

P1. D

We can use historical data as training data, ask a ML algorithm to learn from it and make future commercial behavior prediction.

P2. E

$$y_{n(t)} \underline{w_{t+1}^T} x_{n(t)} > 0 \quad (\text{代入 } [\epsilon] \text{ 选项})$$

$$\Rightarrow y_{n(t)} \left(w_t + y_{n(t)} x_{n(t)} \cdot \left[\frac{-y_{n(t)} w_t^T x_{n(t)}}{\|x_{n(t)}\|^2} + 1 \right] \right)^T x_{n(t)}$$

$$\Rightarrow y_{n(t)} w_t^T x_{n(t)} + y_{n(t)}^2 \cdot \|x_{n(t)}\|^2 \cdot \left[\frac{-y_{n(t)} w_t^T x_{n(t)}}{\|x_{n(t)}\|^2} + 1 \right]$$

令 $0 \leq \epsilon < 1$, 代入 ϵ 与 $y_{n(t)}^2 = 1$

$$\Rightarrow y_{n(t)} w_t^T x_{n(t)} + \|x_{n(t)}\|^2 \cdot \left(\frac{-y_{n(t)} w_t^T x_{n(t)}}{\|x_{n(t)}\|^2} + 1 - \epsilon \right)$$

$$\Rightarrow y_{n(t)} w_t^T x_{n(t)} + (-y_{n(t)} w_t^T x_{n(t)}) + \|x_{n(t)}\|^2 \cdot (1 - \epsilon)$$

$$\Rightarrow \|x_{n(t)}\|^2 \cdot (1 - \epsilon)$$

由於 $0 \leq \epsilon < 1$, 所以 $(1 - \epsilon) \geq 0$, 且 $\|x_{n(t)}\|^2$ 也大於零

$$\Rightarrow \|x_{n(t)}\|^2 \cdot (1 - \epsilon) > 0$$

$$\text{得證 } \underline{y_{n(t)} w_{t+1}^T x_{n(t)}} > 0 \quad \#$$

P3. C

設 $\eta(t)$ 為第 t 次更新時的更新率

inner product

$$w_f^T w_{t+1} \geq w_f^T w_t + \underbrace{(\min_n y_n w_f^T x_n)}_P \times \eta(t)$$

length of $\|w_t\|$

$$\|w_{t+1}\|^2 \geq \|w_t\|^2 + \underbrace{\max_n \|x_n\|^2}_R^2 \times \eta(t)$$

Magic Chain!

$$w_f^T w_1 \geq w_f^T w_0 + P \times \eta(0)$$

$$w_f^T w_2 \geq w_f^T w_1 + P \times \eta(1)$$

\vdots

$$w_f^T w_T \geq w_f^T w_{T-1} + P \times \eta(T-1)$$

$$\Rightarrow w_f^T w_T \geq w_f^T w_0 + P \cdot \sum_{i=0}^{T-1} \eta(i)$$

Start from $\vec{w}_0 = \vec{0}$

$$1 \geq \frac{w_f^T}{\|w_f\|} \cdot \frac{w_T}{\|w_T\|} > \frac{P \cdot \sum_{i=0}^{T-1} \eta(i)}{1 \cdot R \sqrt{\sum_{i=0}^{T-1} (\eta(i))^2}} \dots \dots \dots \textcircled{1}$$

[b] $\eta(t) = 0.621$

$$\text{代入式 } \textcircled{1} \Rightarrow 1 \geq \frac{P \cdot 0.621 T}{R \cdot \sqrt{0.621^2 \cdot T}} \Rightarrow 1 \geq \frac{P \cdot \sqrt{T}}{R} \Rightarrow \frac{R^2}{P^2} \geq T \text{ 收斂}$$

(第三題續)

$$[d] \eta(t) = \frac{1}{1+t}$$

$$\sum_{i=0}^{T-1} \eta(i) = 1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{T} \dots\dots\dots (2)$$

$$\sum_{i=0}^{T-1} \eta^2(i) = 1^2 + \left(\frac{1}{2}\right)^2 + \left(\frac{1}{3}\right)^2 + \dots + \left(\frac{1}{T}\right)^2 \dots\dots\dots (3)$$

