

Αναγνώριση Προτύπων - 2<sup>η</sup> Σειρά Ασκήσεων

Ακαδημαϊκό Έτος 2021-2022

Όνοματεπώνυμο: Χρήστος Δημόπουλος

Αρ. Μητρώου: 031 17 037

e-mail: chrisdim1999@gmail.com

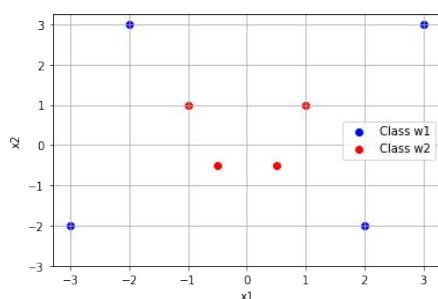
### Άσκηση 2.1: (SVM)

Μας δίνονται  $N=8$  διανύσματα χαρακτηριστικών  $\tilde{x}_n = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$  που προέρχονται από δύο κλάσεις  $w_1$  και  $w_2$ :

$$w_1: \left\{ \begin{pmatrix} 3 \\ 3 \end{pmatrix}, \begin{pmatrix} 2 \\ -2 \end{pmatrix}, \begin{pmatrix} -3 \\ -2 \end{pmatrix}, \begin{pmatrix} -2 \\ 3 \end{pmatrix} \right\}, z_1 = z_2 = z_3 = z_4 = -1$$

$$w_2: \left\{ \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 0.5 \\ -0.5 \end{pmatrix}, \begin{pmatrix} -0.5 \\ -0.5 \end{pmatrix}, \begin{pmatrix} -1 \\ 1 \end{pmatrix} \right\}, z_5 = z_6 = z_7 = z_8 = +1$$

α) Αρχικά, σχεδιάζοντας τα παραπάνω σημεία σε ένα χάρσημα, παρατηρούμε ότι οι δύο κλάσεις  $w_1$  και  $w_2$  δεν είναι γραμμικά διαχωρίσιμες.



Γι' αυτόν τον λόγο, μετατρέπουμε τα διανύσματα  $\tilde{x}_n$  σε έναν χώρο υψηλότερων διαστάσεων,  $y_n = \varphi(\tilde{x}_n)$ , χρησιμοποιώντας την εξής μορφή  $\varphi$ -functions 2<sup>ης</sup> τάξης:

$$\varphi(x_1, x_2) = \left[ 1, x_1, x_2, \frac{x_1^2 + x_2^2 - 5}{3} \right]^T$$

Το διάνυσμα βαρών  $W$  είναι, επίσης, επανυστημένο κατά  $w_0$  ( $W = [w_0 \tilde{W}]^T$ ) και καλούμαστε να ελαχιστοποιήσουμε το μήκος του, υπό τους περιορισμούς  $z_n W^T y_n \geq 1$ ,

β) Κατά την εκπαίδευση ενός SVM ταξινομητή, θέλουμε να ελαχιστοποιήσουμε το μέτρο  $\|w\|$ , υπό τους περιορισμούς  $z_n w^T y_n \geq 1$ ,  $n=1, \dots, 8$ . Για τον σκοπό αυτό, ορίζουμε το συνάρτησιακό:

$$L(w, a) = \frac{1}{2} \|w\|^2 - \sum_{k=1}^8 a_k [z_k w^T y_k - 1]$$

το οποίο θέλουμε να ελαχιστοποιήσουμε ως προς το διάνυσμα βαρών  $w$  και να μεγιστοποιήσουμε ως προς τους άγνωστους πολλαπλασιαστές  $a_k \geq 0$ . Το παραπάνω πρόβλημα είναι ισοδύναμο με τη μεγιστοποίηση ως προς τα  $a$  της ποσότητας:

$$L(a) = \sum_{k=1}^8 a_k - \frac{1}{2} \sum_{i=1}^8 \sum_{j=1}^8 a_i a_j z_i z_j y_i^T y_j$$

υπό τους περιορισμούς:

$$\sum_{k=1}^8 z_k a_k = 0 \quad \text{και} \quad a_k \geq 0, k=1, \dots, 8 \quad (I)$$

Το άνωθεν πρόβλημα βελτιστοποίησης, μπορεί να επιλυθεί με πολλαπλασιαστές Lagrange, θέτοντας το συνάρτησιακό:

$$Q(a, \lambda) = L(a) + \lambda \sum_{k=1}^8 z_k a_k$$

και εξισώνοντας τις παραγώγους του ως προς  $a_i$ ,  $i=1, \dots, 8$  και  $\lambda$  ίσες με μηδέν.

Στην πράξη, χρησιμοποιώντας πακέτα βελτιστοποίησης τις ρυθμίσεις βρίσκουμε:

$$\begin{aligned} a_1 &= 0, \quad a_2 = 0.1536, \quad a_3 = 0, \quad a_4 = 0.0934 \\ a_5 &= 0.1665, \quad a_6 = 0.017, \quad a_7 = 0.0175, \quad a_8 = 0.0459 \end{aligned}$$

Η παραπάνω λύση είναι δεκτή, αφού  $a_i \geq 0$ ,  $i=1, \dots, 8$  και επιπλέον

$$\text{ισχύει} \quad \sum_{i=1}^8 z_i a_i = 0.$$

γ) Αρχικά, προσδιορίζουμε το μη-επαυξημένο διάνυσμα βαρών  $w$  ως:

$$\tilde{w}_r = \sum_{k=1}^8 a_k z_k \tilde{y}_{k,r} \Rightarrow \tilde{w}_r = [-5 \times 10^{-5}, 0.2221, -0.6668]$$

Για την εύρεση του bias  $w_0$ , επιλέχουμε ένα οποιοδήποτε δείγμα που λειτουργεί ως support vector (δηλαδή  $a_i > 0$ ). Στην προκειμένη, χρησιμοποιούμε το

διάνυσμα  $y_2 = [1 \ 2 \ -2 \ 1]^T$ , για το οποίο βρήκαμε  $a_2 = 0.1536 > 0$ :

$$z_2 [w_0 \ \tilde{w}] \begin{bmatrix} 1 \\ \tilde{y}_2 \end{bmatrix} - 1 = 0$$

$$\Rightarrow z_2 \cdot (w_0 + \tilde{w} \tilde{y}_2) - 1 = 0 \Rightarrow w_0 = \frac{1}{z_2} - \tilde{w} \tilde{y}_2$$

Απ' όπου βρίσκουμε  $w_0 = 0.1112$ .

Τελικώς  $w = [0.1112, -5 \times 10^{-5}, 0.2221, -0.6668]^T$

Ελέγχουμε τους περιορισμούς:

$$\cdot z_1 w^T y_1 = 2.11 \geq 1$$

$$\cdot z_2 w^T y_2 = 1 \geq 1$$

$$\cdot z_3 w^T y_3 = 2.11 \geq 1$$

$$\cdot z_4 w^T y_4 \approx 1 \geq 1$$

$$\cdot z_5 w^T y_5 \approx 1 \geq 1$$

$$\cdot z_6 w^T y_6 \approx 1 \geq 1$$

$$\cdot z_7 w^T y_7 \approx 1 \geq 1$$

$$\cdot z_8 w^T y_8 \approx 1 \geq 1$$

δ) Το περιθώριο του ταξινομητή δίνεται ως :

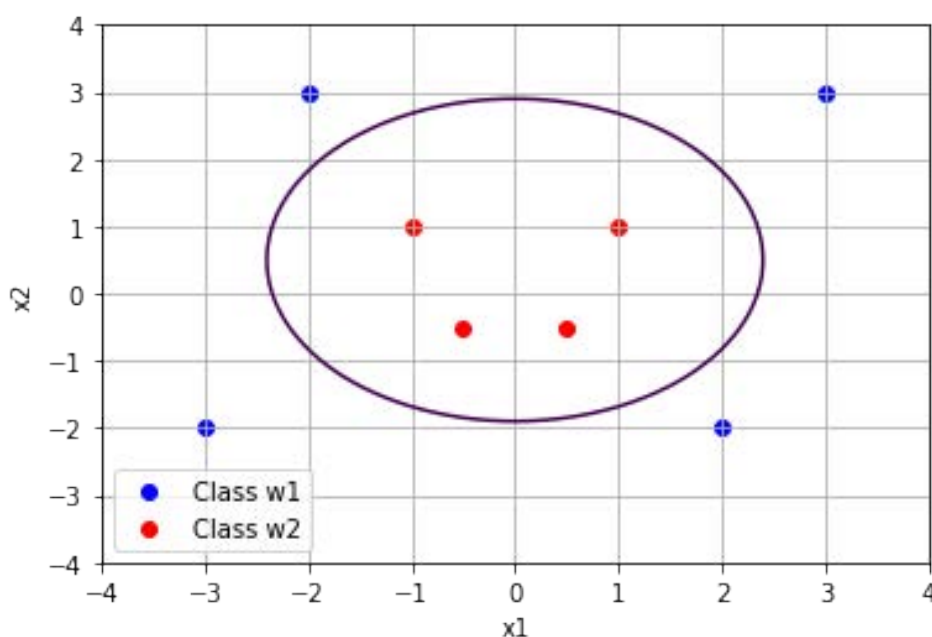
$$\beta = \frac{1}{\|\tilde{w}\|} \Rightarrow \boxed{\beta = 1.423}$$

ε) Η συνάρτηση διαχωρισμού :

$$g(x_1, x_2) = 0 \Rightarrow w^T \phi(x_1, x_2) = 0 \Rightarrow [w_0 \quad \tilde{w}] \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \frac{x_1^2 + x_2^2 - 5}{3} \end{bmatrix} = 0$$

$$\Rightarrow w_0 + w_1 x_1 + w_2 x_2 + w_3 \cdot \frac{x_1^2 + x_2^2 - 5}{3} = 0$$

απεικονίζεται παρακάτω μαζί με τα 8 αρχικά σημεία-δείγματα :



στ) Support Vectors είναι τα διανύσματα  $y_i$  για τα οποία βρήκαμε ότι ισχύει αυστηρώς  $\alpha_i > 0$ , δηλαδή τα διανύσματα  $y_2, y_4, y_5, y_6, y_7, y_8$  (ελαττών των  $y_1, y_3$ ).

ζ) Τα νέα δείγματα  $x' = \begin{pmatrix} \sqrt{2} \\ \sqrt{2} \end{pmatrix}$  και  $x'' = \begin{pmatrix} -2 \\ 3 \end{pmatrix}$  ταξινομούνται στις κατηγορίες

$w_2$  και  $w_1$  αντίστοιχα, καθώς :

$$w^T \cdot \begin{bmatrix} 1 \\ \sqrt{2} \\ \sqrt{2} \\ \frac{2+2-5}{3} \end{bmatrix} = 0.6475 > 0 \quad \downarrow \quad \text{Class } w_2$$

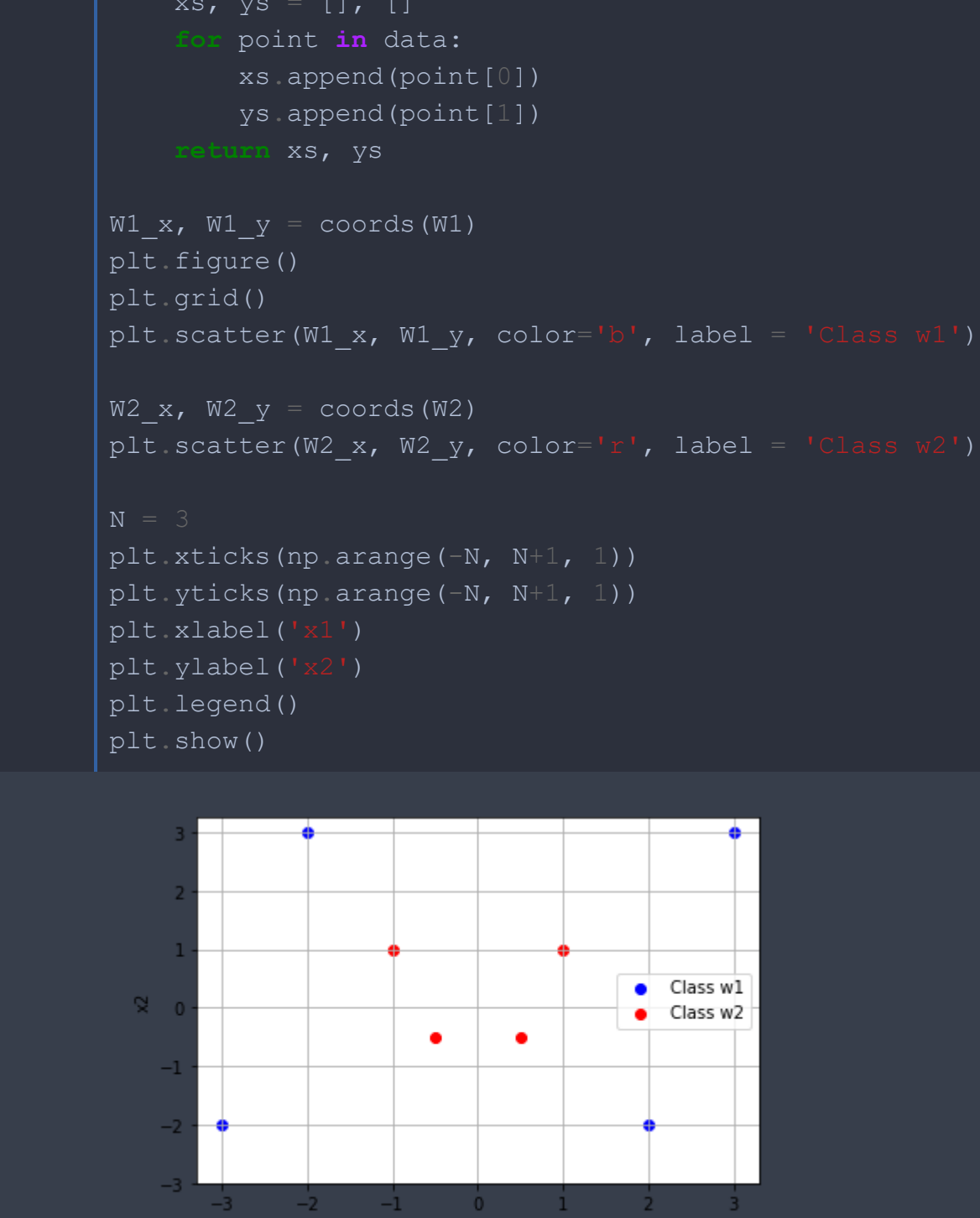
$$w^T \cdot \begin{bmatrix} 1 \\ -2 \\ 3 \\ \frac{4+9-5}{3} \end{bmatrix} = -1 < 0 \quad \downarrow \quad \text{Class } w_1$$

## Άσκηση 2.1: SVM

```
In [ ]: import numpy as np
import matplotlib.pyplot as plt
from scipy.optimize import minimize
```

