

1. (b) I think I can use online learning to solve this problem

Training data can be collected as 1 month history attributes, such as views on website... and the ground truth is the predicted chance of making a purchase in the next η days.

We can collect data every month so that we can improve hypothesis through receiving data instances.

$$2. (E) y_n(t) W_t^T X_n(t) = y_n(t) [W_t^T + y_n(t) X_n(t) \cdot \eta_t] X_n(t) > 0$$

$$\Rightarrow y_n(t) \cdot W_t^T X_n(t) + y_n(t) \cdot y_n(t) \cdot X_n(t) \cdot X_n(t) \cdot \eta_t > 0$$

$$\Rightarrow y_n(t) \cdot \overbrace{y_n(t)}^{=1} \cdot X_n(t) \cdot X_n(t) \cdot \eta_t > -y_n(t) \cdot W_t^T X_n(t) \quad (y_n(t) = \pm 1, \therefore [y_n(t)]^2 = 1)$$

$$\Rightarrow \eta_t > \frac{-y_n(t) \cdot W_t^T X_n(t)}{\|X_n(t)\|^2} \quad \therefore \text{Choose (E).}$$

$$3. W_{t+1} \leftarrow W_t + y_n(t) X_n(t) \cdot \eta_t$$

$$W_f^T W_{t+1} = W_f^T (W_t + y_n(t) X_n(t) \cdot \eta_t)$$

$$\geq W_f^T W_t + \min_n y_n W_f^T X_n \eta_t$$

$$\|W_{t+1}\|^2 = \|W_t + y_n(t) X_n(t) \cdot \eta_t\|^2$$

$$= \|W_t\|^2 + 2 y_n(t) W_t^T X_n(t) \cdot \eta_t + \eta_t^2 \|y_n(t) X_n(t)\|^2$$

$$\leq \|W_t\|^2 + 0 + \max_n \|X_n\|^2 \cdot \eta_t^2$$

$$\text{Let } \min_n y_n W_f^T X_n = \rho, \max_n \|X_n\|^2 = R^2$$

$$W_f^T W_1 \geq W_f^T W_0 + \eta_0 \rho$$

$$\|W_1\|^2 \leq \|W_0\|^2 + R^2 \cdot \eta_0^2$$

$$W_f^T W_2 \geq W_f^T W_1 + \eta_1 \rho$$

$$\|W_2\|^2 \leq \|W_1\|^2 + R^2 \cdot \eta_1^2$$

$$\vdots$$

$$+) W_f^T W_T \geq W_f^T W_{T-1} + \eta_{T-1} \rho$$

$$\vdots$$

$$+) \|W_T\|^2 \leq \|W_{T-1}\|^2 + R^2 \eta_{T-1}^2$$

$$W_f^T W_T \geq \rho \cdot \sum_{t=0}^{T-1} \eta_t$$

$$\|W_T\|^2 \leq R^2 \cdot \sum_{t=0}^{T-1} \eta_t^2$$

$$\frac{W_f^T W_T}{\|W_f\| \|W_T\|} \geq \frac{\rho \cdot \sum_{t=0}^{T-1} \eta_t}{R \cdot \sqrt{\sum_{t=0}^{T-1} \eta_t^2}}$$

$$\text{if } \frac{\sum_{t=0}^{T-1} \eta_t}{\sqrt{\sum_{t=0}^{T-1} \eta_t^2}} \text{ is strongly increasing, that } W_T \text{ can be } W_f.$$

$$\text{For (A), } \frac{\sum_{t=0}^{T-1} 2^{-t}}{\sqrt{\sum_{t=0}^{T-1} 2^{-2t}}} = \frac{1 + \frac{1}{2} + \frac{1}{4} + \dots + \frac{1}{2^{T-1}}}{\sqrt{1 + \frac{1}{4} + \frac{1}{16} + \dots + \frac{1}{4^{T-1}}}}$$

$$= \frac{2 \cdot 1 \cdot (1 - \frac{1}{2^T})}{\sqrt{1 \cdot (1 - \frac{1}{4^T})}} = \sqrt{2} \cdot \frac{(1 - \frac{1}{2^T})}{\sqrt{(1 - \frac{1}{4^T})}}$$

since $(1 - \frac{1}{4^T}) > (1 - \frac{1}{2^T})$, both of them < 1 .