代入①式

$$1 \geq \frac{p \cdot \sum_{i=0}^{T-1} \eta(i)}{1 \cdot R \sqrt{\sum_{i=0}^{T-1} (\eta(i))^2}} \geq \frac{p \cdot (1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{T})}{R \cdot (1^2 + (\frac{1}{2})^2 + (\frac{1}{3})^2 + \dots + (\frac{1}{T})^2)}$$

* 隨著 T 增加, ②. ③ 2 個數列差距會愈來愈大

當 $T \rightarrow \infty$ 時, 式② $\rightarrow \infty$ (發散) 而式③ $\simeq 1.6$

由此可知, 有某 T 使 $\frac{p \cdot (1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{T})}{R \cdot (1^2 + (\frac{1}{2})^2 + \dots + (\frac{1}{T})^2)}$ 超過 1, 因此 pla 會 halt.

(第三題續)

$$[e] \eta(t) = \left\lfloor \frac{-y_{n(t)} w_t^T x_{n(t)}}{\|x_{n(t)}\|^2} + 1 \right\rfloor$$

令 $1 \leq \varepsilon_{(t)} < 2$, 且因為 $w_t^T x_{n(t)}$ 一定與 $y_{n(t)}$ 異號, 故 $-y_{n(t)} w_t^T x_{n(t)} \geq 0$

$$\eta(t) = \frac{|w_t^T x_{n(t)}|}{\|x_{n(t)}\|^2} + \varepsilon_{(t)}$$

$$\begin{aligned} \sum_{i=0}^{T-1} \eta(i) &= \frac{|w_0^T x_{n(0)}|}{\|x_{n(0)}\|^2} + \varepsilon_{(0)} + \frac{|w_1^T x_{n(1)}|}{\|x_{n(1)}\|^2} + \varepsilon_{(1)} + \dots + \frac{|w_{T-1}^T x_{n(T-1)}|}{\|x_{n(T-1)}\|^2} + \varepsilon_{(T-1)} \\ &= \frac{|w_0^T x_{n(0)}|}{\|x_{n(0)}\|^2} + \frac{|w_1^T x_{n(1)}|}{\|x_{n(1)}\|^2} + \dots + \frac{|w_{T-1}^T x_{n(T-1)}|}{\|x_{n(T-1)}\|^2} + \underbrace{\varepsilon_{(0)} + \dots + \varepsilon_{(T-1)}}_{\geq 0} \geq 1 \cdot T \end{aligned}$$

$\geq T$ 因為每一項都為非負

$$\Rightarrow \sum_{i=0}^{T-1} \eta(i) \geq T$$

$$\begin{aligned} \sum_{i=0}^{T-1} \eta^2(i) &= \left(\frac{|w_0^T x_{n(0)}|}{\|x_{n(0)}\|^2} \right)^2 + 2 \varepsilon_{(0)} \cdot \frac{|w_0^T x_{n(0)}|}{\|x_{n(0)}\|^2} + \varepsilon_{(0)}^2 \\ &\quad + \left(\frac{|w_1^T x_{n(1)}|}{\|x_{n(1)}\|^2} \right)^2 + 2 \varepsilon_{(1)} \cdot \frac{|w_1^T x_{n(1)}|}{\|x_{n(1)}\|^2} + \varepsilon_{(1)}^2 \\ &\quad \vdots \\ &\quad + \left(\frac{|w_{T-1}^T x_{n(T-1)}|}{\|x_{n(T-1)}\|^2} \right)^2 + 2 \varepsilon_{(T-1)} \cdot \frac{|w_{T-1}^T x_{n(T-1)}|}{\|x_{n(T-1)}\|^2} + \varepsilon_{(T-1)}^2 \leq T \cdot 4 \end{aligned}$$

令此項為 K , 由於裡面每一項皆為非負, 故 $0 \leq K$

$$\Rightarrow \sum_{i=0}^{T-1} \eta^2(i) \leq K + 4T$$

(第三題續)