### Ερώτημα (α)

Αρχικά, σχεδιάζουμε σε γράφημα τα σημεία εισόδου και παρατηρούμε ότι δεν είναι γραμμικά διαχωρίσιμα. Για τον λόγο αυτό, μετατρέπουμε τα διανύσματα  $\tilde{x}_n$  σε έναν χώρο υψηλότερων διαστάσεων,  $y_n = \phi(\tilde{x}_n)$ , χρησιμοποιώντας την εξής μορφή  $\phi$ -functions 2ης τάξης:  $\phi(x_1, x_2) = [\text{Big}(1 \setminus \setminus; \tilde{x}_1 \setminus \setminus; x_2 \setminus \setminus; \frac{x_1^2 + x_2^2 - 5}{3})\text{Big}]^T$



```
In [ ]: # Map to Higher Dimensionality
def fi_function(x):
    return [1, x[0], x[1], ((x[0])**2 + (x[1])**2 - 5)/3]

new_w1 = []
new_w2 = []
print('W1 Class')
for sample in W1:
    new_w1.append(np.array(fi_function(sample)))
    print(fi_function(sample))

print('\nW2 Class')
for sample in W2:
    new_w2.append(np.array(fi_function(sample)))
    print(fi_function(sample))

W1 Class
[1, 3, 3, 4.333333333333333]
[1, 2, -2, 1.0]
[1, -3, -2, 2.6666666666666665]
[1, -2, 3, 2.6666666666666665]

W2 Class
[1, 1, 1, -1.0]
[1, 0.5, -0.5, -1.5]
[1, -0.5, -0.5, -1.5]
[1, -1, 1, -1.0]
```

### Ερώτημα (β)

Ακολουθώς, με βάση το υποκεφάλαιο 5.11.1 [2], [5, SVM] προβαίνουμε στη μεγιστοποίηση του παρακάτω συναρτησιακού:

$$\mathcal{L}(\alpha) = \sum_{k=1}^N \alpha_j - \frac{1}{2} \sum_{k,j} \alpha_k \alpha_j \mathbf{z}_k^T \mathbf{z}_j \mathbf{y}_k \mathbf{y}_j$$

υπό τους περιορισμούς:

$$\sum_{k=1}^N \mathbf{z}_k^T \mathbf{z}_k \alpha_k = 0 \quad \text{and} \quad \alpha_k \geq 0, \quad k = 1, \dots, N$$

```
In [ ]: # Collec all samples in the new dimensional space
y = new_w1 + new_w2
z = [-1 for i in range(len(new_w1))] + [1 for i in range(len(new_w2))]
```

```
In [ ]: def L(a): # Cost Function
    result = 0
    for i in range(len(z)):
        result += a[i]

    for i in range(len(z)):
        for j in range(len(z)):
            result -= 0.5*a[i]*a[j]*z[i]*z[j]*np.dot(y[i],y[j])

    return -result
```

```
In [ ]: # Constrains on a_i
cons = ({'type': 'eq', 'fun': lambda a: np.dot(z,a)},
        {'type': 'ineq', 'fun': lambda a: a[0]},
        {'type': 'ineq', 'fun': lambda a: a[1]},
        {'type': 'ineq', 'fun': lambda a: a[2]},
        {'type': 'ineq', 'fun': lambda a: a[3]},
        {'type': 'ineq', 'fun': lambda a: a[4]},
        {'type': 'ineq', 'fun': lambda a: a[5]},
        {'type': 'ineq', 'fun': lambda a: a[6]},
        {'type': 'ineq', 'fun': lambda a: a[7]})
```

```
In [ ]: res = minimize(L, (0,0,0,0,0,0,0,0),
                      constraints=cons)
```

```
In [ ]: res
```

```
fun: -0.2469135792544758
jac: array([1.22247892, 0.11110054, 1.22219016, 0.11125482, -0.11112959,
          -0.11106048, -0.11104704, -0.11110269])
message: 'Optimization terminated successfully.'
```

```

nfev: 63
nit: 6
njev: 6
status: 0
success: True
x: array([-1.05270438e-17, 1.53551012e-01, -3.44416410e-16, 9.33680919e-02,
          1.66516910e-01, 1.69953160e-02, 1.74894126e-02, 4.59174647e-02])
```

```
In [ ]: # Optimal a_i
a_opt = np.round(res.x,4)
print(a_opt) # a_i >= 0 for each i
```

```
[-0.         0.1536 -0.         0.0934  0.1665  0.017  0.0175  0.0459]
```

```
In [ ]: 0 == np.round(np.dot(z, a_opt),3) # constraints satisfied
```

```
True
```

Παρατηρούμε ότι ικανοποιούνται οι περιορισμοί του προβλήματος βελτιστοποίησης.

### Ερώτημα (γ)

Το ζητούμενο διάνυσμα βαρών δίνεται από τη σχέση:  $\tilde{w} = \sum_{i=1}^N \alpha_i \tilde{z}_i$

```
In [ ]: weights = []
for j in range(1,len(y[1])):
    w = 0
    for i in range(len(z)):
        w += a_opt[i] * z[i] * y[i][j]
    weights.append(w)

In [ ]: print(weights)
```

```
[-4.999999999997368e-05, 0.22214999999999996, -0.6668166666666666]
```

Ακολουθώς, προσδιορίζουμε το bias  $w_0$  χρησιμοποιώντας το support vector  $y_2$  (για το οποίο ισχύει  $a_2 > 0$ ) στην ακόλουθη σχέση:  $\tilde{w}^T \tilde{y}_2 + w_0 - 1 = 0$  από την οποία προκύπτει:  $w_0 = \frac{1}{z_2} (\tilde{w}^T \tilde{y}_2 - 1)$

```
In [ ]: w0 = (1/z[1]) - np.dot(weights, y[1][1:]) # sample y2 is a Support Vector since a2 > 0
```

```
In [ ]: # Total Weight Vector with bias
weights = [w0] + weights
print(weights)
```

```
[0.11121666666666652, -4.999999999997368e-05, 0.22214999999999996, -0.6668166666666666]
```

Έπειτα, για το διάνυσμα  $w$  που βρήκαμε, ελέγχουμε ότι ισχύουν οι προϋποθέσεις  $\mathbf{z}_n^T \mathbf{w} \geq 1$  ( $n=1, \dots, N$ ):

```
In [ ]: for i in range(len(z)):
    print(z[i]*np.dot(weights, y[i])>=1)
```

```
True
True
True
True
True
True
True
True
```

### Ερώτημα (δ)

To margin  $\beta$  δίνεται από τον τύπο:  $\beta = \frac{1}{\|\tilde{w}\|}$

```
In [ ]: b = 1/(np.linalg.norm(weights[1:]))
print('Margin b =', b)
```

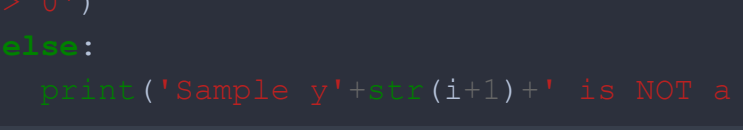
```
Margin b = 1.4227830509919526
```

### Ερώτημα (ε)

Βρίσκουμε τη συνάρτηση διαχωρισμού  $g(x_1, x_2) = \mathbf{w}^T \phi(x_1, x_2)$  στον αρχικό χώρο  $x_1, x_2$  και σχεδιάζουμε την καμπύλη  $g(x_1, x_2) = 0$  μαζί με τα 8 αρχικά σημεία.

```
In [ ]: x1 = np.linspace(-4, 4, 500)
x2 = np.linspace(-4, 4, 500)
X, Y = np.meshgrid(x1, x2)
g = weights[0] + weights[1]*X + weights[2]*Y + weights[3]*(X**2+Y**2-5)/3
```

```
In [ ]: plt.figure()
plt.grid()
plt.scatter(W1_x, W1_y, color='b', label = 'Class w1')
plt.scatter(W2_x, W2_y, color='r', label = 'Class w2')
plt.xlabel('x1')
plt.ylabel('x2')
plt.contour(X, Y, g, levels=[0])
plt.legend()
plt.show()
```



### Ερώτημα (στ)

Τα Support Vectors είναι τα σημεία με συντελεστή Lagrange αυστηρώς θετικό ( $\alpha_i > 0$ ), δηλαδή όλα εκτός από τα  $x_1$  και  $x_3$ .

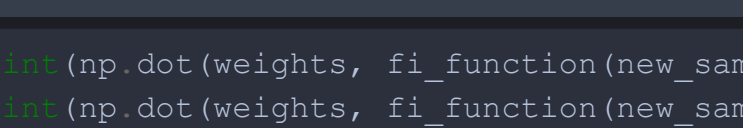
```
In [ ]: for i in range(len(y)):
    if (a_opt[i]>0):
        print('Sample y'+str(i+1)+' is a Support Vector, since a'+str(i+1)+' = '+str(a_opt[i])+' > 0')
    else:
        print('Sample y'+str(i+1)+' is NOT a Support Vector, since a'+str(i+1)+' = 0')
```

```
Sample y1 is NOT a Support Vector, since a1 = 0
Sample y2 is a Support Vector, since a2 = 0.1536 > 0
Sample y3 is NOT a Support Vector, since a3 = 0
Sample y4 is a Support Vector, since a4 = 0.0934 > 0
Sample y5 is a Support Vector, since a5 = 0.1665 > 0
Sample y6 is a Support Vector, since a6 = 0.0175 > 0
Sample y7 is a Support Vector, since a7 = 0.0175 > 0
Sample y8 is a Support Vector, since a8 = 0.0459 > 0
```

### Ερώτημα (ζ)

```
In [ ]: new_samples = [(np.sqrt(2), np.sqrt(2)), (-2, 3)]

In [ ]: plt.figure()
plt.grid()
plt.scatter(new_samples[0][0], new_samples[0][1], color='r', label = 'Class w2')
plt.scatter(new_samples[1][0], new_samples[1][1], color='b', label = 'Class w1')
plt.xlabel('x1')
plt.ylabel('x2')
plt.contour(X, Y, g, levels=[0])
plt.legend()
plt.show()
```



```
In [ ]: print(np.dot(weights, fi_function(new_samples[0]))) # > 0 --> Class w2
print(np.dot(weights, fi_function(new_samples[1]))) # < 0 --> Class w1
```

```
0.6475857210919529
-1.0004111111111111
```

Συνεπώς, ο SVM Classifier ταξινομεί το  $\sqrt{2}, \sqrt{2}$  στην κλάση  $\omega_2$  και το  $(-2, 3)$  στην κλάση  $\omega_1$ .

## Άσκηση 2.2 : (HMM)

Δίνεται HMM  $\lambda = (A, B, \pi)$  με τρεις καταστάσεις 1, 2, 3 και δύο τύπους παρατηρήσεων, H και T.

• Πίνακας Μεταβάσεων:  $A = \begin{pmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{pmatrix}$

• Πιθανότητες παρατηρήσεων B:

$P(o q)$			
$o/q$	1	2	3
H	0.5	0.75	0.25
T	0.5	0.25	0.75

• α-priori πιθανότητες:  $\pi_1 = \pi_2 = \pi_3 = 1/3$ .

• Παρατηρούμενη ακολουθία  $O = (H, T, H)$ .

1) Forward Algorithm: ορίζουμε  $a_t(i) = P(o_0 o_1 \dots o_t | q_t = i; \lambda)$  (forward probability).

• Initialisation ( $t=0$ ):  $a_0(i) = \pi_i b_i(O_0 = H)$

$$a_0(1) = \frac{1}{3} \cdot 0.5 = \frac{1}{6}$$

$$a_0(2) = \frac{1}{3} \cdot 0.75 = \frac{1}{4}$$

$$a_0(3) = \frac{1}{3} \cdot 0.25 = \frac{1}{12}$$

• Iterative Update:  $a_{t+1}(j) = \left[ \sum_{i=1}^N a_t(i) \cdot a_{ij} \right] \cdot b_j(o_{t+1})$

□ Για  $t=1$ :  $O_1 = T$

$$a_1(1) = \frac{1}{3} \cdot \frac{1}{3} \cdot (0.5 + 0.75 + 0.25) \cdot 0.5 = \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{3}{2} \cdot \frac{1}{2} = \frac{1}{12}$$

$$a_1(2) = \frac{1}{3} \cdot \frac{1}{3} \cdot (0.5 + 0.75 + 0.25) \cdot 0.25 = \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{3}{2} \cdot \frac{1}{4} = \frac{1}{24}$$

$$a_1(3) = \frac{1}{3} \cdot \frac{1}{3} \cdot (0.5 + 0.75 + 0.25) \cdot 0.75 = \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{3}{2} \cdot \frac{3}{4} = \frac{1}{8}$$



□ Για  $t=2$ :  $O_2 = H$

$$a_2(1) = \frac{1}{3} \cdot \left( \frac{1}{12} + \frac{1}{24} + \frac{1}{8} \right) \cdot 0.5 = \frac{1}{3} \cdot \frac{1}{4} \cdot \frac{1}{2} = \frac{1}{24}$$

$\frac{6}{24} = \frac{1}{4}$

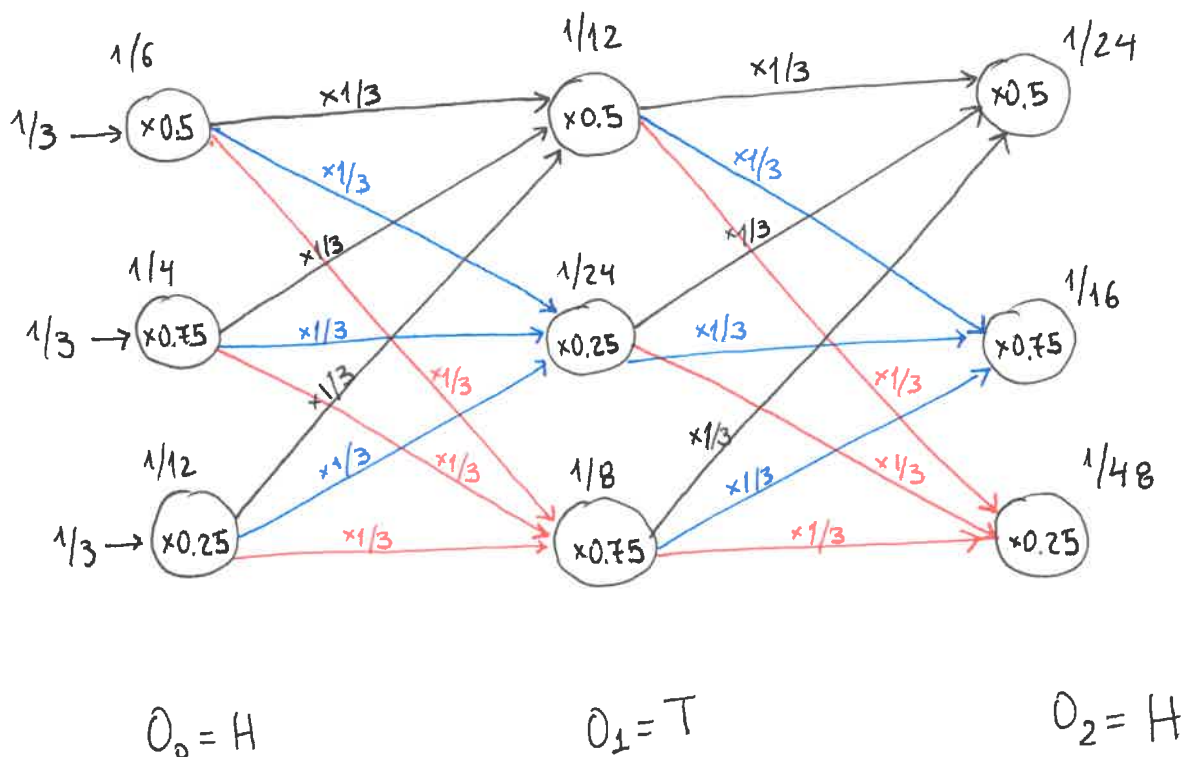
$$a_2(2) = \frac{1}{3} \cdot \left( \frac{1}{4} \right) \cdot 0.75 = \frac{1}{3} \cdot \frac{1}{4} \cdot \frac{3}{4} = \frac{1}{16}$$

$$a_2(3) = \frac{1}{3} \cdot \frac{1}{4} \cdot 0.25 = \frac{1}{3} \cdot \frac{1}{4} \cdot \frac{1}{4} = \frac{1}{48}$$