so, $(1 - \frac{1}{2^T}) < \sqrt{(1 - \frac{1}{4^T})} \Rightarrow \text{not increasing.}$

$$\text{For (B), } \frac{\sum_{t=0}^{T-1} 0.6211}{\sqrt{\sum_{t=0}^{T-1} 0.6211^2}} = \frac{T \cdot 0.6211}{0.6211 \cdot \sqrt{T}}$$

$\Rightarrow T > \sqrt{T} \Rightarrow \text{strongly increasing.}$

$$\text{For (D), } \frac{\sum_{t=0}^{T-1} (\frac{1}{1+t})}{\sqrt{\sum_{t=0}^{T-1} (\frac{1}{1+t})^2}} = \frac{1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{T}}{\sqrt{1 + (\frac{1}{2})^2 + (\frac{1}{3})^2 + \dots + (\frac{1}{T})^2}}$$

$$\text{Let } \sqrt{\quad} = k, k^2 = \frac{1^2 + (\frac{1}{2})^2 + \dots + (\frac{1}{T})^2 + 2(\frac{1}{2} \cdot \frac{1}{3} + \dots + \frac{1}{T(T-1)})}{1 + (\frac{1}{2})^2 + (\frac{1}{3})^2 + \dots + (\frac{1}{T})^2}$$

$\Rightarrow k^2 > 1 \Rightarrow k > 1 \Rightarrow \text{strongly increasing.}$

For (C), (E), from problem 2, we can

easily know that (C) cannot hold.

(E) can hold.

4 According to "Classol Handout" P.52.

start from $w_0 = 0$, after T mistake corrections,

$$1 \geq \frac{w_f^T w_t}{\|w_f\| \|w_t\|} \geq \frac{T\rho}{\|w_f\| \sqrt{T} R} \quad (\rho = \min_n y_n w_f^T X_n, R^2 = \max_n \|X_n\|^2)$$

$$\Rightarrow \|w_f\| \sqrt{T} R \geq T\rho \Rightarrow \|w_f\|^2 T \cdot R^2 \geq T^2 \rho^2 \Rightarrow T \leq \frac{\|w_f\|^2 R^2}{\rho^2}$$

In this problem, $f(x) = w_f^T \cdot X$, w_f can be $[-0.5 \ w_1 \ w_2 \ \dots \ w_d]$

$w_i = 1$ if w_i is spam-like

$= -1$ if w_i is not spam-like

$$X_n = [1, x_1, x_2, x_3, \dots, x_d]$$

$x_i = 1$ if word i is in email X_n

$= 0$ if word i is not in email X_n

$$\rho = \min_n y_n w_f^T X_n = -0.5 \quad \text{since } w_f^T X_n = -0.5 + \sum_{i=1}^d w_i x_i, y_n = \text{sign}(w_f^T X_n)$$

$$R^2 = \max_n \|X_n\|^2 = \max_n (1 + x_1^2 + x_2^2 + \dots + x_d^2) = (1+m) \quad \text{since there are at most } m \text{ distinct words in each email.}$$

$$\|w_f\|^2 = ((-0.5)^2 + d) = 0.25 + d$$

$$\therefore T \leq \frac{(0.25+d) \cdot (1+m)}{(0.5)^2} = (4d+1)(m+1) \#$$

5 Let WPLA's update equation: $w_{t+1} \leftarrow w_t + y_{n(t)} x_{n(t)}$

multiclass PLA update equation: $w_{y^{(t+1)}} \leftarrow w_{y^{(t)}} + x_{n(t)}$

$$w_{y^{(t+1)}} \leftarrow w_{y^{(t)}} + x_{n(t)}$$

Initial weight of WPLA & multiclass PLA are 0

When we find a mistake $x_{n(t)}$,

$$\text{if } y_{n(t)} = \begin{cases} +1 & \text{in binary PLA, update eq: } w_{t+1} \leftarrow w_t + x_{n(t)} \\ 2 & \text{in multiclass PLA} \end{cases}$$

$$\text{if } y_{n(t)} = \begin{cases} -1 & \text{in binary PLA} \\ 1 & \text{in multiclass} \end{cases}$$

$$w_{y=2}^{(t+1)} \leftarrow w_{y=2}^{(t)} + x_{n(t)}$$

$$w_{y=1}^{(t+1)} \leftarrow w_{y=1}^{(t)} - x_{n(t)}$$

$$w_{t+1} \leftarrow w_t - x_{n(t)}$$

$$w_{y=2}^{(t+1)} \leftarrow w_{y=2}^{(t)} - x_{n(t)}$$

$$w_{y=1}^{(t+1)} \leftarrow w_{y=1}^{(t)} + x_{n(t)}$$

conclusion:

Since initial weights are all 0, we do the same calculation as WPLA, we do the opposite calculation from WPLA

Thus, I choose (B) #

6 self-supervised learning is to learn 'physical knowledge' before actual tasks.

The images at similar time stamps should contain similar objects. (knowledge)

As a result, we can pair the images to train our model. (label of an image pair is the difference of time stamps.) \rightarrow self-defined goal

After the model learn that knowledge, images that are taken at similar time stamps can be mapped to similar vectors.

Also, this model can be a pretrained model of other tasks.

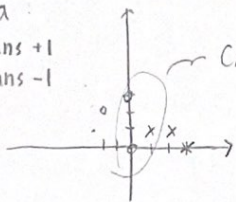
7. Each article can belong to several different categories. \rightarrow multilabel classification.

Having few (126) labeled data \rightarrow semi-supervised learning + batch learning.

Each tag can be labeled as 0 (not belong to an article) or 1 (belong to an article) \rightarrow raw features.

8. Data:

o means +1
x means -1

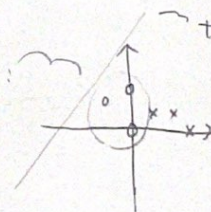


Choose these 3 examples as D . We can easily imagine that after PLA learning, a line that can separate these 3 examples can classify the other 3 examples correctly.

So, the smallest $E_{ots}(g) = 0$.

The largest $E_{ots}(g)$ would occurred when choosing 3 examples are all labeled as "+1" or "-1".

For example, choose 3 examples that labeled as +1:



then, we don't need to update w and we would misclassify all the other examples.

So, the largest $E_{ots}(g)$ would be 1. \neq

9. (A) $E[\hat{\theta}] = E\left(\frac{1}{N} \sum_{n=1}^N [h(X_n) + y_n]\right) = E_{X \sim p}[h(x) + f(x)]$

(B) $E[\hat{\theta}] = E\left(\frac{1}{N} \sum_{n=1}^N X_n\right)$, P is a Bernoulli distribution, $E(X) = \theta$

(D) $E[\hat{\theta}] = E\left(\frac{1}{N} \sum_{n=1}^N X_n^2\right) = E(X^2)$. \because Gaussian distribution, $\therefore E(X^2) = \text{Var}(X) + [E(X)]^2$

$\because P$ is zero-mean, $\therefore E(X^2) = \text{Var}(X) + 0 = \text{Var}(X) = \theta$

(C) Let $Y = \max\{X_1, X_2, \dots, X_N\}$

$$E(Y) = \sum_{y=1}^M y \cdot P(Y=y) = \sum_{y=1}^M y \cdot P(X_1 \leq y, X_2 \leq y, \dots, X_N \leq y) = \sum_{y=1}^M y \cdot \left(\frac{y}{M}\right)^N = \frac{1}{M^N} \sum_{y=1}^M y^{N+1}$$

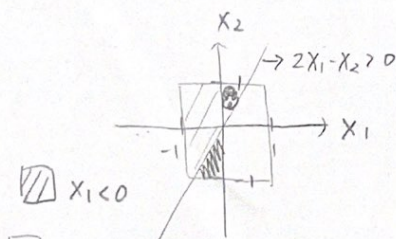
$$= \frac{1}{M^N} (1 + 2^{N+1} + 3^{N+1} + \dots + M^{N+1}) = \frac{1 + 2^{N+1} + \dots + (M-1)^{N+1}}{M^N} + M \neq M$$

