

4.(B)

4. In class, we introduced our version of the cross-entropy error function

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N -\ln \theta(y_n \mathbf{w}^T \mathbf{x}_n).$$

based on the definition of $y_n \in \{-1, +1\}$. If we transform y_n to $y'_n \in \{0, 1\}$ by $y'_n = \frac{y_n + 1}{2}$, which of the following error function is equivalent to E_{in} above? Choose the correct answer; explain your answer.

[a] $\frac{1}{N} \sum_{n=1}^N (y'_n \ln \theta(\mathbf{w}^T \mathbf{x}_n) + (1 - y'_n) \ln(\theta(-\mathbf{w}^T \mathbf{x}_n)))$

[b] $\frac{1}{N} \sum_{n=1}^N (y'_n \ln \theta(-\mathbf{w}^T \mathbf{x}_n) + (1 - y'_n) \ln(\theta(\mathbf{w}^T \mathbf{x}_n)))$ ✓

[c] $\frac{1}{N} \sum_{n=1}^N (y'_n \ln \theta(\mathbf{w}^T \mathbf{x}_n) - (1 - y'_n) \ln(\theta(-\mathbf{w}^T \mathbf{x}_n)))$

[d] $\frac{1}{N} \sum_{n=1}^N (y'_n \ln \theta(-\mathbf{w}^T \mathbf{x}_n) - (1 - y'_n) \ln(\theta(\mathbf{w}^T \mathbf{x}_n)))$

[e] none of the other choices

$\min -\ln \theta(y_n \mathbf{w}^T \mathbf{x}_n)$
 $\Rightarrow \min -\theta(y_n \mathbf{w}^T \mathbf{x}_n)$
 $\Rightarrow \min -(y_n \mathbf{w}^T \mathbf{x}_n)$
 $\Rightarrow \min -y_n \ln \theta(\mathbf{w}^T \mathbf{x}_n)$
 $\rightarrow y_n = -1 \rightarrow \ln \theta(\mathbf{w}^T \mathbf{x}_n) \quad y'_n = 0$
 $y_n = +1 \rightarrow \ln \theta(\mathbf{w}^T \mathbf{x}_n) \quad y'_n = 1$

when $y_n = -$
 when $y'_n = 0$ 对.
 $1 - y'_n$ we got 对式.
 when $y'_n = 1$ 对
 $y_n = 1$
 前式.

5.(E)

$$P(\mu \leq \mu) \geq 1 - \delta$$

$$P(\mu - \mu) > \epsilon \leq \delta$$

$$Pr(|\mu - \mu| > \epsilon) \leq 2 \exp(-2\epsilon^2 N)$$

$$2 \exp(-2\epsilon^2 N) \leq \delta$$

$$-2\epsilon^2 N \leq \ln \frac{\delta}{2}$$

$$-\epsilon \leq \frac{1}{\sqrt{2N}} \ln \frac{\delta}{2}$$

$$\epsilon \leq \sqrt{\frac{1}{2N} \ln \frac{2}{\delta}}$$

$$\epsilon \leq \sqrt{\frac{\ln 2}{2N\delta}}$$

$$\rightarrow \mu \leq v + \sqrt{\frac{1}{2N} \ln \frac{2}{\delta}}$$

Yes

upper bound

by Hoeffding's inequality

$$\frac{dE_{\text{Sqr}}}{d\hat{y}} = \frac{\sum (y^2 - 2\hat{y}y_n + y_n^2)}{d\hat{y}}$$

$$= 2\hat{y} - \frac{1}{N} \sum y_n = 0$$

$$\Rightarrow \hat{y} = \frac{1}{N} \sum y_n \Rightarrow v$$

Yes

$$\text{Likelihood}(\hat{\mu}) = \prod P(x_i) P(y_i | \hat{\mu})$$

$$P(1 | \hat{\mu}) = \hat{\mu}$$

$$P(0 | \hat{\mu}) = 1 - \hat{\mu}$$

$$\rightarrow \max \sum \ln(P(y_i | \hat{\mu}))$$

if $y=1$ has a times $y=0$ has b times

$$\rightarrow \frac{d \sum \ln(P(y_i | \hat{\mu}))}{d\hat{\mu}} = \frac{d(a \ln \hat{\mu} + b \ln(1 - \hat{\mu}))}{d\hat{\mu}}$$