Τελικώς:

$$P(O=(H,T,H)|\lambda) = \sum_{i=1}^3 a_2(i) = \frac{1}{24} + \frac{1}{16} + \frac{1}{48} = \frac{2+3+1}{48} = \frac{6}{48} = \frac{1}{8}$$

$= 0.125.$



2) Backward Algorithm: Ορίζουμε  $\beta_t(i) = P(o_{t+1}, o_{t+2}, \dots, o_T | q_t = i; \lambda)$

• Initialization ( $t = T = 3$ ):  $\beta_T(i) = 1, 1 \leq i \leq N = 3$

$$\beta_3(1) = \beta_3(2) = \beta_3(3) = 1$$

• Iterative Process:  $\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j), 1 \leq i \leq N = 3, 1 \leq t \leq T = 3$

□ Για  $t = 2$ :  $O_3 = H$

$$\beta_2(1) = \frac{1}{3} \cdot 1 \cdot (0.5 + 0.75 + 0.25) = \frac{1}{3} \cdot \frac{3}{2} = \frac{1}{2}$$

$$\beta_2(2) = \frac{1}{3} \cdot 1 \cdot (0.5 + 0.75 + 0.25) = \frac{1}{2}$$

$$\beta_2(3) = \frac{1}{3} \cdot 1 \cdot (0.5 + 0.75 + 0.25) = \frac{1}{2}$$

□ Για  $t = 1$ :  $O_2 = T$

$$\beta_1(1) = \frac{1}{3} \cdot \frac{1}{2} \cdot (0.5 + 0.25 + 0.75) = \frac{1}{3} \cdot \frac{1}{2} \cdot \frac{3}{2} = \frac{1}{4}$$

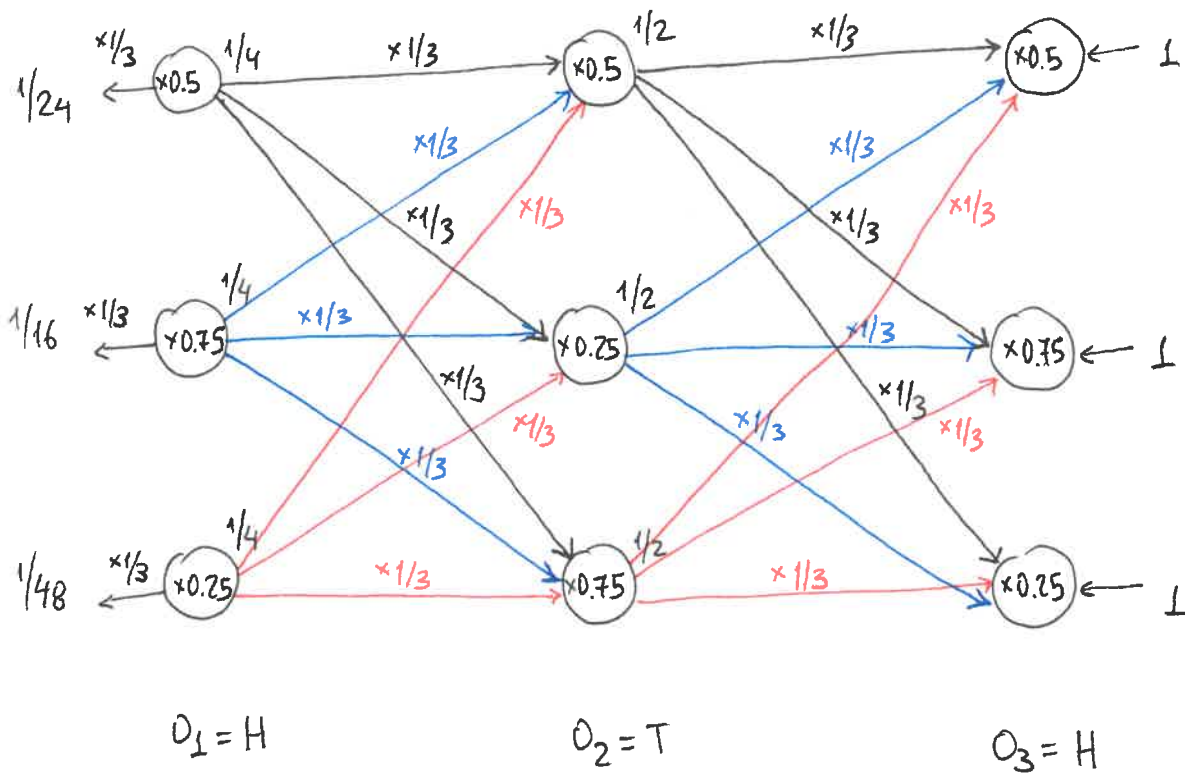
$$\beta_1(2) = \frac{1}{3} \cdot \frac{1}{2} \cdot \frac{3}{2} = \frac{1}{4}$$

$$\beta_1(3) = \frac{1}{3} \cdot \frac{1}{2} \cdot \frac{3}{2} = \frac{1}{4}$$

Τελικώς: ( $O_1 = H$ )

$$\begin{aligned} P(O = (H, T, H) | \lambda) &= \sum_{j=1}^3 \pi_j b_j(o_1) \beta_1(j) = \frac{1}{3} \times 0.5 \times \frac{1}{4} + \frac{1}{3} \times 0.75 \times \frac{1}{4} + \frac{1}{3} \times 0.25 \times \frac{1}{4} \\ &= \frac{1}{3} \times \frac{1}{4} \times \frac{3}{2} = \frac{1}{8} = \underline{\underline{0.125}} \end{aligned}$$





### 3) Αλγόριθμος Viterbi:

• Initialisation:  $\delta_0(i) = \pi_i b_i (O_0 = H)$

$$\delta_0(1) = 1/6, \quad \delta_0(2) = \frac{1}{4}, \quad \delta_0(3) = \frac{1}{12}$$

• Iterative Process:  $\delta_{t+1}(j) = \max_i (\delta_t(i) a_{ij}) \cdot b_j(O_{t+1})$

□  $\Gamma_{12} t=1$ :  $O_1 = T$

$$\delta_1(1) = \frac{1}{4} \cdot \frac{1}{3} \cdot 0.5 = \frac{1}{24}, \quad S_1(1) = 2$$

$$\delta_1(2) = \frac{1}{4} \cdot \frac{1}{3} \cdot 0.25 = \frac{1}{48}, \quad S_1(2) = 2$$

$$\delta_1(3) = \frac{1}{4} \cdot \frac{1}{3} \cdot 0.75 = \frac{1}{16}, \quad S_1(3) = 2$$

□  $\Gamma_{12} t=2$ :  $O_2 = H$

$$\delta_2(1) = \frac{1}{16} \cdot \frac{1}{3} \cdot 0.5 = \frac{1}{96}, \quad S_2(1) = 3$$

$$\delta_2(2) = \frac{1}{16} \cdot \frac{1}{3} \cdot 0.75 = \frac{1}{64}, \quad S_2(2) = 3$$

$$\delta_2(3) = \frac{1}{16} \cdot \frac{1}{3} \cdot 0.25 = \frac{1}{192}, \quad S_2(3) = 3$$

Συνεπώς, η πιο πιθανή ακολουθία καταστάσεων είναι η  $2 \rightarrow 3 \rightarrow 2$  με πιθανότητα

$$\delta_2(2) = \frac{1}{64} = 0.015625.$$

#### 4) Viterbi Training:

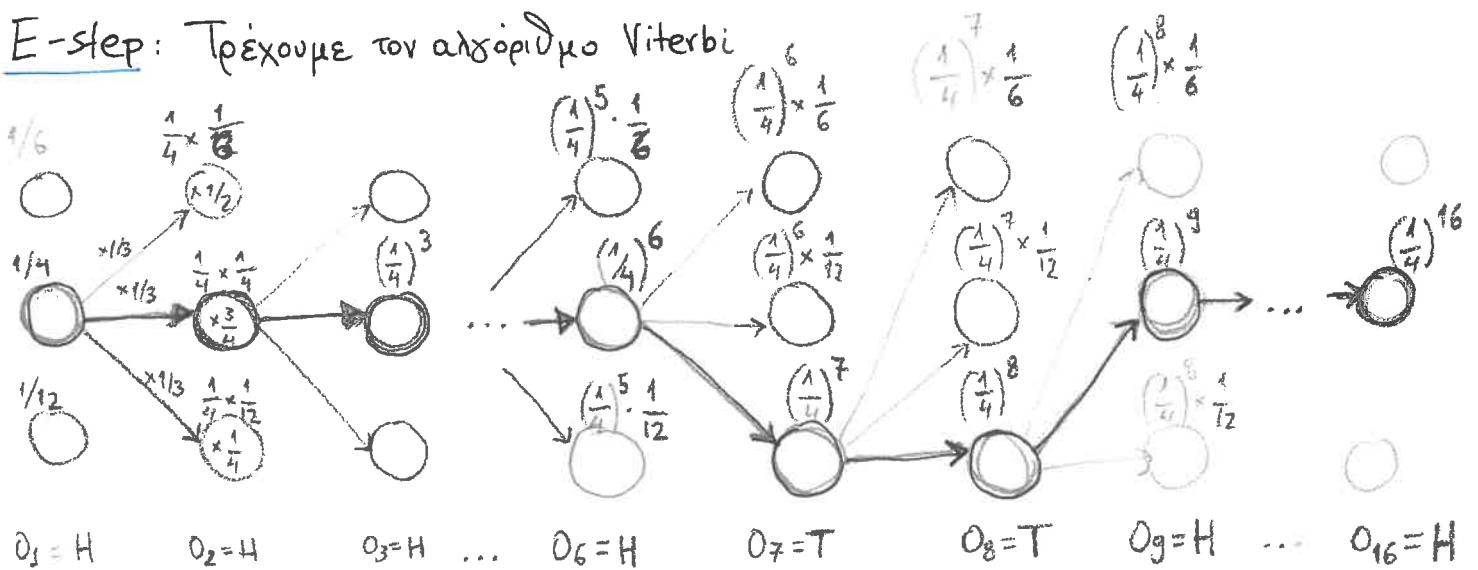
Παρατηρούμενη ακολουθία:  $O = H H H H H H T T H H H H H H H H$

##### • 1<sup>η</sup> Επανάληψη:

Αρχικοποιούμε το μοντέλο HMM με βάση τα δεδομένα της άσκησης, δηλαδή:

$$\lambda^{(0)} = (A, B, \pi)$$

E-step: Τρέχουμε τον αλγόριθμο Viterbi



Η πιο πιθανή ακολουθία καταστάσεων είναι:

2 2 2 2 2 2 3 3 2 2 2 2 2 2 2 2

με πιθανότητα  $(\frac{1}{4})^{16}$ .

M-step: Ενημερώνουμε τις ακόλουθες τιμές:

$$a_{22} = \frac{\#(2 \rightarrow 2)}{\#(2 \rightarrow *)} = \frac{12}{13}, \quad a_{23} = \frac{\#(2 \rightarrow 3)}{\#(2 \rightarrow *)} = \frac{1}{13}, \quad a_{21} = 0.$$

$$a_{33} = \frac{\#(3 \rightarrow 3)}{\#(3 \rightarrow *)} = \frac{1}{2}, \quad a_{32} = \frac{\#(3 \rightarrow 2)}{\#(3 \rightarrow *)} = \frac{1}{2}, \quad a_{31} = 0.$$

Επιπλέον:

$$P(O=H | q=2) = 1, \quad P(O=\pi | q=2) = 0$$

$$P(O=\pi | q=3) = 1, \quad P(O=H | q=3) = 0$$

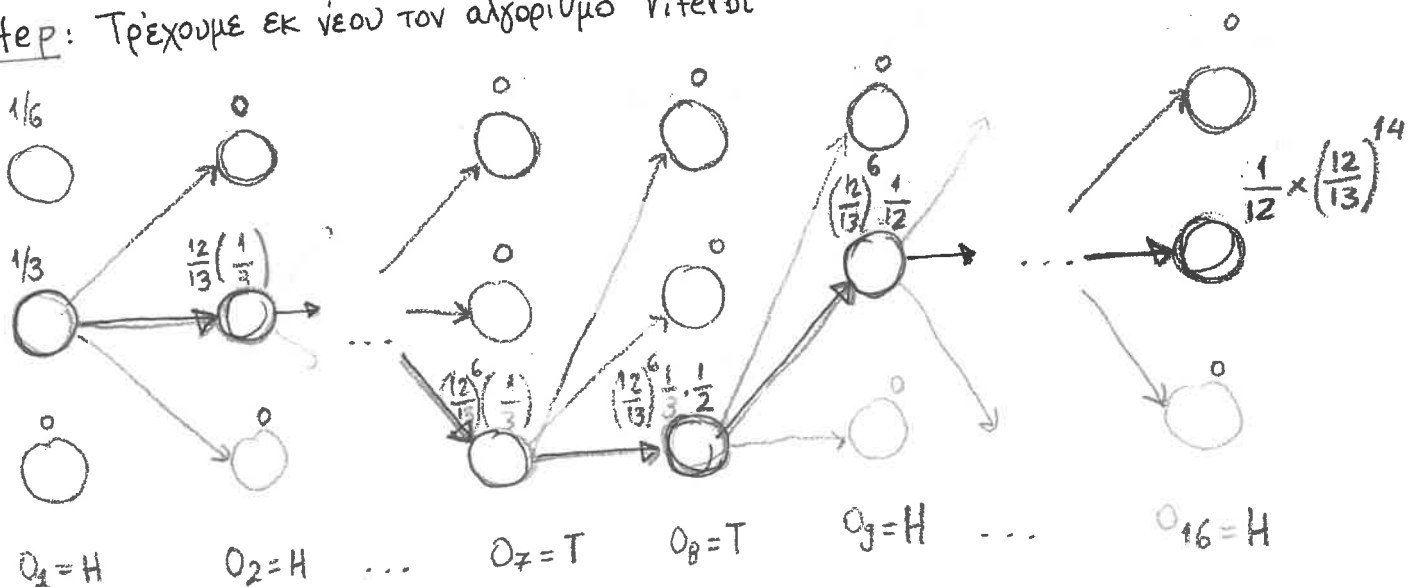
2<sup>η</sup> Επανάληψη:

Οι ανανεωμένες παράμετροι του HMM θα είναι  $\lambda^{(1)} = (A^{(1)}, B^{(1)}, \pi)$ , όπου:

$$A^{(1)} = \begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 0 & 12/13 & 1/13 \\ 0 & 1/2 & 1/2 \end{bmatrix} \text{ και } B^{(1)}:$$

$P(O q)$			
$O/q$	1	2	3
H	0.5	1	0
T	0.5	0	1

E-step: Τρέχουμε εκ νέου τον αλγόριθμο Viterbi



Η πιο πιθανή ακολουθία καταστάσεων είναι και πάλι  $\pi$ :

$$2 \rightarrow 2 \rightarrow 2 \rightarrow 2 \rightarrow 2 \rightarrow 2 \rightarrow 3 \rightarrow 3 \rightarrow 2 \rightarrow 2 \rightarrow 2 \rightarrow 2 \rightarrow 2 \rightarrow 2 \rightarrow 2 \rightarrow 2$$

με πιθανότητα  ~~$\left(\frac{12}{13}\right)^{16}$~~   $\cdot \frac{1}{12} \times \left(\frac{12}{13}\right)^{14}$ .

