p} w_j$ , we can rewrite as

$$\frac{\hat{x}_j}{J_p} x_{ij} R_i$$

which is proportional to the expression for  $\frac{d \hat{x}_j}{dt}$  in (4) at the stationary point. Similarly, we can show that the expression for  $\hat{x}_j$  at the stationary point

is proportional to the solution of the normal equations for the least square problem.

Therefore, at the stationary point (training completes), the weights  $\hat{\alpha}_j$  and  $\hat{\beta}_j$  have converged to the values that minimize the least square loss function, and the neural network  $\hat{f}(x)$  coincides with the least square estimator.

3. Take the ratio of (3) to (4), we have

$$\frac{d\hat{\alpha}_j}{d\hat{\beta}_j} = \sum_i \frac{X_{ij} R_i}{\cancel{T_p \hat{\alpha}_j}}$$

Integrating both sides wrt  $\hat{\beta}_j$ , we obtain

$$\int_{\hat{\alpha}_j^{(0)}}^{\hat{\alpha}_j} \frac{d\hat{\alpha}_j'}{\cancel{T_p \hat{\alpha}_j}} = \int_{\hat{\beta}_j^{(0)}}^{\hat{\beta}_j} \frac{\sum_i X_{ij} R_i}{\hat{\beta}'_j} d\hat{\beta}'_j$$

We have the result above.

Let  $r_j = \hat{\alpha}_j - \hat{\beta}_j$ , then

$$\begin{aligned}
 \frac{d}{dt} r_j &= \frac{d}{dt} (\alpha_j^* - \beta_j^*) \\
 &= 2\alpha_j^* \left( \sum_{i=1}^n (Y_i - f(x_i)) B_j x_i \right) - 2\beta_j^* \left( \sum_{i=1}^n (Y_i - f(x_i)) \alpha_j x_i \right) \\
 &= -2 \sum_{i=1}^n (Y_i - f(x_i)) r_j x_i^*
 \end{aligned}$$

$\Rightarrow$  for each  $j$ ,  $\frac{d}{dt} r_j$  is not depends on  $j$   
 $\Rightarrow \alpha_j^* - \beta_j^*$  remain constant for each  $j$ .

4. At initialization,  $\alpha_j^* = \alpha^0$ ,  $\beta_j^* = \beta^0$

$$\text{Let } E^0 = E(\alpha^0, \beta^0).$$

At stationary point, we have

$$\frac{d\alpha_j^*}{dt} = \frac{d\beta_j^*}{dt} = 0 \quad \text{for all } j's.$$

$$\begin{aligned}
 \Rightarrow & \left. \begin{aligned} &\sum_{i=1}^n (Y_i - f(x_i)) \beta_j x_i \\ &\sum_{i=1}^n (Y_i - f(x_i)) \alpha_j x_i \end{aligned} \right\} \quad \text{for all } j's \\
 +) & \frac{\sum_{i=1}^n (Y_i - f(x_i)) \beta_j x_i}{\sum_{i=1}^n (Y_i - f(x_i)) \alpha_j x_i}
 \end{aligned}$$

$$\begin{aligned}
 \frac{d}{dt} (\alpha_j^* + \beta_j^*) &= \frac{2}{\sqrt{p}} \sum_{i=1}^n (\alpha_j^* + \beta_j^*) x_i R_i \\
 &= \frac{2}{\sqrt{p}} (\alpha_j^* + \beta_j^*) S
 \end{aligned}$$

By observation, we have the general solution  
 $\alpha_j^* + \beta_j^* = Ce^{2tS/\sqrt{p}}$ , where  $C$  is constant that

depends on the initial values of  $\alpha_j^*$  and  $\beta_j^*$ .

```

1 import numpy as np
2
3 # Define the linear activation function
4 def sigma(x):
5     return x
6
7 # Define the loss function
8 def loss(X, Y, alpha, beta):
9     Z = sigma(np.outer(alpha, X))
10    f = np.sum(beta * Z, axis=1) / np.sqrt(len(alpha))
11    return np.sum((Y - f) ** 2)
12
13 # Define the gradient of the loss function with respect to alpha
14 def grad_alpha(X, Y, alpha, beta):
15     Z = sigma(np.outer(alpha, X))
16     f = np.sum(beta * Z, axis=1) / np.sqrt(len(alpha))
17     R = Y - f
18     return -2 * np.sum(np.outer(R, X) * beta * np.sqrt(len(alpha)) * Z, axis=0)
19
20 # Define the gradient of the loss function with respect to beta
21 def grad_beta(X, Y, alpha, beta):
22     Z = sigma(np.outer(alpha, X))
23     f = np.sum(beta * Z, axis=1) / np.sqrt(len(alpha))
24     R = Y - f
25     return -2 * np.sum(np.outer(R, Z) / np.sqrt(len(alpha)), axis=0)
26
27 # Initialize the weights
28 p = 1000
29 alpha = np.random.normal(loc=0, scale=1, size=p)
30 beta = np.random.normal(loc=0, scale=1, size=p)
31
32
33 # Generate some toy data
34 n = 1000
35 X = np.random.normal(loc=0, scale=1, size=n)
36 Y = 2 * X + np.random.normal(loc=0, scale=0.5, size=n)
37
38 # Set the learning rate and the number of iterations
39 eta = 1e-3
40 num_iters = 100
41
42 # Train the neural network using gradient descent
43 for i in range(num_iters):
44     alpha -= eta * grad_alpha(X, Y, alpha, beta)
45     beta -= eta * grad_beta(X, Y, alpha, beta)
46
47 # Compute the final estimate of the function
48 def f_hat(x):
49     Z = sigma(alpha * x)
50     return np.sum(beta * Z) / np.sqrt(p)
51
52 # Test the final estimate on some new data
53 X_test = np.random.normal(loc=0, scale=1, size=n)
54 Y_test = 2 * X_test + np.random.normal(loc=0, scale=0.5, size=n)
55 predictions = np.array([f_hat(x) for x in X_test])
56 mse = np.mean((Y_test - predictions) ** 2)
57 print('MSE:', mse)