$$\Rightarrow \frac{a}{\hat{\mu}} - \frac{b}{1 - \hat{\mu}} = 0$$

$$a(1 - \hat{\mu}) - b\hat{\mu} = 0$$

$$\hat{\mu} = \frac{a}{a+b} = v$$

Yes

$$\frac{dE_{\text{ce}}}{d\hat{y}} = \frac{\frac{1}{N} \sum (y_n \ln \hat{y} + (1 - y_n) \ln(1 - \hat{y}))}{d\hat{y}}$$

$$= \left[\frac{\frac{1}{N} \sum y_n}{\hat{y}} + \frac{-1}{1 - \hat{y}} + \frac{\frac{1}{N} \sum y_n}{1 - \hat{y}} \right]$$

$$= \left[\frac{\frac{1}{N} \sum y_n - \frac{\hat{y}}{N} \sum y_n - \hat{y} + \frac{\hat{y}}{N} \sum y_n}{\hat{y}(1 - \hat{y})} \right]$$

$$= \left[\frac{\frac{1}{N} \sum y_n - \hat{y}}{\hat{y}(1 - \hat{y})} \right] = \frac{\hat{y} - v}{\hat{y}(1 - \hat{y})} = 0$$

$$\Rightarrow \hat{y} - v = 0 \Rightarrow \hat{y} = v \neq \text{Yes}$$

6. (A)

Stochastic Gradient Descent

6. In the perceptron learning algorithm, we find one example $(\mathbf{x}_{n(t)}, y_{n(t)})$ that the current weight vector \mathbf{w}_t mis-classifies, and then update \mathbf{w}_t by

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_{n(t)} \mathbf{x}_{n(t)}.$$

The algorithm can be viewed as optimizing some $E_{\text{in}}(\mathbf{w})$ that is composed of one of the following point-wise error functions with stochastic gradient descent (neglecting any non-differentiable points of the error function). What is the error function? Choose the correct answer; explain your answer.

- [a] $\text{err}(\mathbf{w}, \mathbf{x}, y) = \max(0, -y\mathbf{w}^T \mathbf{x})$
- [b] $\text{err}(\mathbf{w}, \mathbf{x}, y) = -\max(0, -y\mathbf{w}^T \mathbf{x})$
- [c] $\text{err}(\mathbf{w}, \mathbf{x}, y) = \max(y\mathbf{w}^T \mathbf{x}, -y\mathbf{w}^T \mathbf{x})$
- [d] $\text{err}(\mathbf{w}, \mathbf{x}, y) = -\max(y\mathbf{w}^T \mathbf{x}, -y\mathbf{w}^T \mathbf{x})$
- [e] none of the other choices

if error $(y, \mathbf{w}^T \mathbf{x})$ is negative
- $(y\mathbf{w}^T \mathbf{x})$

if correct
err = 0
that is $y\mathbf{w}^T \mathbf{x} \geq 0$
→ if
(as $\max(0, -y\mathbf{w}^T \mathbf{x})$)

7.(A)

$$\sum_{k=0} \ln[y=k] \frac{e^{w^T x}}{\sum_{j=0} e^{w^T x}}$$

取完
gradient
要加负号

→ 正确分类的 k 时, $\ln 1 = 0$, err 没有贡献

$$\Rightarrow \text{展开} \quad h_k = \sum_{k=0} [y=k] \ln e^{w^T x} - [y=k] \ln \sum_{j=0} e^{w^T x}$$

$$= \sum_{k=0} [y=k] w^T x - \sum_{k=0} [y=k] \ln \sum_{j=0} e^{w^T x}$$

$$\sum [y=k] w^T x = (y=0)w^T x + (y=1)w^T x + \dots (y=k)w^T x + \dots \Rightarrow (y=k)w^T x$$