Ως εκ τούτου, οι παράμετροι του μοντέλου δεν ανανεώνονται και ο αλγόριθμος

pseudo-EM έχει συγκλίνει με  $\lambda = (A^{(1)}, B^{(1)}, \pi)$ .

5) Ο forward-backward αλγόριθμος χρησιμοποιεί κάθε πιθανό μονοπάτι καταστάσεων για την εκπαίδευση του HMM, σε αντίθεση με τον pseudo-EM αλγόριθμο, ο οποίος χρησιμοποιεί το πιο πιθανό μονοπάτι καταστάσεων και μόνο. Το πιθανότερο μονοπάτι προσδιορίζεται με χρήση του αλγορίθμου Viterbi με αποδοτικό τρόπο, πολυπλοκότητας  $O(N \times T)$ , όπου  $N$ : number of hidden states,  $T$ : number of ~~possible~~ observations

Συνοψίως:

- Ο forward-backward αλγόριθμος εκπαίδευσης είναι πιο αργός, αλλά δίνει πιο ακριβή αποτελέσματα στην εκπαίδευση του HMM, σε αντίθεση με τον pseudo-EM (Viterbi training), ο οποίος είναι πιο γρήγορος, αλλά λιγότερο ακριβής.

## Άσκηση 2.3: (CART)

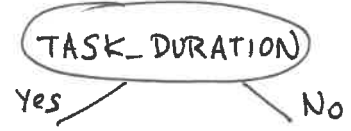
SATISFACTION	WORD ACCURACY	TASK_COMPLETION	TASK_DURATION
Y	100	Y	3
Y	100	Y	2
Y	90	N	4
N	95	N	2
N	80	Y	5
Y	85	Y	5
Y	80	Y	1
N	85	N	3
Y	95	N	4

1) Οι δύο κατηγορίες του προβλήματος ταξινόμησης είναι οι εξής:

$W_1$ : Τα δείγματα-χρήστες που έμειναν ευχαριστημένοι από το σύστημα διαλόγου,  
δηλαδή  $SATISFACTION = Y$

$W_2$ : τα δείγματα-χρήστες που δεν έμειναν ευχαριστημένοι από το σύστημα διαλόγου,  
δηλαδή  $SATISFACTION = N$

Οι παράμετροι του δέντρου απόφασης είναι ερωτήσεις - κόμβοι της εξής μορφής:



όπου  $\alpha, \beta \in \mathbb{R}$ .

2) Ελέγχουμε τις ερωτήσεις του δέντρου απόφασης με βάση το ποίες δίνουν την ελάχιστη εντροπία:

• Ερώτηση Task-Completion = Y:



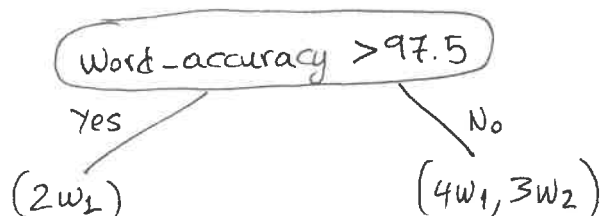
$$i_e(N_L) = -\left(\frac{4}{5} \log_2 \left(\frac{4}{5}\right) + \frac{1}{5} \log_2 \left(\frac{1}{5}\right)\right) = 0.722$$

$$i_e(N_R) = -(-0.5 \log_2 0.5 + 0.5 \log_2 0.5) = 1$$

$$i_e(\text{Task-completion} = Y) = \frac{5}{9} i_e(N_L) + \frac{4}{9} i_e(N_R) = 0.8455$$

• Ερώτηση Word-accuracy > a:

$$\triangleright a = \frac{95+100}{2} = 97.5$$

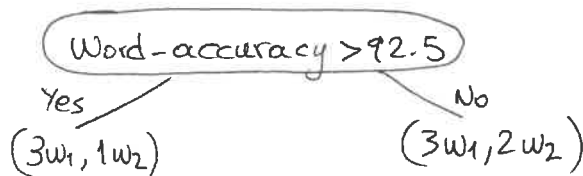


$$i_e(N_L) = 0$$

$$i_e(N_R) = -\left(\frac{4}{7} \log_2 \left(\frac{4}{7}\right) + \frac{3}{7} \log_2 \left(\frac{3}{7}\right)\right)$$

$$i_e(\text{Word-acc} > 97.5) = 0 + \frac{7}{9} i_e(N_R) = 0.7663$$

$$\triangleright a = \frac{90+95}{2} = 92.5$$

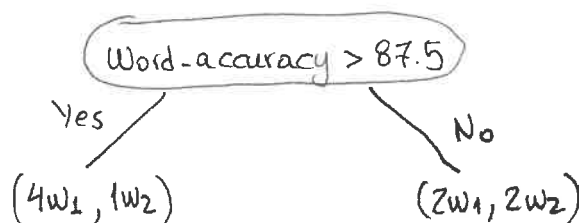


$$i_e(N_L) = -\left(\frac{3}{4} \log_2 \left(\frac{3}{4}\right) + \frac{1}{4} \log_2 \left(\frac{1}{4}\right)\right)$$

$$i_e(N_R) = -\left(\frac{3}{5} \log_2 \left(\frac{3}{5}\right) + \frac{2}{5} \log_2 \left(\frac{2}{5}\right)\right)$$

$$i_e(\text{Word-acc} > 92.5) = \frac{4}{9} i_e(N_L) + \frac{5}{9} i_e(N_R) = 0.9$$

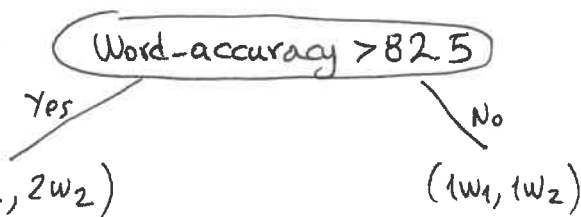
$$\triangleright a = \frac{85+90}{2} = 87.5$$



$$i_e(\text{Word-acc} > 87.5) = 0.8455$$



$$a = \frac{80+85}{2} = 82.5$$



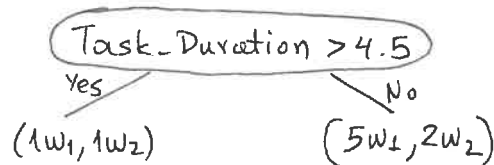
$$i_c(N_L) = -\left(\frac{5}{7}\log_2\left(\frac{5}{7}\right) + \frac{2}{7}\log_2\left(\frac{2}{7}\right)\right) \quad (5w_1, 2w_2)$$

$$i_c(N_R) = -(0.5\log_2(0.5) \times 2) = 1 \quad (1w_1, 1w_2)$$

$$i_c(\text{Word-acc} > 82.5) = \frac{7}{9} i_c(N_L) + \frac{2}{9} i_c(N_R) = 0.8935$$

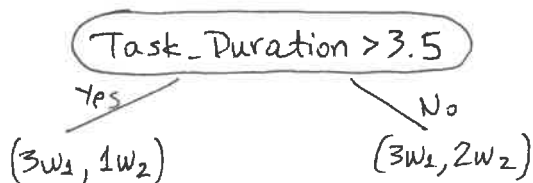
Ερώτηση Task-Duration > a:

$$a = 4.5$$



$$i_c(\text{Task-Duration} > 4.5) = 0.8935$$

$$a = 3.5$$



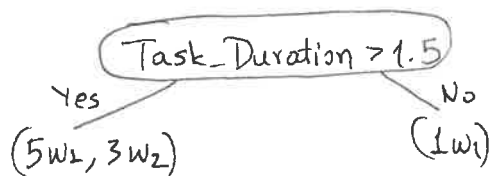
$$i_c(\text{Task-Duration} > 3.5) = 0.9$$

$$a = 2.5$$



$$i_c(\text{Task-Duration} > 2.5) = 0.9183$$

$$a = 1.5$$



$$i_c(\text{Task-Duration} > 1.5) = 0.8484$$

Συνεπώς, η χαμηλότερη εντροπία δίνεται για την ερώτηση  $\text{Word-accuracy} > 97.5$

( $i_c = 0.7663$ ), γι' αυτό και την επιλέγουμε ως τη ρίζα του δέντρου απόφασης.

Τώρα ακολουθούμε την ίδια διαδικασία, ώστε να βρούμε τη λιγότερο impure ερώτηση για τα δείγματα  $(4w_1, 3w_2)$  του δεξιού κλαδιού. (όπου  $\text{Word-acc} \leq 97.5$ ).

Epiwrtion Task-completion = Y :



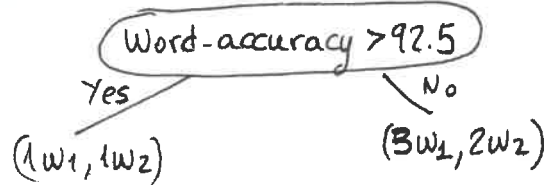
$$i_e(N_L) = -\left(\frac{2}{3}\log_2\left(\frac{2}{3}\right) + \frac{1}{3}\log_2\left(\frac{1}{3}\right)\right) = 0.918$$

$$i_e(N_R) = -2(0.5\log_2 0.5) = 1$$

$$I(\text{task-compl.} = Y) = \frac{3}{7} i_e(N_L) + \frac{4}{7} i_e(N_R) = 0.9648$$

Epiwrtion Word-accuracy > a :

$$\triangleright a = \frac{90+95}{2} = 92.5 :$$

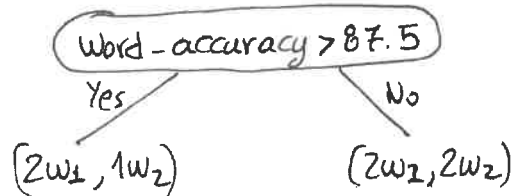


$$i_e(N_L) = 1$$

$$i_e(N_R) = -\left[\frac{3}{5}\log_2\left(\frac{3}{5}\right) + \frac{2}{5}\log_2\left(\frac{2}{5}\right)\right] = 0.971$$

$$i_e(\text{Word-accuracy} > 92.5) = \frac{2}{7} i_e(N_L) + \frac{5}{7} i_e(N_R) = 0.9792$$

$$\triangleright a = \frac{85+90}{2} = 87.5 :$$



$$i_e(\text{Word-acc} > 87.5) = 0.9648$$

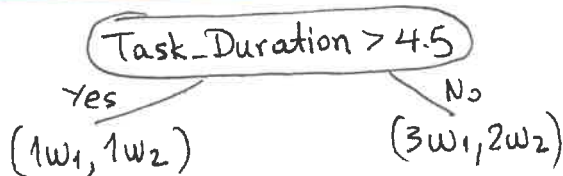
$$\triangleright a = \frac{80+85}{2} = 82.5 :$$



$$i_e(\text{Word-acc} > 82.5) = 0.971$$

Epiwrtion Task-Duration > a :

$$\triangleright a = 4.5 :$$



$$i_e(T.D. > 4.5) = 0.971$$

$$\triangleright a = 3.5 :$$

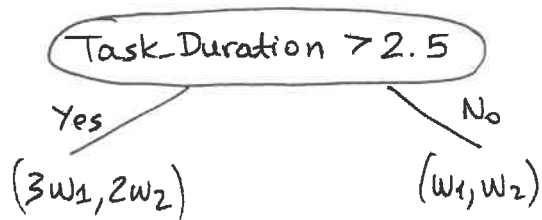


$$i_e(N_L) = -\left(\frac{3}{4}\log_2\frac{3}{4} + \frac{1}{4}\log_2\frac{1}{4}\right) = 0.811$$

$$i_e(N_R) = -\left(\frac{1}{3}\log_2\frac{1}{3} + \frac{2}{3}\log_2\frac{2}{3}\right) = 0.918$$

$$i_e(\text{Task-Duration} > 3.5) = \frac{4}{7} i_e(N_L) + \frac{3}{7} i_e(N_R) = 0.8568$$

$\triangleright a = 2.5 :$

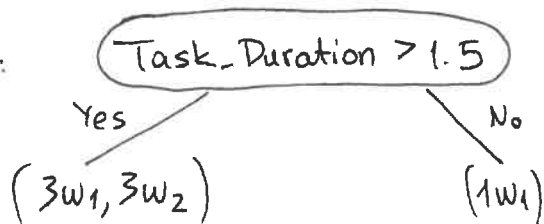


$$i_e(N_L) = 0.971$$

$$i_e(N_R) = 1$$

$$i_e(T.D. > 2.5) = \frac{5}{7} \times 0.971 + \frac{2}{7} = 0.9792$$

$\triangleright a = 1.5 :$



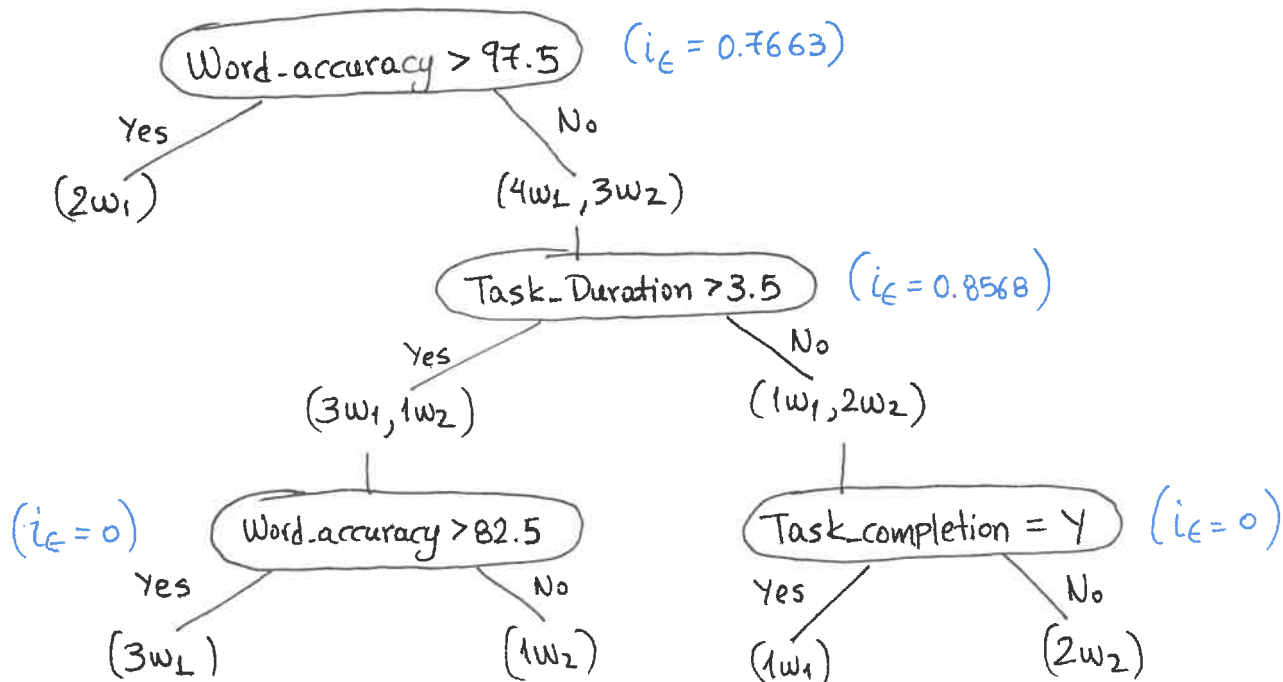
$$i_e(N_L) = -2 \times 0.5 \log_2 0.5 = 1$$

$$i_e(N_R) = 0$$

$$i_e(T.D. > 1.5) = \frac{6}{7} \times 1 = 0.8571$$

Συνεπώς, η χαμηλότερη εντροπία - άρα και το υψηλότερο purity - δίνεται από την ερώτηση  $Task\_Duration > 3.5$  ( $i_e = 0.8568$ ), χι' αυτό και την είχαμε στο δέντρο απόφασης.