Choose C $\#$

10. For $E_{\text{out}}(h_2)$, error occur when $\begin{cases} X_1 > 0, X_2 < 0 \rightarrow \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} \\ X_1 < 0, X_2 > 0 \rightarrow \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} \end{cases}$

$$\therefore E_{\text{out}}(h_2) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2}$$

For $E_{\text{out}}(h_1)$, error occur when $\begin{cases} 2X_1 - X_2 > 0, X_1 < 0 \text{ --- (1)} \\ 2X_1 - X_2 < 0, X_1 > 0 \text{ --- (2)} \end{cases}$

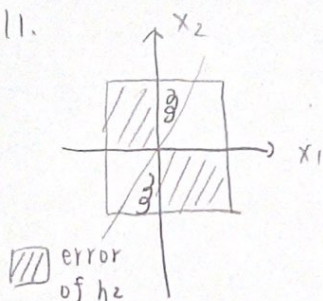


$$\text{Prob: } \frac{\frac{1}{2} \times \frac{1}{2}}{4} = \frac{1}{16} \text{ --- (1)}$$

$$\text{Prob: } \frac{\frac{1}{2} \times \frac{1}{2}}{4} = \frac{1}{16} \text{ --- (2)}$$

$$\therefore E_{\text{out}}(h_1) = \frac{1}{16} + \frac{1}{16} = \frac{1}{8} \#$$

11.



$E_{\text{in}}(h_2) = E_{\text{in}}(h_1)$ means the number of error is the same

When 0 error occur: $\left(1 - \frac{1}{2} - \frac{1}{8}\right)^4 = \frac{81}{4096}$

When 1 error occur: $\frac{1}{8} \times \frac{1}{2} \times \left(1 - \frac{1}{2} - \frac{1}{8}\right)^2 \times \frac{4!}{2!} = \frac{9 \times 4}{4096} \times 12 = \frac{432}{4096}$

When 2 error occur: $\frac{1}{8} \times \frac{1}{8} \times \frac{1}{2} \times \frac{1}{2} \times \frac{4!}{2! \times 2!} = \frac{16 \times 6}{4096} = \frac{96}{4096}$

Since there's no overlap, 3 and 4 errors are impossible

$$\text{Prob: } \frac{81}{4096} + \frac{432}{4096} + \frac{96}{4096} = \frac{609}{4096} \#$$

12. number	dices (green)
1	B (50)
2	A, B, D
3	D (50)
4	A (5A)
5	D (50)
6	ABC.

5A, 5B, 5C, 5D can get some number are purely green. $\rightarrow \left(\frac{1}{4}\right)^5 \times 4 = \frac{4}{1024}$

For number "2" is purely green:

$$4A \left[\begin{array}{l} 1B : \left(\frac{1}{4}\right)^5 \cdot \frac{5!}{4!} = \frac{5}{1024} \\ 1D : \left(\frac{1}{4}\right)^5 \cdot \frac{5!}{4!} = \frac{5}{1024} \end{array} \right.$$

$$3A \left[\begin{array}{l} 2B-0D : \left(\frac{1}{4}\right)^5 \cdot \frac{5!}{3!2!} = \frac{10}{1024} \\ 1B-1D : \left(\frac{1}{4}\right)^5 \cdot \frac{5!}{3!2!} = \frac{20}{1024} \\ 0B-2D : \left(\frac{1}{4}\right)^5 \cdot \frac{5!}{3!2!} = \frac{10}{1024} \end{array} \right.$$

$$2A \left[\begin{array}{l} 3B-0D : \left(\frac{1}{4}\right)^5 \cdot \frac{5!}{2!3!} = \frac{10}{1024} \\ 2B-1D : \left(\frac{1}{4}\right)^5 \cdot \frac{5!}{2!3!} = \frac{20}{1024} \\ 1B-2D : \left(\frac{1}{4}\right)^5 \cdot \frac{5!}{2!3!} = \frac{20}{1024} \\ 0B-3D : \left(\frac{1}{4}\right)^5 \cdot \frac{5!}{2!3!} = \frac{10}{1024} \end{array} \right.$$

$$1A \left[\begin{array}{l} 4B-0D = \frac{5}{1024} \\ 3B-1D = \frac{20}{1024} \\ 2B-2D = \frac{30}{1024} \\ 1B-3D = \frac{20}{1024} \\ 0B-4D = \frac{5}{1024} \end{array} \right.$$

$$0A \left[\begin{array}{l} 4B-1D = \frac{5}{1024} \\ 3B-2D = \frac{10}{1024} \\ 2B-3D = \frac{10}{1024} \\ 1B-4D = \frac{5}{1024} \end{array} \right.$$

$$\text{total: } \frac{240}{1024}$$

For num "6", also 3 kinds of dices (A, B, C).