[只有 $y=k$ 有值, that is 样本正确
类别, 其它=0]

$$\text{同理} \quad \sum [y=k] \ln \sum_{j=0} e^{w^T x} = \ln \sum_{j=0} e^{w^T x}$$

$$h_k = (y=k)w^T x - \ln \sum_{j=0} e^{w^T x}$$

$$\frac{\partial h_k}{\partial w_t} = \left(x - \frac{1}{\sum e^{w^T x}} e^{w^T x} \cdot x \right) \quad (t=C)$$

$$= \left(0 - \frac{1}{\sum e^{w^T x}} e^{w^T x} \cdot x \right) \quad (t \neq C) = -x(1 - h_y)$$

$$= -x \left(-\frac{e^{w^T x}}{\sum e^{w^T x}} \right)$$

$$= x(h_y)$$

$$\nabla \frac{\partial h_k}{\partial w_t} = (h_y - 1(y=k))x$$

$$W \Rightarrow W + (1(y=k) - h_y)x$$

→ 每步都要更新

$$\rightarrow (1 - h_y)x \rightarrow A$$

8. (E)

Nonlinear Transformation

8. Given the following training data set:

$$\mathbf{x}_1 = (0, 1), y_1 = -1 \quad \mathbf{x}_2 = (0, -1), y_2 = -1 \quad \mathbf{x}_3 = (-1, 0), y_3 = +1 \quad \mathbf{x}_4 = (1, 0), y_4 = +1$$

Use the quadratic transform $\Phi_2(\mathbf{x}) = (1, x_1, x_2, x_1^2, x_1x_2, x_2^2)$ and take $\text{sign}(0) = 1$. Which of the following weight vector $\tilde{\mathbf{w}}$ represents a linear classifier in the \mathcal{Z} -space that can separate all the transformed examples perfectly? Choose the correct answer; explain your answer.

- ☒ [a] $(0, -1, 0, 0, 0, 0)$
- ☒ [b] $(0, 0, -1, 0, 0, 0)$
- ☒ [c] $(0, 0, 0, -1, 0, 0)$
- ☒ [d] $(0, 0, 0, 0, -1, 0)$
- ☐ [e] $(0, 0, 0, 0, 0, -1)$

$$\mathbf{x}_1 = [1, 0, 1, 0, 0, 0]$$

$$\mathbf{x}_2 = [1, 0, -1, 0, 0, 0]$$

$$\mathbf{x}_3 = [1, -1, 0, 1, 0, 0]$$

$$\mathbf{x}_4 = [1, 1, 0, 1, 0, 0]$$

4 of 7

W a · x₁ = 1 wrong
W b · x₂ = 1 wrong
W c · x₃ = 1 wrong
W d correct
W e =

$$\tilde{\mathbf{w}} = \mathbf{x}^T \mathbf{y}$$

9. (B)

$$\mathbf{w} = \mathbf{X}^T \mathbf{P} \cdot \mathbf{y}$$

$$\hat{\mathbf{w}} = (\mathbf{P} \mathbf{X})^T \mathbf{y}$$

$$(\mathbf{P} \mathbf{X})^T \hat{\mathbf{w}} = \mathbf{y}$$

$$\mathbf{X}^T \mathbf{w} = \mathbf{y} \rightarrow \mathbf{X}^T \mathbf{w} = (\mathbf{P} \mathbf{X})^T \hat{\mathbf{w}}$$

$$\mathbf{X}^T \mathbf{w} = \mathbf{X}^T \mathbf{P}^T \hat{\mathbf{w}}$$

$$\mathbf{w} = \mathbf{P}^T \hat{\mathbf{w}}$$

10. (C)

10. After "visualizing" the data and noticing that all $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ are distinct, Dr. Trans magically decides the following transform

$$\Phi(\mathbf{x}) = ([\mathbf{x} = \mathbf{x}_1], [\mathbf{x} = \mathbf{x}_2], \dots, [\mathbf{x} = \mathbf{x}_N])$$

$$\Phi(\mathbf{x}) = ([\mathbf{x} = \mathbf{x}_1], [\mathbf{x} = \mathbf{x}_2], \dots, [\mathbf{x} = \mathbf{x}_N])$$

if $\mathbf{x} = \mathbf{x}_i$ then $\rightarrow 1$
otherwise $\rightarrow 0$

That is, $\Phi(\mathbf{x})$ is a N -dimensional vector whose n -th component is 1 if and only if $\mathbf{x} = \mathbf{x}_n$. If we run linear regression after applying this transform, what is the optimal $\tilde{\mathbf{w}}$? Choose the correct answer; explain your answer.

[a] 1, the vector of all 1s.

[b] 0, the vector of all 0s.