Έπειτα, για το αριστερό κλαδί μπορούμε να χρησιμοποιήσουμε την ερώτηση  $Word\_accuracy > 82.5$  που θα διαχωρίσει πλήρως τα δείγματα ( $i_e = 0$ ), ενώ για το δεξιό κλαδί, μπορούμε να διαχωρίσουμε πλήρως τα δείγματα με την ερώτηση  $Task\_completion = Y$  ( $i_e = 0$ ). Το δέντρο απόφασης, λοιπόν, θα είναι το εξής:



3) Για να ταξινομήσουμε τις ερωτήσεις του δέντρου ως προς τη σημαντικότητα τους, χρησιμοποιούμε τη μετρική Gini Importance [1]:

$$n_{ij} = W_j C_j - W_{\text{left}(j)} C_{\text{left}(j)} - W_{\text{right}(j)} C_{\text{right}(j)}$$

όπου  $W_j$ : weighted number of samples reaching node  $j$

$C_j$ : impurity value of node  $j$

Επομένως:

$$\cdot n_i(\text{Word\_acc} > 97.5) = 1 \cdot 0.7663 - \frac{2}{9} \cdot 0 - \frac{7}{9} \cdot 0.8568 = 0.0999$$

$$\cdot n_i(\text{Task\_Duration} > 3.5) = \frac{7}{9} \cdot 0.8568 - \frac{4}{7} \cdot 0 - \frac{3}{7} \cdot 0 = 0.6664$$

$$\cdot n_i(\text{Word\_accuracy} > 82.5) = n_i(\text{Task\_completion} = \gamma) = 0$$

Δηλαδή, η σειρά σημαντικότητας των ερωτήσεων είναι:

$$\bullet (\text{Task\_Duration} > 3.5) > (\text{Word\_acc} > 97.5) > (\text{Task\_compl.} = \gamma) = (\text{Word\_acc} > 82.5)$$

Η σημαντικότητα των features, ορίζεται ως:

$$f_i = \frac{\sum_{j: \text{node } j \text{ splits on feature } i} n_{ij}}{\sum_{k \in \text{all nodes}} n_{ik}}$$

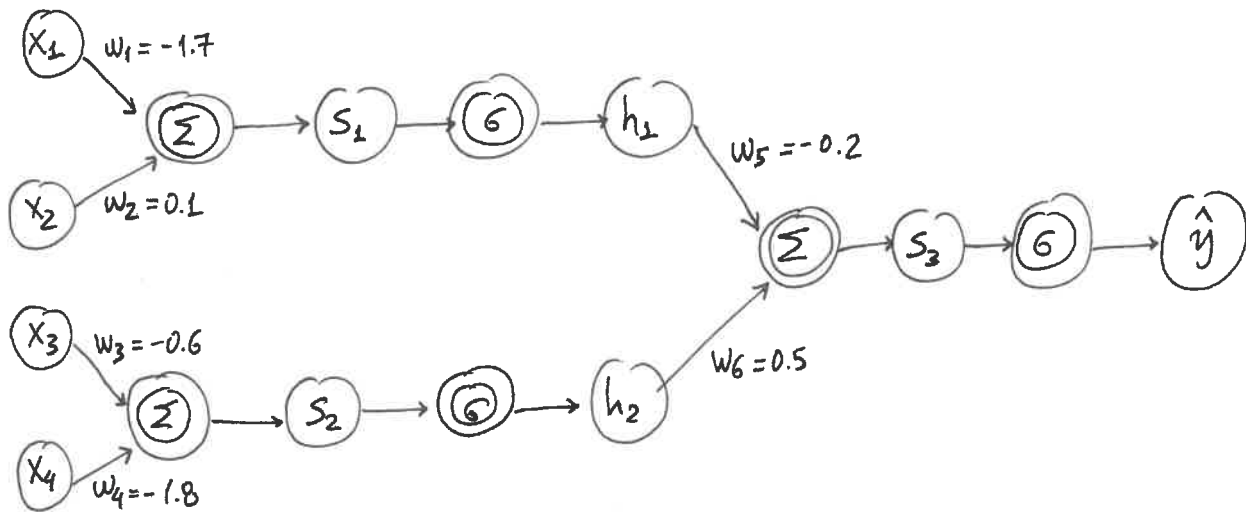
$$\text{Άρα: } f_i(\text{Word\_acc}) = \frac{0.0999 + 0}{0.0999 + 0.664 + 0 + 0} = \frac{0.0999}{0.7639} = 0.1307$$

$$f_i(\text{Task\_Duration}) = \frac{0.664}{0.7639} = 0.8692$$

$$f_i(\text{Task\_completion}) = 0$$

Συμπερασματικά:  $(\text{Task\_Duration}) > (\text{Word\_accuracy}) > (\text{Task\_completion})$ .

## Aufgabe 2.4: (MLP Backpropagation)



Activation Function:  $\sigma(x) = \frac{1}{1 + e^{-x}}$  (logistic)

Loss Function:  $L(y, \hat{y}) = \|y - \hat{y}\|_2^2$

Input Vector:  $(x_1, x_2, x_3, x_4) = (-0.5, 1.4, 0.9, -3)$

True Label:  $y = 0.5$

Forward Pass:

$$S_1 = w_1 x_1 + w_2 x_2 \Rightarrow S_1 = 0.99$$

$$h_1 = \frac{1}{1 + e^{-0.99}} \Rightarrow h_1 = 0.729$$

$$S_2 = w_3 x_3 + w_4 x_4 \Rightarrow S_2 = 4.86$$

$$h_2 = \frac{1}{1 + e^{-4.86}} \Rightarrow h_2 = 0.9923$$

$$S_3 = w_5 h_1 + w_6 h_2 \Rightarrow S_3 = 0.35$$

$$\hat{y} = \sigma(S_3) = \frac{1}{1 + e^{-0.35}} \Rightarrow \hat{y} = 0.5867$$

$$\text{Zuweis: } L(y, \hat{y}) = \|y - \hat{y}\|^2 \Rightarrow L(y, \hat{y}) = 0.007517$$

### Backward Pass:

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial s_3} \cdot \frac{\partial s_3}{\partial h_1} \cdot \frac{\partial h_1}{\partial s_1} \cdot \frac{\partial s_1}{\partial w_1} \quad (\text{chain rule})$$

Όπου :

$$\frac{\partial L}{\partial \hat{y}} = +2 \|y - \hat{y}\| = 0.086699 \times 2 = 0.1728$$

$$\frac{\partial \hat{y}}{\partial s_3} = \frac{\partial}{\partial s_3} \left[ \frac{1}{1 + e^{-s_3}} \right] = \hat{y}(1 - \hat{y}) = 0.5867 \cdot (1 - 0.5867) = 0.242483$$

$$\frac{\partial s_3}{\partial h_1} = \frac{\partial}{\partial h_1} [w_5 h_1 + w_6 h_2] = w_5 = -0.2$$

$$\frac{\partial h_1}{\partial s_1} = \frac{\partial}{\partial s_1} \left[ \frac{1}{1 + e^{-s_1}} \right] = h_1(1 - h_1) = 0.19752$$

$$\frac{\partial s_1}{\partial w_1} = \frac{\partial}{\partial w_1} [w_1 x_1 + w_2 x_2] = x_1 = -0.5$$

Συνολικά :

$$\frac{\partial L}{\partial w_1} = 2 \times 0.086699 \times 0.242483 \times (-0.2) \times (0.19752) \times (-0.5)$$

$$\Rightarrow \boxed{\frac{\partial L}{\partial w_1} = 8,3048 \times 10^{-4}}$$

και το ανανεωμένο βάρος θα είναι:

$$w_1' = w_1 + \eta \frac{\partial L}{\partial w_1}$$



## Άσκηση 2.5 : (KLT-PCA)

1) Υποθέτουμε τυχαία διανύσματα  $x \in \mathbb{C}^d$ . Ο αντίστοιχος χώρος Hilbert είναι εξορισμένος με το εσωτερικό γινόμενο :

$$\langle x, y \rangle = E\{y^H \cdot x\}$$

Θέλουμε να ορίσουμε τον PCA μετασχηματισμό με βάση την ελαχιστοποίηση ενός λάθους προβολής. Για τον σκοπό αυτό, ελάχουμε ένα πλήρες ορθοκανονικό σύνολο από  $d$ -διάστατα διανύσματα βάσης  $\{u_i\}$ , με  $i=1, 2, \dots, d$  για τα οποία :

$$\langle u_i, u_j \rangle = E\{u_i^H u_j\} = \delta_{ij} = \begin{cases} 1, & \text{if } i=j \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

Κάθε ένα από τα τυχαία διανύσματα  $x_n \in \mathbb{C}^d$ , λοιπόν, μπορεί να γραφεί ως γραμμικός συνδυασμός των διανυσμάτων βάσης :

$$x_n = \sum_{i=1}^d y_{ni} \cdot u_i \quad (2)$$

όπου οι συντελεστές  $y_{ni}$  θα είναι διαφορετικοί για κάθε σημείο  $x_n$ . Κάθε τέτοιος συντελεστής, μάλιστα, θα είναι η προβολή του σημείου  $x_n$  πάνω στο διάνυσμα βάσης  $u_i$ , δηλαδή το εσωτερικό γινόμενο τους :

$$y_{ni} = \langle x_n, u_i \rangle = E\{x_n^H u_i\}$$

Οπότε, η (2) μπορεί να γραφεί :

$$x_n = \sum_{i=1}^d \langle x_n, u_i \rangle \cdot u_i \quad (3)$$

Στόχος μας είναι να προσεγγίσουμε τα διανύσματα  $x_n$  με αναπαράστασεις  $p$ -διαστάσεων, όπου  $p < d$ , δηλαδή να λάβουμε τις προβολές τους σε έναν υποχώρο χαμηλότερης διαστατικότητας. Χωρίς βλάβη της γενικότητας, ο  $p$ -διάστατος αυτός χώρος Hilbert, μπορεί να αναπαράσται με τα  $p$  πρώτα διανύσματα βάσης. Η προσέγγιση του  $x_n$  θα είναι :

$$\tilde{x}_n = \sum_{i=1}^p z_{ni} u_i + \sum_{i=p+1}^d b_i u_i \quad (4)$$

Ορίζουμε, στη συνέχεια, τη συνάρτηση εσφαλματος λόγω μείωσης διαστατικότητας ως:

$$J = E \{ \|x - \tilde{x}\|^2 \} \stackrel{(4)}{=} E \left\{ \left\| x - \sum_{i=1}^p z_i u_i - \sum_{i=p+1}^d b_i u_i \right\|^2 \right\} \quad (5)$$

το οποίο θέλουμε να ελαχιστοποιήσουμε ως προς τα  $z_i, b_i$ . Χρησιμοποιώντας τις σχέσεις ορθοκανονικότητας (1) έχουμε:

$$\nabla_{z_i} J = 0 \Rightarrow \boxed{z_j = \langle x, u_j \rangle}, \quad j=1, 2, \dots, p$$

$$\nabla_{b_i} J = 0 \Rightarrow \boxed{b_j = \langle \bar{x}, u_j \rangle}, \quad j=p+1, \dots, d, \quad \text{όπου } \bar{x} = E\{x\} \text{ η μέση τιμή των τυχαίων διανυσμάτων } x \in \mathbb{C}^d.$$

Αντικαθιστώντας στη σχέση (5):

$$\begin{aligned} J &= E \left\{ \left\| x - \sum_{i=1}^p \langle x, u_i \rangle u_i - \sum_{i=p+1}^d \langle \bar{x}, u_i \rangle u_i \right\|^2 \right\} \stackrel{(3)}{=} \\ &= E \left\{ \left\| \sum_{i=1}^d \langle x, u_i \rangle u_i - \sum_{i=1}^p \langle x, u_i \rangle u_i - \sum_{i=p+1}^d \langle \bar{x}, u_i \rangle u_i \right\|^2 \right\} \\ &= E \left\{ \left\| \sum_{i=p+1}^d (\langle x, u_i \rangle - \langle \bar{x}, u_i \rangle) \cdot u_i \right\|^2 \right\} \quad (6) \end{aligned}$$

Λαμβάνοντας υπόψη ότι ο covariance matrix για τυχαία διανύσματα στον χώρο  $\mathbb{C}^d$

είναι:

$$S = E \{ (x - \bar{x})(x - \bar{x})^H \}$$

η σχέση (6) μπορεί να γραφεί ως:

$$J = \sum_{i=p+1}^d u_i^H S u_i \quad (7)$$

Θέλουμε να ελαχιστοποιήσουμε την παραπάνω ποσότητα ως προς  $u_i$ , υπό τους περιορισμούς ορθοκανονικότητας (1). Για τον σκοπό αυτό, ορίζουμε το Lagrangian

συναρτησιακό:

$$\tilde{J} = \sum_{i=p+1}^d u_i^H S u_i + \sum_{i=p+1}^d \lambda_i (\langle u_i, u_i \rangle - 1) \quad (8)$$

Έχοντας ολοκληρώσει τον φορμαλισμό του προβλήματος για τυχαία διανύσματα, προβαίνουμε στη λύση με Μαθηματική Επαγωγή:

□ Βάση Επαγωγής ( $p=1$ ):

Έστω, δηλαδή, ότι μεταχειριστούμε τα διανύσματα  $x \in \mathbb{C}^d$  σε 1-διάστατες αναπαραστάσεις.