\therefore similar to this situation, only minus duplicate scenario.

(only have "A", "B")

$$\text{Num "6"} : \frac{240 - 5 - 10 - 10 - 5}{1024} = \frac{210}{1024} \quad (\text{scenario that contains "0D"})$$

$$\text{Total: } \frac{4 + 240 + 210}{1024} = \frac{454}{1024} \#$$

```

import random
import numpy as np
from tqdm import tqdm

train_data = [np.array([1.0]+[float(x) for x in i.strip().split()]) for i in
open("hw1_train.dat").readlines())

#Q13
W_pla = []
for i in tqdm(range(1000)):
    r_seed = i*11
    random.seed(r_seed)
    w = [0.0] * 11
    w = np.array(w)
    accuracy = 0
    while True:
        data_id = random.randint(0,99)
        sign = 1.0 if np.sum(w * train_data[data_id][:11]) > 0 else -1.0
        if sign != train_data[data_id][-1]:
            w += train_data[data_id][-1]*train_data[data_id][:11]
            accuracy = 0
        else:
            accuracy += 1
            if accuracy == 500:
                W_pla.append(w)
                break
    length = 0
    for w_ in W_pla:
        length += np.sum(np.power(w_,2))
    print(length/1000)

#Q14
for idx in range(100):
    train_data[idx][:11] = train_data[idx][:11]*2
W_pla = []
for i in tqdm(range(1000)):
    r_seed = i*11
    random.seed(r_seed)
    w = [0.0] * 11
    w = np.array(w)
    accuracy = 0
    while True:
        data_id = random.randint(0,99)
        sign = 1.0 if np.sum(w * train_data[data_id][:11]) > 0 else -1.0
        if sign != train_data[data_id][-1]:
            w += train_data[data_id][-1]*train_data[data_id][:11]
            accuracy = 0
        else:
            accuracy += 1
            if accuracy == 500:
                W_pla.append(w)
                break
    length = 0
    for w_ in W_pla:
        length += np.sum(np.power(w_,2))
    print(length/1000)

#Q15
for idx in range(100):
    l = np.sqrt(np.sum(np.power(train_data[idx][:11],2)))
    train_data[idx][:11] = train_data[idx][:11]/l
W_pla = []
for i in tqdm(range(1000)):
    r_seed = i*11
    random.seed(r_seed)
    w = [0.0] * 11
    w = np.array(w)
    accuracy = 0
    while True:
        data_id = random.randint(0,99)
        sign = 1.0 if np.sum(w * train_data[data_id][:11]) > 0 else -1.0
        if sign != train_data[data_id][-1]:
            w += train_data[data_id][-1]*train_data[data_id][:11]
            accuracy = 0
        else:
            accuracy += 1
            if accuracy == 500:
                W_pla.append(w)
                break
    length = 0
    for w_ in W_pla:
        length += np.sum(np.power(w_,2))
    print(length/1000)

#Q16
train_data = [np.array([0.0]+[float(x) for x in i.strip().split()]) for i in
open("hw1_train.dat").readlines())

W_pla = []
for i in tqdm(range(1000)):
    r_seed = i*11
    random.seed(r_seed)
    w = [0.0] * 11
    w = np.array(w)
    accuracy = 0
    while True:
        data_id = random.randint(0,99)
        sign = 1.0 if np.sum(w * train_data[data_id][:11]) > 0 else -1.0
        if sign != train_data[data_id][-1]:
            w += train_data[data_id][-1]*train_data[data_id][:11]
            accuracy = 0
        else:
            accuracy += 1
            if accuracy == 500:
                W_pla.append(w)
                break
    length = 0
    for w_ in W_pla:
        length += np.sum(np.power(w_,2))
    print(length/1000)

```