代入式

$$1 \geq \frac{\rho \cdot \sum_{i=0}^{T-1} \eta(i)}{1 \cdot R \sqrt{\sum_{i=0}^{T-1} (\eta(i))^2}} \geq \frac{\rho \cdot T}{R \cdot \sqrt{K+4T}}$$

$$\Rightarrow K+4T \geq \frac{\rho^2 T^2}{R^2}$$

$$\Rightarrow 0 \geq \frac{\rho^2}{R^2} T^2 - 4T - K$$

$$\Rightarrow 0 \geq T^2 - \frac{4R^2}{\rho^2} T - \frac{R^2}{\rho^2} K$$

$$\Rightarrow 0 \geq \underbrace{T^2 - \frac{4R^2}{\rho^2} T + \left(\frac{2R^2}{\rho^2}\right)^2}_{\text{perfect square}} - \left(\frac{2R^2}{\rho^2}\right)^2 - \frac{R^2 K}{\rho^2}$$

$$\Rightarrow \frac{2R^2}{\rho^2} + \frac{R^2 K}{\rho^2} \geq \left(T - \frac{2R^2}{\rho^2}\right)^2$$

$$\Rightarrow \frac{2R^2}{\rho^2} + \sqrt{\frac{2R^2}{\rho^2} + \frac{R^2 K}{\rho^2}} \geq T \quad \text{收敛} \quad \#$$

P4, C

By Slide P.41 $T \leq \left(\frac{R}{\rho}\right)^2$, where $R = \sqrt{\max_n \|x_n\|^2}$, $\rho = \min_n y_n w_f^T x_n$

求 R, 考慮一封 email 會使得 x_n 裡有最多 1 or -1, 這會是擁有最多 distinct word 的 email, 而一封 email 最多只能有 M 個 distinct word.

$$x_n = (\underbrace{1, 1, 1, 1, \dots, 1}_{x_0, \text{ } m \text{ 個 distinct word}}, 0, 0, 0, \dots, 0)$$

$$\|x_n\| = \sqrt{(m+1) \cdot 1^2} = \sqrt{m+1}, \quad R = \sqrt{m+1} = \underline{\sqrt{m+1}} \#$$

求 ρ , 考慮 w_f 根據題意,

$$w_f = (\underbrace{-0.5}_b, \underbrace{1, 1, 1, \dots, 1}_{d+}, \underbrace{-1, -1, -1, \dots, -1}_{d-})$$

↓ normalize

$$w_f' = w_f \times \frac{1}{\sqrt{\frac{1}{4} + d}}$$

考慮一封空白 email, 意即沒有任何 distinct word

$$x_n = (\underbrace{1}_{x_0}, \underbrace{0, 0, \dots, 0}_d)$$

$$w_f'^T \cdot x_n = (-0.5) \cdot \frac{1}{\sqrt{\frac{1}{4} + d}}, \quad y_n w_f'^T x_n = (-1) \cdot (-0.5) \cdot \frac{1}{\sqrt{\frac{1}{4} + d}} = \underline{\frac{\frac{1}{2}}{\sqrt{\frac{1}{4} + d}}} \#$$

$$\frac{R^2}{\rho^2} = \frac{m+1}{\frac{\frac{1}{4}}{\frac{1}{4} + d}} = 4(m+1)(d + \frac{1}{4}) = \underline{(m+1)(4d+1)} \#$$

P5, B

舉例而言, 假設 $y_{n(1)} = 1$, $y_{n(2)} = 2$ 且 $t=2$ 之後 PLA 就是 perfect line 了。

則 $\tilde{y}_{n(1)} = -1$, $\tilde{y}_{n(2)} = 1$ * $y'_{(1)} = 2$, $y'_{(2)} = 1$

$$\begin{cases} W_{PLA} = \vec{0} + \tilde{y}_{n(1)} X_{n(1)} + \tilde{y}_{n(2)} X_{n(2)} \\ W_1^* = \vec{0} + X_{n(1)} - X_{n(2)} \\ W_2^* = \vec{0} - X_{n(1)} + X_{n(2)} \end{cases} \Rightarrow \begin{cases} W_{PLA} = -X_{n(1)} + X_{n(2)} \\ W_1^* = X_{n(1)} - X_{n(2)} \\ W_2^* = -X_{n(1)} + X_{n(2)} \end{cases}$$