Τότε, η (8) γράφεται:

$$\tilde{J}_1 = \sum_{i=2}^d u_i^H S u_i + \sum_{i=2}^d \lambda_i (-\langle u_i, u_i \rangle + 1)$$

Θέτοντας την παράγωγο ως προς  $u_i$  ίση με μηδέν:

$$S u_i - \lambda_i \cdot 2 u_i = 0 \Rightarrow S u_i = \lambda_i u_i$$

Δηλαδή, τα  $u_i$  μπορούν να είναι οποιαδήποτε από τα ιδιοδιανύσματα του μητρώου  $S$  με αντίστοιχη ιδιοτιμή  $\lambda_i$ . Αν θεωρήσουμε ότι ο πίνακας  $S$  έχει ιδιοτιμές  $\lambda_1 > \lambda_2 > \dots > \lambda_d$  με αντίστοιχα ιδιοδιανύσματα  $e_1, e_2, \dots, e_d$ , τότε για να ελαχιστοποιήσουμε τη συνάρτηση  $J$ , ορίζουμε ως  $u_i$  τα ιδιοδιανύσματα του  $S$  που αντιστοιχούν στις  $d-1$  μικρότερες ιδιοτιμές. Με άλλα λόγια, προβάλλουμε τα  $x \in \mathbb{C}^d$  στη διεύθυνση του ιδιοδιανύσματος  $e_1$  (με την μεγαλύτερη ιδιοτιμή  $\lambda_1$ ) και το ερώτημα θα είναι:

$$\tilde{J}_1 = \sum_{i=2}^d e_i^H S e_i = \sum_{i=2}^d \lambda_i$$

□ Επαγωγική Υπόθεση ( $p:=p) < d$

Υποθέτουμε ότι τα  $u_i$  είναι τα ιδιοδιανύσματα  $e_i$  του πίνακα  $S$  που αντιστοιχούν στις  $d-p$  μικρότερες ιδιοτιμές (ισοδυναμώς προβάλλουμε τα  $x$  στον υποχώρο που ορίζουν τα ιδιοδιανύσματα του  $S$  που αντιστοιχούν στις  $p$  μεγαλύτερες ιδιοτιμές). Τότε:

$$\tilde{J}_p = \sum_{i=p+1}^d e_i^H S e_i = \sum_{i=p+1}^d \lambda_i$$

□ Βήμα Επάρχως ( $p: p+1$ )  $< d$

Θέλουμε να ελαχιστοποιήσουμε το συναρτησιακό :

$$J_{p+1} = \sum_{i=p+2}^d u_i^H S u_i + \sum_{i=p+2}^d \lambda_i (1 - \langle u_i, u_i \rangle)$$

ως προς τα  $u_i$ . Παραγωγίζοντας ως προς  $u_i$  και θέτοντας με μηδέν :

$$S u_i = \lambda_i u_i$$

Δηλαδή, τα  $u_i$  είναι οποιαδήποτε από τα ιδιοδιανύσματα του  $S$  με ιδιοτιμές  $\lambda_i$ . Αν αντικαταστήσουμε στην  $J$ , τότε προκύπτει :

$$J_{p+1} = \sum_{i=p+2}^d e_i^H S e_i = \sum_{i=p+2}^d \lambda_i \quad (\text{για κάποια ιδιοδιανύσματα } e_i)$$

Οπότε, για να ελαχιστοποιείται η  $J$  επιλέγουμε ως  $u_i$  τα ιδιοδιανύσματα  $e_i$  του πίνακα  $S$  που αντιστοιχούν στις  $d - (p+1) = d - p - 1$  μικρότερες ιδιοτιμές. Αν

θεωρήσουμε, λοιπόν, ότι  $\lambda_1 > \lambda_2 > \dots > \lambda_d$ , θα ισχύει :

$$J_{p+1} = \sum_{i=p+2}^d \lambda_i = \left( \sum_{i=p+1}^d \lambda_i \right) - \lambda_{p+1}$$

$$\Rightarrow \boxed{J_{p+1} = J_p - \lambda_{p+1}}$$

2) Το δυναμνησιακό προς ελαχιστοποίηση που ορίσαμε στο προηγούμενο ερώτημα ήταν :

$$\tilde{J} = \sum_{i=p+1}^d u_i^H S u_i + \sum_{i=p+1}^d \lambda_i (1 - \langle u_i, u_i \rangle) \quad (2.1)$$

με παραμέτρους τα  $u_i$ . Αν ορίσουμε τον πίνακα  $U$ , διαστάσεων  $d \times (d-p)$ , του οποίου οι στήλες είναι τα  $u_i$  και τον πίνακα  $H$  ως τον διαγώνιο πίνακα με τους Lagrangian συντελεστές  $\lambda_i$ , και επίσης μετονομάσουμε τον πίνακα συνδιασποράς  $S$  σε  $R$ , τότε η παραπάνω εξίσωση δράφεται :

$$\tilde{J} = \text{trace} \{ U^H R U \} + \text{trace} \{ H (I - U^H U) \} \quad (2.2)$$

Δείξαμε ότι παραγωγίζοντας το  $\tilde{J}$  ως προς  $u_i$  και θέτοντας με μηδέν, προκύπτουν εξισώσεις  $R u_i = \lambda_i u_i$ . Βάζοντας τες όλες μαζί ως μια εξίσωση, προκύπτει :

$$\boxed{R U = U H} \quad (2.3)$$

η οποία είδαμε ότι δίνει ως λύση, ο πίνακας  $U$  να έχει ως στήλες ιδιοδιανύσματα του μητρώου συνδιασποράς  $R$  και το μητρώο  $H$  να είναι διαγώνιος πίνακας με τις αντίστοιχες ιδιοτιμές.

Για να αποκτήσουμε μια γενική λύση, θεωρούμε ότι ο  $H$  μπορεί να είναι ένας συμμετρικός Hermitian πίνακας. Τότε, σύμφωνα με το Φασματικό Θεώρημα [2] μπορεί να παραχοντοποιηθεί ως :

$$H = F D F^H$$

όπου  $F$  unitary matrix και  $D$  diagonal matrix.

Λαμβάνοντας υπόψη ότι τα ιδιοδιανύσματα του  $H$  είναι orthogonal, προβαίνουμε στη φασματική του διάσπαση :

$$H = (f_1, \dots, f_n) \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \lambda_n \end{pmatrix} \begin{pmatrix} f_1^H \\ \vdots \\ f_n^H \end{pmatrix} = (\lambda_1 f_1, \dots, \lambda_n f_n) \begin{pmatrix} f_1^H \\ \vdots \\ f_n^H \end{pmatrix} = \sum_{i=1}^n \lambda_i f_i f_i^H$$

με  $n = d-p$ . Σημειώνεται ότι τα  $f_i f_i^H$  είναι rank-one πίνακες.

Αντικαθιστώντας, λοιπόν, στην (2.2) έχουμε:

$$\tilde{J} = \text{trace}\{U^H R U\} + \text{trace}\left\{\sum_{i=1}^{d-p} \lambda_i f_i f_i^H \cdot (I - U^H U)\right\}$$

Παραγωγίζοντας ως προς ένα εκ των  $u_i$  και θέτοντας με μηδέν έχουμε:

$$\cancel{Z}^H R u_i - \lambda_i f_i f_i^H \cancel{Z} u_i = 0 \Rightarrow \boxed{R u_i = (f_i f_i^H) \cdot \lambda_i u_i} \quad (2.4)$$

πρόκειται, δηλαδή, για ένα πρόβλημα γενικευμένων ιδιοδιανυσμάτων. Όμως, <sup>και τα  $f_i$  orthogonal</sup> ο πίνακας  $f_i f_i^H$  είναι rank-one, οπότε η λύση της (2.4) θα είναι ίδια με τη λύση της εξίσωσης  $R u_i = \lambda_i u_i$  (απλό ιδιοδιάνυσμα) του μητρώου  $R$ .

Συνεπώς, η γενική λύση είναι ίδια με την ειδική, σύμφωνα με την οποία τα  $u_i$  είναι ιδιοδιανύσματα του μητρώου  $R$ .

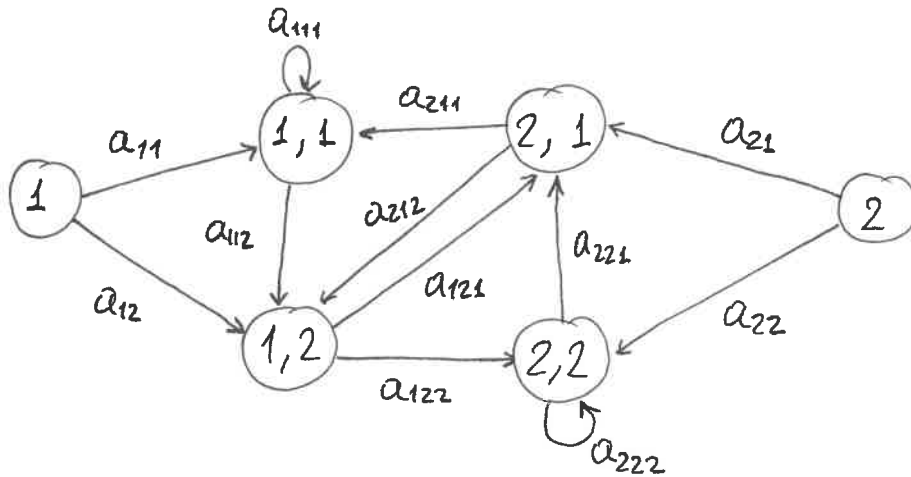


## Άσκηση 2.6 : (Graphical Models)

Δίνεται Markov model με μνήμη 2 (MM-2) με δύο καταστάσεις 1 και 2. Ισχύει η σχέση ανεξαρτησίας:

$$P(q_t | q_{t-1}, q_{t-2}, q_{t-3}, \dots) = P(q_t | q_{t-1}, q_{t-2})$$

1) Bayesian Net:



2) Παρατηρούμε την ακόλουθη σειρά από καταστάσεις:

Seq = (1, 1, 1, 2, 1, 2, 2, 2, 2, 1, 2, 2, 2, 2, 1, 1, 1, 1, 2)

Δίνεται ότι:  $P(q_0=1) = \pi_1 = 0.5$

$P(q_0=2) = \pi_2 = 0.5$

Με Maximum Likelihood έχουμε:

(δίνεται)

$$\begin{aligned} \cdot a_{11} &= P(q_t=1 | q_{t-1}=1) = \frac{1}{4} = 0.25 \\ \cdot a_{12} &= P(q_t=2 | q_{t-1}=1) = \frac{3}{4} = 0.75 \\ \cdot a_{21} &= P(q_t=1 | q_{t-1}=2) = \frac{1}{4} = 0.25 \\ \cdot a_{22} &= P(q_t=2 | q_{t-1}=2) = \frac{3}{4} = 0.75 \end{aligned}$$

$$\cdot a_{111} = P(q_t=1 | q_{t-1}=1, q_{t-2}=1) = \frac{\#(11 \rightarrow 1)}{\#(11 \rightarrow *)} = \frac{3}{5} = 0.6$$

$$\cdot a_{112} = P(q_t=2 | q_{t-1}=1, q_{t-2}=1) = \frac{\#(11 \rightarrow 2)}{\#(11 \rightarrow *)} = \frac{2}{5} = 0.4$$

$$a_{121} = P(q_t=1 | q_{t-1}=2, q_{t-2}=1) = \frac{\#(12 \rightarrow 1)}{\#(12 \rightarrow *)} = \frac{1}{2} = 0.5$$

$$a_{122} = P(q_t=2 | q_{t-1}=2, q_{t-2}=1) = \frac{\#(12 \rightarrow 2)}{\#(12 \rightarrow *)} = 0.5$$

$$a_{211} = P(q_t=1 | q_{t-1}=1, q_{t-2}=2) = \frac{\#(21 \rightarrow 1)}{\#(21 \rightarrow *)} = \frac{1}{3} = 0.333...$$

$$a_{212} = P(q_t=2 | q_{t-1}=1, q_{t-2}=2) = \frac{\#(21 \rightarrow 2)}{\#(21 \rightarrow *)} = \frac{2}{3} = 0.666...$$

$$a_{221} = P(q_t=1 | q_{t-1}=2, q_{t-2}=2) = \frac{\#(22 \rightarrow 1)}{\#(22 \rightarrow *)} = \frac{2}{6} = \frac{1}{3} = 0.333...$$

$$a_{222} = P(q_t=2 | q_{t-1}=2, q_{t-2}=2) = \frac{\#(22 \rightarrow 2)}{\#(22 \rightarrow *)} = \frac{4}{6} = \frac{2}{3} = 0.666...$$

$$\begin{aligned} 3) P(q_1=1, q_2=1, \dots, q_{10}=1, q_{11}=2 | q_0=1) &= \frac{P(\overbrace{11 \dots 1}^{11 \text{ times}} 2)}{P(q_0=1)} = \frac{\cancel{\pi_1} \cdot a_{11} \cdot a_{11}^9 \cdot a_{112}}{\cancel{\pi_1}} \\ &= \frac{1}{4} \cdot (0.6)^9 \cdot 0.4 = \underline{0.004031} \times \frac{1}{4} \simeq 0.00100775 \end{aligned}$$

4) Θεωρούμε μοντέλο Markov με μνήμη 1 (MM-1):

$$\begin{aligned} (a') \quad a_{11} &= \frac{\#(1 \rightarrow 1)}{\#(1 \rightarrow *)} = \frac{2+3}{9} = \frac{5}{9} & a_{21} &= \frac{\#(2 \rightarrow 1)}{\#(2 \rightarrow *)} = \frac{3}{9} = \frac{1}{3} \\ a_{12} &= \frac{\#(1 \rightarrow 2)}{\#(1 \rightarrow *)} = \frac{4}{9} & a_{22} &= \frac{\#(2 \rightarrow 2)}{\#(2 \rightarrow *)} = \frac{6}{9} = \frac{2}{3} \end{aligned}$$

Συνεπώς, ο πίνακας μετάβασης είναι  $A = \begin{bmatrix} 5/9 & 4/9 \\ 1/3 & 2/3 \end{bmatrix}$ .

$$(b) P(q_1=1, \dots, q_{10}=1, q_{11}=2 | q_0=1) = \frac{\cancel{\pi_1} \cdot a_{11}^{10} \cdot (1-a_{11})}{\cancel{\pi_1}} = \left(\frac{5}{9}\right)^{10} \cdot \frac{4}{9}$$

$$= \underline{0.0012448} > 0.00100775 \quad \text{λόγω της υποθέσεως στο MM-2} \\ \text{ότι } a_{11} = \frac{1}{4} = 0.25.$$

## Άσκηση 2.3: (Linear Discriminant Analysis)

Δίνονται τα μητρώα:

$$S_W = \sum_{i=1}^{|Classes|} E_{x|x \in \omega_i} [(\vec{x} - \vec{\mu})(\vec{x} - \vec{\mu})^T]$$

$$S_B = \sum_{i=1}^{|Classes|} P(\omega_i) (\vec{\mu}_i - \vec{\mu})(\vec{\mu}_i - \vec{\mu})^T$$