☒ [c] \mathbf{y}

[d] $-\mathbf{y}$

[e] none of the other choices

(Note: Be sure to also check what $E_{in}(\tilde{\mathbf{w}})$ is!)

$$\mathbf{X} = [14, 15, 10] \rightarrow \text{that is}$$

$$\Phi(15) = [0, 1, 0]$$

the data after transform
 $\rightarrow [0, 0, 0, 0, \dots, 1, 0, \dots, 0]$

$$\tilde{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\tilde{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \rightarrow \tilde{\mathbf{w}} = \mathbf{y}$$

Assume that we could do linear regression with one-versus-all decomposition for multi-class classification.

11.(C)

11. Assume that we couple linear regression with one-versus-all decomposition for multi-class classification, and get K weight vectors $\mathbf{w}_{[k]}^*$. Assume that the squared error $E_{in}^{sqr}(\mathbf{w}_{[k]}^*)$ for the k -th binary classification problem is e_k . What is the tightest upper bound of $E_{in}^{0/1}(g)$, where g is the multi-class classifier formed by the one-versus-all decomposition? Choose the correct answer; explain your answer.

linear one-versus-all → 不是线性分类

[a] $2 \sum_{k=1}^K e_k$
 [b] $\sum_{k=1}^K e_k$
 [c] $\frac{1}{2} \sum_{k=1}^K e_k$ (circled)
 [d] $\frac{1}{K} \sum_{k=1}^K e_k$
 [e] $\frac{1}{2K} \sum_{k=1}^K e_k$

OVA 最好的 prediction 来自于

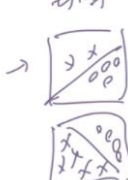
最少有幾個 classifier 犯錯

$E_{in}^{sqr}(\mathbf{w}_{[k]}^*)$
 $\Rightarrow \frac{1}{L} \sum_{n=1}^L (\mathbf{w}_{[k]}^{*T} \mathbf{x}_n - y_n)^2 = e_k$
 $\text{sign}[\mathbf{w}_{[k]}^T \mathbf{x}_n \times y_n]$

$N \cdot E_{in}^{0/1}(g) \leq N \sum_i \frac{E_{in}^{0/1}(\mathbf{w}_i^*)}{2}$
 $\leq N \frac{K}{2} \frac{E_{in}^{sqr}(\mathbf{w}_{[k]}^*)}{2}$

TA hour !!

每個 $\mathbf{w}_{[k]}^T \mathbf{x}$
 → argmax_k
 → 選分類最好的
 機率提供最大的



12. 13.14.15.16.

```
import numpy as np
import random

def getdata(a):
    text = []
    if a == 1:
        path = 'hw3_train.dat.txt'
    elif a == 2:
        path = 'hw3_test.dat.txt'
    with open(path) as f:
        for line in f:
            text.append([float(i) for i in line.split()])
    mm = np.asarray(text)
    X = mm[:, :-1]
    Y = mm[:, -1]

    return X, Y

    return X, Y
def functionQ(X, q):
    function = []
    func_temp = [1]
    for i in X:
        for j in range(1, q+1):
            for x in i[0:]:
                func_temp.append(x**j)
            function.append(func_temp)
            func_temp = [1]
    function_array = np.array(function)
    return function_array

def functionFullQ(X, q):
    function = []
    func_temp = [1]
    row = X.shape[1]
    a = 1
    b = 0
    for i in X:
        for j in range(1, q+1):
            for x in i[0:]:
                func_temp.append(x**j)
            for _ in range(a, row):
                func_temp.append(x*X[b][_])
                a += 1
            function.append(func_temp)
            func_temp = [1]
            a = 1
            b += 1
    function_array = np.array(function)
    return function_array
```



```

def functionLower(X, q):
    function = []
    func_temp = [1]
    for i in X:
        for x in i[0:q]:
            func_temp.append(x)
            function.append(func_temp)
            func_temp = [1]
    function_array = np.array(function)
    return function_array

def functionRandom(X, n):
    random.seed(n)
    function = []
    func_temp = [1]
    list1 = [0, 1, 2, 3, 4, 5, 6, 7, 8, 9]
    list2 = [5, 6, 7, 8, 9]
    random.shuffle(list1)
    list3 = np.delete(list1, list2)
    for i in range(X.shape[0]):
        newX = np.delete(X[i], list3)
        for _ in range(newX.shape[0]):
            func_temp.append(newX[_])
            function.append(func_temp)
            func_temp = [1]

    function_array = np.array(function)
    return function_array

def linear_regression(x, y):
    x_t = np.transpose(x)
    t = x_t.dot(x)
    t_inv = np.linalg.inv(t)
    w_lin = (t_inv.dot(x_t)).dot(y)
    return w_lin

def E_in_minus_E_out_bin(w, x, y, x_test, y_test, n):
    # E_in_bin
    E_in_bin = 0
    N = y.size
    for xn, yn in zip(x, y):
        E_in_bin += (np.sign(xn.dot(w)) != yn)/N
    if n:
        print("E_in_bin:", E_in_bin)
    # E_out_bin
    E_out_bin = 0
    N_test = y_test.size
    for xn, yn in zip(x_test, y_test):
        E_out_bin += (np.sign(xn.dot(w)) != yn)/N_test
    if n:
        print("E_out_bin:", E_out_bin)
        print("/E_in_bin - E_out_bin/", abs(E_in_bin - E_out_bin))
    if n == 0:
        return abs(E_in_bin - E_out_bin)