$$\Rightarrow W_{PLA} = -W_1^* = W_2^*$$

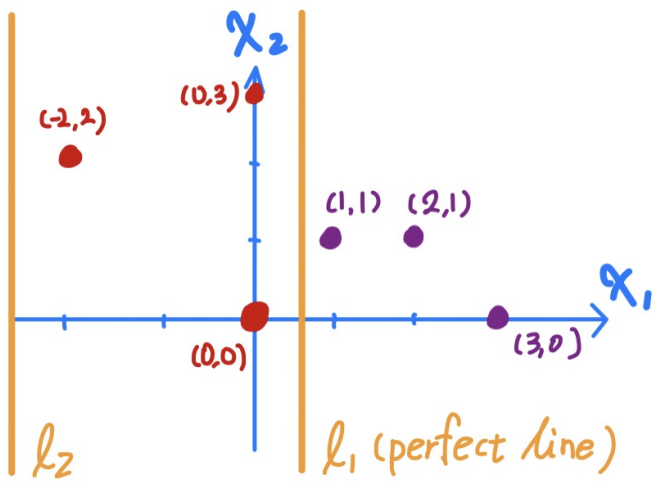
P6, D

The process isn't specifically labeled; moreover, it can be fine-tuned after learning from a simple task. Therefore, it match with the definition of self-supervised learning.

P7, C

It need to generate tag from an article, which is an application of multi-label classification. It's dataset is partially labeled, which means it's a semi-supervised learning. The algorithm will learn from all the data in a single training, which means its batch learning. The data is utf8 encoded strings which is raw features

P8, B



$$\text{設 } D = \{(1,1), (2,1), (3,0)\}$$

$$U/D = \{(-2,2), (0,3), (0,0)\}$$

$$\hat{y} = \text{sign}(w_1 x_1 + w_2 x_2 + b)$$

令 $w_1 = -1$, $w_2 = 0$, $b = 0.5$, 此為圖中 l_1

$$\hat{y} = \text{sign}(-w_1 + 0.5)$$

$$\left. \begin{array}{l} \text{代入 } (1,1), (2,1), (3,0), \hat{y} \text{ 皆為 } -1 \\ (0,3), (-2,2), (0,0), \hat{y} \text{ 皆為 } +1 \end{array} \right] \text{perfect line}$$

$$\underline{E_{ots} = 0} \#$$

令 $w_1 = -1$, $w_2 = 0$, $b = -3$, 此為圖中 l_2

$$\hat{y} = \text{sign}(-w_1 - 3)$$

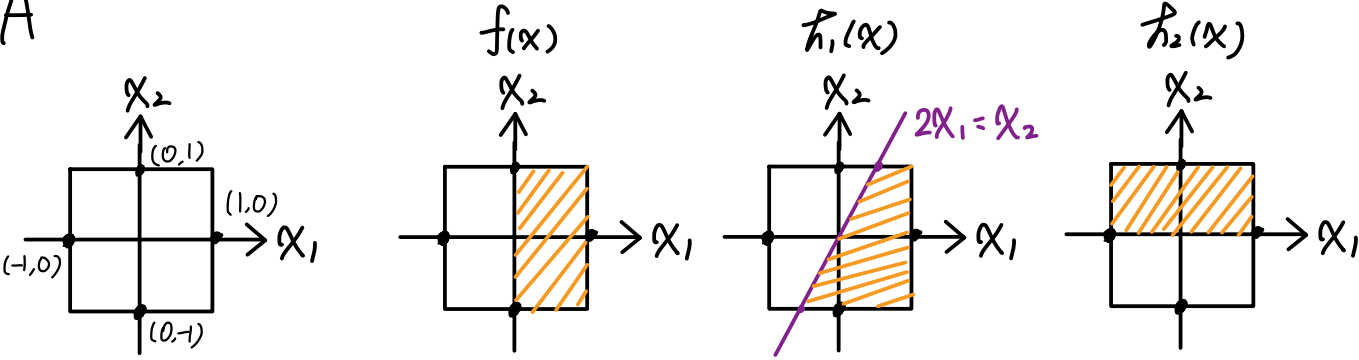
$$\left. \begin{array}{l} \text{代入 } (1,1), (2,1), (3,0), \hat{y} \text{ 皆為 } -1 \\ (0,3), (-2,2), (0,0), \hat{y} \text{ 皆為 } -1 \end{array} \right] U/D \text{ 全部做錯}$$

$$\underline{E_{ots} = 1} \#$$

P9. D

高斯分佈的 variance 為 $E[(X-\mu)^2]$, 題目說該分佈為 zero-mean ($\mu=0$)
 故 $\theta = E[X^2] = \frac{1}{N} \sum_{i=1}^N x_i^2 = \hat{\theta}$, where x 是從高斯分佈抽樣出來的