όπου  $\omega_i$ :  $i$ -οστή κλάση με μέση τιμή  $\vec{\mu}_i$

$|Classes|$ : το πλήθος των κλάσεων

$\vec{\mu}$ : η μέση τιμή όλων των δειγμάτων

α) Έστω η περίπτωση όπου  $|Classes| = 2$ . Τότε:

$$S_B = P(\omega_1) \cdot (\vec{\mu}_1 - \vec{\mu})(\vec{\mu}_1 - \vec{\mu})^T + P(\omega_2) (\vec{\mu}_2 - \vec{\mu})(\vec{\mu}_2 - \vec{\mu})^T$$

Λαμβάνοντας υπόψη ότι  $\vec{\mu} = P(\omega_1) \cdot \vec{\mu}_1 + P(\omega_2) \vec{\mu}_2$  και  $P(\omega_1) = 1 - P(\omega_2)$ :

$$\begin{aligned} S_B &= P(\omega_1) \cdot (\vec{\mu}_1 - P(\omega_1)\vec{\mu}_1 - P(\omega_2)\vec{\mu}_2) \cdot (\vec{\mu}_1 - P(\omega_1)\vec{\mu}_1 - P(\omega_2)\vec{\mu}_2)^T + \\ &\quad + P(\omega_2) \cdot (\vec{\mu}_2 - P(\omega_1)\vec{\mu}_1 - P(\omega_2)\vec{\mu}_2) \cdot (\vec{\mu}_2 - P(\omega_1)\vec{\mu}_1 - P(\omega_2)\vec{\mu}_2)^T \\ &= P(\omega_1) \cdot [(1 - P(\omega_1))\vec{\mu}_1 - P(\omega_2)\vec{\mu}_2] [ \vec{\mu}_1 \cdot (1 - P(\omega_1)) - P(\omega_2)\vec{\mu}_2 ]^T + \\ &\quad + P(\omega_2) \cdot [-P(\omega_1)\vec{\mu}_1 + (1 - P(\omega_2))\vec{\mu}_2] \cdot [(1 - P(\omega_2))\vec{\mu}_2 - P(\omega_1)\vec{\mu}_1]^T \\ &= P(\omega_1) \cdot P(\omega_2)^2 \cdot (\vec{\mu}_1 - \vec{\mu}_2)(\vec{\mu}_1 - \vec{\mu}_2)^T + P(\omega_2) \overset{2}{P(\omega_1)} \cdot (\vec{\mu}_2 - \vec{\mu}_1) \cdot (\vec{\mu}_2 - \vec{\mu}_1)^T \\ &= P(\omega_1) \cdot P(\omega_2) \cdot (\vec{\mu}_1 - \vec{\mu}_2)(\vec{\mu}_1 - \vec{\mu}_2)^T \cdot (\cancel{P(\omega_1)} + P(\omega_2))^T \\ &= P(\omega_1)P(\omega_2) \cdot (\vec{\mu}_1 - \vec{\mu}_2)(\vec{\mu}_1 - \vec{\mu}_2)^T \end{aligned}$$

β) Το πρόβλημα που καλούμαστε να λύσουμε είναι το εξής :

$$\bar{S}_w^{-1} S_B \vec{V} = \lambda \vec{V} , \vec{V} \neq \vec{0}$$

Όπως :

$$\begin{aligned} \bar{S}_w^{-1} S_B \vec{V} &= \bar{S}_w^{-1} P(w_1) P(w_2) (\vec{\mu}_1 - \vec{\mu}_2) \cdot \underbrace{(\vec{\mu}_1 - \vec{\mu}_2)^T \cdot \vec{V}}_{\text{Βαθμωτό}} \\ &= \underbrace{P(w_1) P(w_2) (\vec{\mu}_1 - \vec{\mu}_2)^T \vec{V}}_{\substack{\lambda \in \mathbb{C} \\ \text{Βαθμωτό}}} \cdot \bar{S}_w^{-1} (\vec{\mu}_1 - \vec{\mu}_2) \end{aligned}$$

Η ποσότητα  $\bar{S}_w^{-1} S_B \vec{V}$  έχει την ίδια διεύθυνση με το διάνυσμα  $e = \bar{S}_w^{-1} (\vec{\mu}_1 - \vec{\mu}_2)$ , για κάθε  $\vec{V} \neq \vec{0}$ . Το ιδιοδιάνυσμα, λοιπόν, του μητρώου  $\bar{S}_w^{-1} S_B$  θα είναι το:

$$\boxed{e = \bar{S}_w^{-1} (\vec{\mu}_1 - \vec{\mu}_2)}$$

Ενώ η αντίστοιχη ιδιοτιμή θα είναι :

$$\begin{aligned} \lambda &= P(w_1) P(w_2) (\vec{\mu}_1 - \vec{\mu}_2)^T \cdot \vec{e} \\ \Rightarrow \quad &\boxed{\lambda = P(w_1) P(w_2) (\vec{\mu}_1 - \vec{\mu}_2)^T \cdot \bar{S}_w^{-1} (\vec{\mu}_1 - \vec{\mu}_2)} \end{aligned}$$

## Άσκηση 2.8 : (ICA)

α) Έστω τυχαιές μεταβλητές  $x, y$  μη μηδενικές μέσες τιμές. Τότε :

$$\begin{aligned} \cdot E[(x+y)^4] &= E[x^4 + 4x^3y + 6x^2y^2 + 4xy^3 + y^4] \\ &= E[x^4] + 4E[x^3y] + 6E[x^2y^2] + 4E[xy^3] + E[y^4] \end{aligned}$$

Λόγω ανεξαρτησίας, εν γένει ισχύει  $E[x \cdot y] = E[x] \cdot E[y]$ . Οπότε :

$$\begin{aligned} E[(x+y)^4] &= E[x^4] + 4E[x^3] \cdot \overset{0}{E[y]} + 6E[x^2y^2] + 4\overset{0}{E[x]} \cdot E[y^3] + E[y^4] \\ &= E[x^4] + 6E[x^2] \cdot E[y^2] + E[y^4] \end{aligned}$$

$$\cdot E[(x+y)^2] = E[x^2 + 2xy + y^2] = E[x^2] + 2E[xy] + E[y^2]$$

$$\overset{x, y \text{ ανεξ.}}{=} E[x^2] + 2\overset{0}{E[x] \cdot E[y]} + E[y^2] = E[x^2] + E[y^2].$$

Επομένως :

$$\begin{aligned} \text{kurt}(x+y) &= E[(x+y)^4] - 3(E[x+y]^2)^2 \\ &= E[x^4] + 6E[x^2]E[y^2] + E[y^4] - 3 \cdot (E[x^2] + E[y^2] + 2E[x]E[y])^2 \\ &= \underbrace{E[x^4] - 3(E[x^2])^2}_{\text{kurt}(x)} + \underbrace{E[y^4] - 3(E[y^2])^2}_{\text{kurt}(y)} \end{aligned}$$

$$\Rightarrow \boxed{\text{kurt}(x+y) = \text{kurt}(x) + \text{kurt}(y)}$$

B1) Θεωρούμε  $N$  στατιστικώς ανεξάρτητες Τ.Μ.  $s_i$ , με μηδενικές μέσες τιμές, μοναδιαίες διασπορές και κυρτώσεις με τιμές  $a_i$ , όπου  $-a \leq a_i \leq a$ , για κάποια άγνωστη αλλά σταθερή τιμή  $a$ . Δηλαδή:

$$E[s_i] = 0, \quad \text{Var}(s_i) = 1, \quad \text{kurt}(s_i) = a_i, \quad \text{για } i = 1, 2, \dots, N.$$

Αναμειγνύουμε τις παραπάνω Τ.Μ. μέσω σταθερών βαρών  $w_i$  ως εξής:

$$x := \sum_{i=1}^N w_i s_i$$

Τότε:

$$E[x] = E\left[\sum_{i=1}^N w_i s_i\right] = \sum_{i=1}^N w_i \overset{0}{E[s_i]} = 0$$

Έπειτα:

$$\text{Var}(x) = E[x^2] - \overset{0}{(E[x])^2} = E\left[\left(\sum_{i=1}^N w_i s_i\right)^2\right]$$

$$= E\left[\sum_{i=1}^N w_i^2 s_i^2 + \sum_{i=1}^N \sum_{\substack{j=1, \\ j \neq i}}^N w_i w_j s_i s_j\right]$$

$$= \sum_{i=1}^N w_i^2 E[s_i^2] + \sum_{i=1}^N \sum_{\substack{j=1, \\ j \neq i}}^N w_i w_j \overset{0}{E[s_i] E[s_j]}, \quad \text{διότι } s_i, s_j \text{ ανεξάρτητα} \\ \text{εφόσον } i \neq j$$

$$= \sum_{i=1}^N w_i^2 E[s_i^2]$$

$$\text{Όπως } \text{Var}(s_i) = E[s_i^2] - \overset{0}{(E[s_i])^2} \Rightarrow E[s_i^2] = 1.$$

Άρα:  $\text{Var}(x) = \sum_{i=1}^N w_i^2$ . Για να έχει, λοιπόν, η τ.μ.  $x$  μοναδιαία διασπορά

πρέπει να ισχύει:

$$\boxed{\sum_{i=1}^N w_i^2 = 1}$$



B2) Θεωρούμε ένα ομοιόμορφα σταθμισμένο μείγμα, δηλαδή  $w_i = w_j \forall i, j, N$  τυχαίων μεταβλητών. Θεωρούμε, επίσης, ότι το μείγμα έχει μοναδιαία διασπορά, οπότε από την απαίτηση του (β1):

$$\sum_{i=1}^N w_i^2 = 1 \xRightarrow{\forall i,j} \sum_{i=1}^N w^2 = 1 \Rightarrow N \cdot w^2 = 1 \Rightarrow \boxed{w^2 = \frac{1}{N}}$$

Επιπλέον, ισχύει:

$$\begin{aligned} \text{kurt}(ax) &= E[(ax)^4] - 3(E[(ax)^2])^2 = a^4 \cdot [E[x^4] - 3(E[x^2])^2] \\ &= a^4 \cdot \text{kurt}(x), \quad \forall a \in \mathbb{R}. \end{aligned}$$

Οπότε, σύμφωνα με το παραπάνω και την ιδιότητα του ερωτήματος (α) έχουμε:

$$\text{kurt}(x) = \text{kurt}\left(\sum_{i=1}^N w s_i\right) = w^4 \cdot \text{kurt}\left(\sum_{i=1}^N s_i\right) \stackrel{(a)}{=} w^4 \sum_{i=1}^N \text{kurt}(s_i)$$

$$\begin{aligned} w^2 &= \frac{1}{N} \\ &= \frac{1}{N^2} \sum_{i=1}^N \text{kurt}(s_i) = \frac{1}{N^2} \sum_{i=1}^N a_i \end{aligned}$$

Οι επιμέρους κυρτώσεις των τ.μ.  $s_i$ , όμως, δηλαδή οι τιμές  $a_i$  είναι φραγμένες εντός του διαστήματος  $[-a, a]$ , όπου  $a$  μια άγνωστη αλλά σταθερή τιμή. Ος

Εκ τούτου:

$$\boxed{\lim_{N \rightarrow +\infty} \frac{1}{N^2} \sum_{i=1}^N a_i = 0}, \quad \text{καθώς} \quad \lim_{N \rightarrow +\infty} \sum_{i=1}^N a_i \in \mathbb{R}.$$

$$\text{Δηλαδή} \quad \text{kurt}(x) \xrightarrow{N \rightarrow +\infty} 0.$$

## Άσκηση 2.9: (Logistic Regression)

Θεωρούμε το πρόβλημα logistic regression για ένα σύνολο δεδομένων  $\{\varphi_n, t_n\}$ , όπου  $t_n \in \{0, 1\}$  και  $\varphi_n = \varphi(x_n)$  είναι οι κατηγορίες και οι συναρτήσεις βάσης, αντίστοιχα, για δείγματα  $n=1, 2, \dots, N$ .

Η έξοδος του LR μοντέλου στο διάνυσμα εισόδου  $x_n$  δηλώνει την a-posterior πιθανότητα το δείγμα να ανήκει στην κατηγορία  $C_1$ , δηλαδή:

$$y_n = p(C_1 | \varphi_n) = \sigma(w^T \varphi_n)$$

$$\text{και } p(C_2 | \varphi_n) = 1 - p(C_1 | \varphi_n) = 1 - \sigma(w^T \varphi_n)$$

όπου  $w$ : το διάνυσμα βαρών και  $\sigma(x) = \frac{1}{1 + e^{-x}}$  η σιγμοειδής συνάρτηση.

Τότε, η συνάρτηση πιθανοφάνειας γράφεται:

$$p(t|w) = \prod_{n=1}^N y_n^{t_n} \cdot \{1 - y_n\}^{1-t_n}$$

από την οποία, προκύπτει το log-likelihood ή cross-entropy error:

$$E(w) = -\ln p(t|w) = -\sum_{n=1}^N \{t_n \ln y_n + (1-t_n) \ln(1-y_n)\}$$

α) Σε πρώτο στάδιο, προσδιορίζουμε την παράγωγο του cross-entropy εσφαλματος ως προς τα βάρη  $w$ :

$$\begin{aligned} \nabla_w E(w) &= -\nabla_w \sum_{n=1}^N \{t_n \ln y_n + (1-t_n) \ln(1-y_n)\} \\ &= -\sum_{n=1}^N \frac{\partial}{\partial y_n} \{t_n \ln y_n + (1-t_n) \ln(1-y_n)\} \cdot \frac{\partial y_n}{\partial (w^T \varphi_n)} \cdot \frac{\partial (w^T \varphi_n)}{\partial w} \quad (\text{chain rule}) \\ &= -\sum_{n=1}^N \left( \frac{t_n}{y_n} - \frac{1-t_n}{1-y_n} \right) \cdot y_n(1-y_n) \cdot \varphi_n \quad \left( \text{δεδομένου } \frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1-\sigma(x)) \right) \end{aligned}$$

$$\Rightarrow \nabla_w E(w) = \sum_{n=1}^N (y_n - t_n) \cdot \varphi_n$$

Δεδομένου ότι τα target values  $t_i$  παίρνουν τιμές  $\{0, 1\}$ , το likelihood θα λαμβάνει τη μέγιστη τιμή του για :

$$y = 0.5 \Rightarrow \sigma(w^T \varphi) = 0.5 \Rightarrow \frac{1}{1 + e^{-w^T \varphi}} = 0.5 \Rightarrow 0/5 = 0.5 e^{-w^T \varphi} \\ \Rightarrow e^{-w^T \varphi} = 1 \Rightarrow \boxed{w^T \varphi = 0}$$

Με άλλα λόγια, η διαχωριστική επιφάνεια θα είναι η  $w^T \varphi(x) = 0$ , έτσι τα δείγματα  $x_i$  που ανήκουν στην κατηγορία  $C_1$  θα έχουν  $w^T \varphi(x_i) > 0$ , ενώ τα δείγματα  $x_m$  που ανήκουν στην κατηγορία  $C_2$  θα έχουν  $w^T \varphi(x_m) < 0$ .