```

```

x_train,y_train= getdata(1)
x_test,y_test = getdata(2)

print('q12')
x_train_trans_2 = functionQ(x_train , 2)
x_test_trans_2 =functionQ(x_test, 2)
wlin_train_2 = linear_regression(x_train_trans_2, y_train)
E_in_minus_E_out_bin(wlin_train_2, x_train_trans_2, y_train, x_test_trans_2, y_test,1)

print('q13')
x_train_trans_8 = functionQ(x_train , 8)
x_test_trans_8 =functionQ(x_test, 8)
wlin_train_8 = linear_regression(x_train_trans_8, y_train)
E_in_minus_E_out_bin(wlin_train_8, x_train_trans_8, y_train, x_test_trans_8, y_test,1)

print('q14')
x_train_fulltrans_2 = functionFullQ(x_train, 2)
x_test_fulltrans_2 = functionFullQ(x_test, 2)
wlin_fulltrain_2 = linear_regression(x_train_fulltrans_2, y_train)
E_in_minus_E_out_bin(wlin_fulltrain_2, x_train_fulltrans_2, y_train, x_test_fulltrans_2, y_test,1)

print('q15')
x_train_compose =[]
for i in range(1,11):
    x_train_c = functionLower(x_train, i)
    x_test_c = functionLower(x_test, i)
    w_c = linear_regression(x_train_c, y_train)
    temp = E_in_minus_E_out_bin(w_c, x_train_c, y_train, x_test_c, y_test,0)
    x_train_compose.append(temp)
print("The minimum of i is ",x_train_compose.index(min(x_train_compose))+1)
print("|E_in_bin - E_out_bin|:",min(x_train_compose))

print('q16')
temprandom=0
for _ in range(200):
    n = random.random()
    x_train_random =functionRandom(x_train, n)
    x_test_random = functionRandom(x_test, n)
    w_random = linear_regression(x_train_random, y_train)
    temprandom += E_in_minus_E_out_bin(w_random, x_train_random, y_train, x_test_random, y_test, 0)
print("the average |E_in_bin - E_out_bin|over 200 experiments :", temprandom/200)

```

```

In [145]: runfile('C:/Users/paddy/Desktop/NTU_course/2021_fall/HTML/homework3
wdir='C:/Users/paddy/Desktop/NTU_course/2021_fall/HTML/homework3')
q12
E_in_bin: 0.18100000000000013
E_out_bin: 0.5073333333333344
|E_in_bin - E_out_bin|: 0.32633333333333438
q13
E_in_bin: 0.05200000000000004
E_out_bin: 0.5096666666666677
|E_in_bin - E_out_bin|: 0.4576666666666677
q14
E_in_bin: 0.16800000000000012
E_out_bin: 0.5066666666666674
|E_in_bin - E_out_bin|: 0.3386666666666672
q15
The minimum of i is 3
|E_in_bin - E_out_bin|: 0.132333333333334057
q16
the average |E in bin - E out bin|over 200 experiments : 0.207596666666667237

```