P10. A

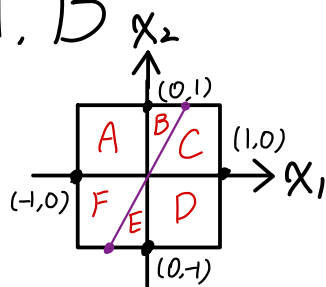


此為 $x = [x_1, x_2] \in \mathbb{R}^2$ Orange parts are $\text{sign}(\vec{x}) = 1$, blank parts are $\text{sign}(\vec{x}) = -1$

$$E_{\text{out}}(h_1) = \frac{\text{XOR of } f \text{ and } h_1}{\text{OR of } f \text{ and } h_1} = \frac{\frac{1}{2}}{\frac{2}{4}} = \frac{1}{8} \#$$

$$E_{\text{out}}(h_2) = \frac{\text{XOR of } f \text{ and } h_2}{\text{OR of } f \text{ and } h_2} = \frac{2}{4} = \frac{1}{2} \#$$

P11, B



將 $[-1, 1] \times [-1, 1]$ 切成 6 個區域, ABCDEF

	$[\hat{h}_1(x) = f(x)]$	$[\hat{h}_2(x) = f(x)]$	$P(x)$
A	1	-1	$1/4$
B	-1	1	$1/16$
C	1	1	$3/16$
D	1	-1	$1/4$
E	-1	1	$1/16$
F	1	1	$3/16$

$[\hat{h}_1(x) = f(x)]$	$[\hat{h}_2(x) = f(x)]$	$P(x)$
1	1	$3/8$
1	-1	$4/8$
-1	1	$1/8$
-1	-1	0

* $P(x)$ 為抽樣的 sample 在該區域的機率

考慮 $E_{in}(\hat{h}_1) = E_{in}(\hat{h}_2) = 0 \Rightarrow (\frac{3}{8})^4 = \frac{81}{4096}$

$E_{in}(\hat{h}_1) = E_{in}(\hat{h}_2) = 0.25 \Rightarrow (\frac{3}{8})^2 \times \frac{4}{8} \times \frac{1}{8} \times \frac{4!}{2!} = \frac{432}{4096} \Rightarrow \frac{81+432+96}{4096} = \frac{609}{4096} \#$

$E_{in}(\hat{h}_1) = E_{in}(\hat{h}_2) = 0.5 \Rightarrow (\frac{4}{8})^2 \times (\frac{1}{8})^2 \times \frac{4!}{2!2!} = \frac{96}{4096}$

P12, B

	1	2	3	4	5	6
A	X	O	X	O	X	O
B	O	O	X	X	X	O
C	X	X	X	X	X	O
D	X	O	O	X	O	X

* O 代表 green, X 代表 orange

考慮抽到的骰子的所有可能種類

一種: A, B, C, D

二種: AB, AC, AD, BC, BD, CD

三種: ABC, ABD, ACD, BCD

四種: ABCD

其中 ACD, BCD, CD, ABCD 不能產生 purely green number 故刪除, 所以剩下

一種: A, B, C, D $\Rightarrow 4$

二種: AB, AC, AD, BC, BD $\Rightarrow 5 \times (\frac{5!}{4!} + \frac{5!}{3!2!} + \frac{5!}{3!2!} + \frac{5!}{4!}) = 150$

三種: ABC, ABD $\Rightarrow 2 \times (3 \times \frac{5!}{3!} + 3 \times \frac{5!}{2!2!}) = 300$