```

This code generates a random toy dataset of  $n = 1000$  samples, where  $Y$  is a noisy version of  $2X$ . It then initializes the weights of the neural network with random normal variables and trains the network using gradient descent with a learning rate of  $1e-3$  for 10000 iterations. Finally, it computes the final estimate of the function  $f(x)$  and tests it on a new set of data, computing the mean squared error between the true  $Y$  values and the predicted values.

81

No, it does not necessarily contradict the fact that  $f(x)$  changes noticeably during training. While each individual weight does not change much during training, the overall effect of changing all the weights together can still result in a significant change in the output function. In other words, even small changes in each weight can accumulate and result in a noticeable change in the overall function.

To be more specific, the fact that individual weights change very little during training does not necessarily imply that the overall function  $f(x)$  does not change significantly. In this case, the change in  $f(x)$  from  $x$  to  $2x$  is due to the combined effect of all the weights and inputs.

As  $p \rightarrow \infty$ , the number of weights in the network becomes very large, and each weight contributes only a small amount to the overall function  $f(x)$ . However, since there are many weights, the combined effect of all the weights can still be significant, leading to a large change in  $f(x)$ .

Moreover, the fact that  $c \rightarrow 1$  as  $p \rightarrow \infty$  indicates that the weights are becoming increasingly balanced, with the magnitudes of the  $a_j$  and  $b_j$  becoming closer together. This means that the contribution of each weight to the overall function  $f(x)$  becomes more equal, further emphasizing the combined effect of all the weights.

the definition of the function  $f(x)$ . Specifically, we have:

$$\begin{aligned} f(x) &= \frac{1}{\sum_{j=1}^P b_j} \sum_{j=1}^P b_j z_j = \frac{1}{P} \sum_{j=1}^P b_j \alpha_j x = \frac{1}{P} \sum_{j=1}^P \frac{b_j}{\alpha_j} \alpha_j^2 x \\ &= \frac{1}{P} \sum_{j=1}^P \frac{b_j}{\alpha_j} z'_j \end{aligned}$$

where  $z'_j := \alpha_j^2 x$  is the activation of  $j$ th neuron after applying the linear transformation  $\alpha_j$ .

9. Using the same argument as in part 7, we can show that during training,  $\hat{\alpha}_j - \hat{B}_j$  remains constant for all  $j$ , and that  $\hat{\alpha}_j + \hat{B}_j$  converges to some constant  $C > 0$ .

To compute value of  $C$ , we can take the sum of the equations for  $\frac{d\alpha_j}{dt}$  and  $\frac{dB_j}{dt}$

$$\sum_{j=1}^P \left( \frac{d\alpha_j}{dt} + \frac{dB_j}{dt} \right) = -2 \sum_{i=1}^n R_i X_i = -2 \sum_{i=1}^n (Y_i - f(X_i)) X_i$$

$$= -2 \sum_{i=1}^n \left( Y_i - \frac{1}{P} \sum_{j=1}^P \frac{\alpha_j}{\hat{\alpha}_j} \hat{z}_j \right) X_i = -2 \sum_{i=1}^n (Y_i - f'(X_i)) X_i$$

Where  $f'(x) = \frac{1}{P} \sum_{j=1}^P \frac{\alpha_j}{\hat{\alpha}_j} \hat{z}_j$  is the intermediate function we defined above.

Taking the derivative of the loss function with respect to  $\alpha_j$  and  $B_j$ , we get.

$$\begin{aligned} \frac{d\epsilon}{d\alpha_j} &= -2 \sum_{i=1}^n (Y_i - f(X_i)) \frac{B_j}{\hat{\alpha}_j} \alpha_j X_i \\ &= -2 \sum_{i=1}^n (Y_i - f(X_i)) \hat{\alpha}_j X_i \end{aligned}$$



Setting rhs of these equations to zero, we get:

$$\sum_{i=1}^n (Y_i - f(x_i)) \beta_j x_i = 0$$