$4 + 150 + 300 = 454$, 此為分子

$(4)^5 = 4096$, 此為分母

故 answer = $\frac{454}{4096} \#$

P13.B

```
import random
import numpy as np
from numpy import linalg as LA

N = 100 # Number of training data

# Load training data
X = [] # Training data
Y = [] # Labels
with open('hw1_train.dat', 'r') as f:
    for i in f.readlines():
        l = [1] # X_0
        for s in i.split('\t')[:-1]:
            l.append(float(s))
        X.append(np.array(l))
        Y.append(float(i.split('\t')[-1].split('\n')[0]))

def sign(x):
    if x > 0:
        return 1
    else:
        return -1

s = 0.0
for _ in range(1000): # Do 1000 times pla
    W = np.array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]) # weight [b, w1, w2, ..., w10]
    while True:
        is_done = True
        for _ in range(5*N): # Randomly find wrong point 5*N times
            # Randomly pick a point
            idx = random.randint(0, N-1)
            # Update weight
            if not ( sign(np.inner(W, X[idx])) == sign(Y[idx]) ):
                W = W + Y[idx]*X[idx]
                is_done = False
                break
        if is_done:
            break
    s += LA.norm(W)**2
print("Square norm of W = " + str(s/1000))
```

P14.C

```
import random
import numpy as np
from numpy import linalg as LA

N = 100 # Number of training data

# Load training data
X = [] # Training data
Y = [] # Labels
with open('hw1_train.dat', 'r') as f:
    for i in f.readlines():
        l = [2] # X_0
        for s in i.split('\t')[:-1]:
            l.append(float(s)*2)

        X.append(np.array(l))
        Y.append(float(i.split('\t')[-1].split('\n')[0]))

def sign(x):
    if x > 0:
        return 1
    else:
        return -1

s = 0.0
for _ in range(1000): # Do 1000 times pla
    W = np.array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0]) # weight [b, w1, w2, ..., w10]
    while True:
        is_done = True
        for _ in range(5*N): # Randomly find wrong point 5*N times
            # Randomly pick a point
            idx = random.randint(0, N-1)
            # Update weight
            if not ( sign(np.inner(W, X[idx])) == sign(Y[idx]) ):
                W = W + Y[idx]*X[idx]
                is_done = False
                break
        if is_done:
            break
    s += LA.norm(W)**2
print("Square norm of W = " + str(s/1000))
```

P15_E

```
import random
import numpy as np
from numpy import linalg as LA

N = 100 # Number of training data

# Load training data
X = [] # Training data
Y = [] # Labels
with open('hw1_train.dat', 'r') as f:
    for i in f.readlines():
        l = [1] # X_0
        for s in i.split('\t')[:-1]:
            l.append(float(s))

        # Use norm to normalize
        X.append(np.array(l)/LA.norm(np.array(l)))
        Y.append(float(i.split('\t')[-1].split('\n')[0]))

def sign(x):
    if x > 0:
        return 1
    else:
        return -1

s = 0.0
for _ in range(1000): # Do 1000 times pla
    W = np.array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0]) # weight [b, w1, w2, ..., w10]
    while True:
        is_done = True
        for _ in range(5*N): # Randomly find wrong point 5*N times
            # Randomly pick a point
            idx = random.randint(0,N-1)
            # Update weight
            if not ( sign(np.inner(W,X[idx])) == sign(Y[idx]) ):
                W = W + Y[idx]*X[idx]
                is_done = False
                break
        if is_done:
            break
    s += LA.norm(W)**2
print("Square norm of W = " + str(s/1000))
```

P16.A

```
import random
import numpy as np
from numpy import linalg as LA

N = 100 # Number of training data

# Load training data
X = [] # Training data
Y = [] # Labels
with open('hw1_train.dat', 'r') as f:
    for i in f.readlines():
        l = [0] # X_0
        for s in i.split('\t')[:-1]:
            l.append(float(s))

        X.append(np.array(l))
        Y.append(float(i.split('\t')[-1].split('\n')[0]))

def sign(x):
    if x > 0:
        return 1
    else:
        return -1

s = 0.0
for _ in range(1000): # Do 1000 times pla
    W = np.array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]) # weight [b, w1, w2, ..., w10]
    while True:
        is_done = True
        for _ in range(5*N): # Randomly find wrong point 5*N times
            # Randomly pick a point
            idx = random.randint(0, N-1)
            # Update weight
            if not ( sign(np.inner(W, X[idx])) == sign(Y[idx]) ):
                W = W + Y[idx]*X[idx]
                is_done = False
                break
        if is_done:
            break
    s += LA.norm(W)**2
print("Square norm of W = " + str(s/1000))
```