Ταυτόχρονα, αυξάνοντας το μέτρο του διανύσματος βαρών  $w$ , δηλαδή για  $|w| \rightarrow +\infty$ , θα έχουμε:  $w^T \varphi(x_i) \rightarrow +\infty$ ,  $w^T \varphi(x_m) \rightarrow -\infty$

$$\cdot \lim_{|w| \rightarrow \infty} P(C_1 | \varphi(x_i)) = \lim_{|w| \rightarrow \infty} \sigma(w^T \varphi(x_i)) = \lim_{|w| \rightarrow \infty} \frac{1}{1 + e^{-w^T \varphi(x_i)}} = \frac{1}{1 + 0} = 1.$$

$$\cdot \lim_{|w| \rightarrow \infty} P(C_2 | \varphi(x_m)) = \lim_{|w| \rightarrow \infty} [1 - P(C_1 | \varphi(x_m))] = \lim_{|w| \rightarrow \infty} \left[ 1 - \frac{1}{1 + e^{-w^T \varphi(x_m)}} \right] = 1.$$

Επομένως, για  $|w| \rightarrow +\infty$ , καθένας από τους όρους  $y_n = P(C_1 | \varphi(x_n))$  και  $1 - y_n = P(C_2 | \varphi(x_n))$  θα λάβει τη μέγιστη τιμή του, δηλαδή 1, οδηγώντας στη μεγιστοποίηση του likelihood  $P(t|w) = \prod_{n=1}^N y_n^{t_n} \cdot (1 - y_n)^{1-t_n}$ .

Η παραπάνω προσέγγιση, ωστόσο, θα επιφέρει τεράστιο overfitting για γραμμικά διαχωρίσιμα datasets, καθώς η sigmoid function θα γίνει απείρως steep στο feature space, αντιστοιχώς με την Heaviside step function. [3].

β) Η Hessian μήτρα για το logistic regression δίνεται από τη σχέση:

$$H = [\Phi^T R \Phi]$$

όπου  $\Phi$  είναι ο πίνακας των χαρακτηριστικών και  $R$  είναι ένας διαγώνιος πίνακας με στοιχεία  $y_n(1-y_n)$ .

Για να είναι ο πίνακας  $H$  θετικά ορισμένος, πρέπει  $\forall u \in \mathbb{R}^n / \{0\}$  να ισχύει:

$$u^T H u > 0$$

Οστόσο, για τον διαγώνιο πίνακα  $R$  έχουμε:

$$\det(R) = \text{tr}(R) = \sum_{n=1}^N y_n(1-y_n)$$

$$\text{ενώ } 0 < y_n < 1 \Rightarrow 0 < 1-y_n < 1 \text{ οπότε και } 0 < (1-y_n)y_n < 1$$

Συνεπώς:

$$\begin{aligned} u^T H u &= u^T \Phi^T R \Phi u = u^T \sum_{n=1}^N \Phi_n y_n(1-y_n) \Phi_n^T \cdot u = \\ &= \sum_{n=1}^N y_n(1-y_n) \cdot (\Phi_n^T \cdot u)^T (\Phi_n^T \cdot u) = \sum_{n=1}^N y_n(1-y_n) \cdot b_n^2, \text{ όπου } b_n = \Phi_n^T \cdot u \\ &\text{ με } u \in \mathbb{R}^n / \{0\}. \end{aligned}$$

Αν υποθέσουμε ότι υπάρχει τουλάχιστον ένα μη-μηδενικό χαρακτηριστικό  $\Phi_n \neq 0$ ,

ώστε  $b_n \neq 0$ , τότε λόγω του ότι  $0 < y_n(1-y_n) < 1$ , τελικώς:

$$u^T H u = \sum_{n=1}^N y_n(1-y_n) \cdot b_n^2 > 0, \text{ για κάθε } u \in \mathbb{R}^n / \{0\}$$

και το μήτριο  $H$  θα είναι θετικά ορισμένο.

$$(H = \nabla \nabla E)$$

Έχοντας, λοιπόν,  $H > 0$ , εύκολα προκύπτει ότι η  $L$  παράγωγος της συνάρτησης βελάματος  $L$  ως προς τα βάρη  $w$  είναι παντού θετική. Ως εκ τούτου, η

$L$  παράγωγος της  $L$  θα είναι γνησίως αύξουσα και θα έχει το πλὺν ένα σημείο μηδενισμού, στο οποίο θα αντιστοιχεί το μοναδικό ελάχιστο της  $L$ . Ισοδυναμώς,

η συνάρτηση  $L$  θα είναι κυρτή.

## Άσκηση 2.9(γ): Iterative Reweighted Least Squares

Σε πρώτο στάδιο, θα υλοποιήσουμε τον αλγόριθμο iterative reweighted least squares (IRWLS) για multiclass logistic regression.

```
In [ ]: from numpy as np
import matplotlib.pyplot as plt
from sklearn.preprocessing import OneHotEncoder
from sklearn.metrics import accuracy_score
```

Αρχικά, διαβάσουμε τα δεδομένα εισόδου που αποτελούνται από διανύσματα δύο διαστάσεων,  $x_1$  και  $x_2$ , τα οποία κατηγοριοποιούνται σε τρεις κλάσεις: 0, 1 και 2. Στα διανύσματα εισόδου  $X$  προσυζάνουμε την διάσταση κατά ένα, ώστε  $x[0] = 1$ , προκειμένου να διευκολυνθούμε στις πράξεις με το bias των decision surfaces.

Επιπλέον, ορίζουμε:

- $N$ : τον αριθμό των δειγμάτων εισόδου. Στην προκειμένη περίπτωση, έχουμε 120 δείγματα.
- $M$ : τον αριθμό των χαρακτηριστικών κάθε δείγματος. Στην προκειμένη περίπτωση, έχουμε 3 χαρακτηριστικά, με το πρώτο να είναι πάντα μονάδα για το bias.  $-K$ : τον αριθμό των κλάσεων του προβλήματος, στην προκειμένη περίπτωση 3.

```
In [ ]: # Read Data
f = open('MLR.data','r')
lines = f.readlines()
x1 = []
x2 = []
y = []
for x in lines:
    x1.append(float(x.split(' ')[0]))
    x2.append(float(x.split(' ')[1]))
    y.append(int(x.split(' ')[2][0]))
f.close()
```

```
In [ ]: X = np.column_stack((np.ones((len(x1)),x1,x2)) # concatenate 1 for bias multiplication
y = np.array(y)
print(X.shape)
N, M = X.shape
K = len(np.unique(y)) # 3 classes
```

(120, 3)

Έπειτα, για να ορίσουμε τις a-posteriori πιθανότητες του Multiclass Logistic Regression Classifier, χρησιμοποιούμε τη συνάρτηση softmax:

$$p(C_k|\varphi_n) = y_{nk}(\varphi_n) = \frac{exp^{w_k^T \varphi_n}}{\sum_{j=1}^K exp^{w_j^T \varphi_n}}$$

όπου  $\varphi_n$ : η συνάρτηση βάσης πάνω στο δείγμα  $X_n$ .

```
In [ ]: def softmax(x, w, Class):
    numerator = np.exp(np.dot(w[Class,:], x))
    denominator = 0.
    for i in range(np.shape(w)[0]):
        denominator += np.exp(np.dot(w[i,:], x))
    return numerator/denominator
```

```
In [ ]: def softmax_all(X, w, K = 3):
    N = X.shape[0]
    Y = np.zeros((N,K)) # samples x classes = N x K
    for j in range(K):
        for n in range(N):
            Y[n,j] = softmax(X[n], w, Class = j)
    return Y
```

```
In [ ]: w = np.random.rand(K, M) # weight vectors KxM - one for each class
```

Έπειτα, χρησιμοποιούμε ένα 1-of-K coding scheme, με βάση το οποίο το target vector  $t_n$  για ένα διάνυσμα χαρακτηριστικών  $\varphi_n$  το οποίο ανήκει στην κλάση  $C_k$  είναι ένα δυαδικό διάνυσμα με μηδέν όλα τα στοιχεία του, εκτός του k-οστού στοιχείου που είναι μονάδα.

```
In [ ]: # One-Hot Encode Labels
enc = OneHotEncoder()
enc.fit(np.array(y).reshape((-1), 1))
T = enc.transform(np.array(y).reshape((-1), 1)).toarray() # T: y one-hot encoded N x K
```

Ακολούθως, ορίζουμε τη συνάρτηση σφάλματος που λειτουργεί ως μετρική της επίδοσης του Logistic Regression Classifier. Η εν λόγω συνάρτηση  $E(w)$ , γνωστή και ως cross-entropy, είναι στην ουσία των log-likelihood των δειγμάτων εισόδου και ορίζεται ως:

$$E(w_1, \dots, w_K) = -\ln p(T|w_1, \dots, w_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

```
In [ ]: # Log-Likelihood Error Function
def log_likelihood(x, y_enc, w, classes = 3):
    E = 0
    for n in range(x.shape[0]):
        for k in range(classes):
            ynk = softmax(x[n], w, Class = k)
            E -= y_enc[n,k]*np.log(ynk)
    return E
```

```
In [ ]: w = np.random.rand(K, M)
print('Random Weights: log-Likelihood E = ', log_likelihood(X, T, w))
```

Random Weights: Log-Likelihood E = 127.04051855009493

Για την εκτέλεση του IRLS αλγόριθμου, ωστόσο, χρειάζαστε την **Hessian μήτρα** της συνάρτησης σφάλματος  $E$ . Η μήτρα αυτή είναι διαστάσεων  $KM * KM$  και αποτελείται από επαναλαμβανόμενα μπλοκ διαστάσεων  $M * M$ , με το μπλοκ j,k να δίνεται ως:

$$\nabla_{w_k} \nabla_{w_j} E(w) = \sum_{n=1}^N y_{nk} (I_{kj} - y_{nj}) \varphi_n \varphi_n^T = \Phi^T R \Phi$$

$R$ : διαγώνιος πίνακας με στοιχεία  $y_{nk} (I_{kj} - y_{nj}) \setminus I$ : ο μοναδιαίος πίνακας

```
In [ ]: def Hessian(X, w, K = 3):
    N, M = X.shape
    H = np.zeros((K*M,K*M))
    I = np.eye(K)
    Y = softmax_all(X, w)
    for k in range(K):
        for j in range(K):
            # Diagonal R
            R = np.zeros((N,N))
            for n in range(N):
                R[n][n] = Y[n][k] * (I[k,j] - Y[n][j])

            H[k*M:k*M+M, j*M:j*M+M] = np.dot( np.dot(X.T, R), X )
    return H # KM x KM
```

Ομοίως, προσδιορίζουμε και την πρώτη παράγωγο της συνάρτησης  $E$  ως:

$$\nabla_{w_j} E(w) = \sum_{n=1}^N (y_{nj} - t_{nj}) \varphi_n$$

```
In [ ]: def GradientE(X, w, K = 3):
    N, M = X.shape
    grad_E = np.zeros((K,M)) # K x M
    for j in range(K):
        for n in range(N):
            ynj = softmax(X[n],w, Class = j)
            grad_E[j,:] += (ynj - T[n,j])*X[n]
```

return grad\_E

Στο σημείο αυτό, τρέχουμε τον αλγόριθμο IRLS για Multiclass LR. Ο αλγόριθμος τρέχει επαναληπτικά (εδώ για 20 επαναλήψεις), και τα βάρη ανανεώνονται ως εξής:

$$w^{new} = w^{old} - H^{-1} \nabla_w E$$
$$\rightarrow w^{new} = w^{old} - (\Phi^T R \Phi)^{-1} \Phi^T (y - t)$$

```
In [ ]: # IRLS Algorithm
w = np.ones((K, M)) # initialize with ones
total_error = [] # Log-Likelihood Function

for iteration in range(20):
    gradE = GradientE(X, w) # K, M
    H = Hessian(X, w) # KM x KM

    # Newton-Raphson update
    w = np.reshape(w,(K*M,1))
    gradE = np.reshape(gradE, (K*M,1))
    w -= np.dot(np.linalg.pinv(H), gradE)
    w = np.reshape(w, (K,M))

    total_error.append(log_likelihood(X,T,w)) # collect loss
```

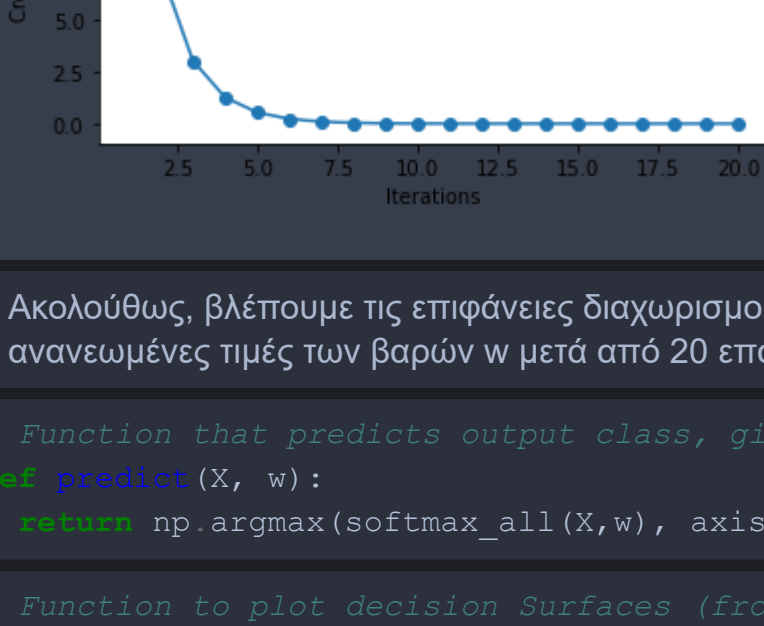
```
In [ ]: print('Weights After 20 Iterations of IRLS:')
print(w)
```

Weights After 20 Iterations of IRLS:  
[[ -0.55576221 -6.52705988 2.66298619]  
[ -12.35573972 1.31499084 -4.16030782]  
[ 15.04372858 -4.48760453 -8.97792324]]

```
In [ ]: print('Cross-Entropy Error after 20 Iterations of IRLS:', log_likelihood(X,T,w))
```

Cross-Entropy Error after 20 Iterations of IRLS: 2.864611554360134e-07

```
In [ ]: plt.plot([i for i in range(1,21)], total_error, '-o')
plt.xlabel('Iterations')
plt.ylabel('Cross-Entropy Loss')
plt.show()
```



Ακολούθως, βλέπουμε τις επιφάνειες διαχωρισμού που χωρίζουν τα δείγματα στις τρεις κατηγορίες, με βάση τις ανανεωμένες τιμές των βαρών  $w$  μετά από 20 επαναλήψεις του IRLS αλγορίθμου.

```
In [ ]: # Function that predicts output class, given weights w
def predict(X, w):
    return np.argmax(softmax_all(X,w), axis = 1)
```

```
In [ ]: # Function to plot decision Surfaces (from python helper lab)
def plot_clf(w, X, y, labels, clf = 'IRLS'):
    fig, ax = plt.subplots()
    # title for the plots
    title = ('Decision surface of LR Classifier')
    # Set-up grid for plotting.
    X0, X1, X2 = X[:,0], X[:,1], X[:,2]

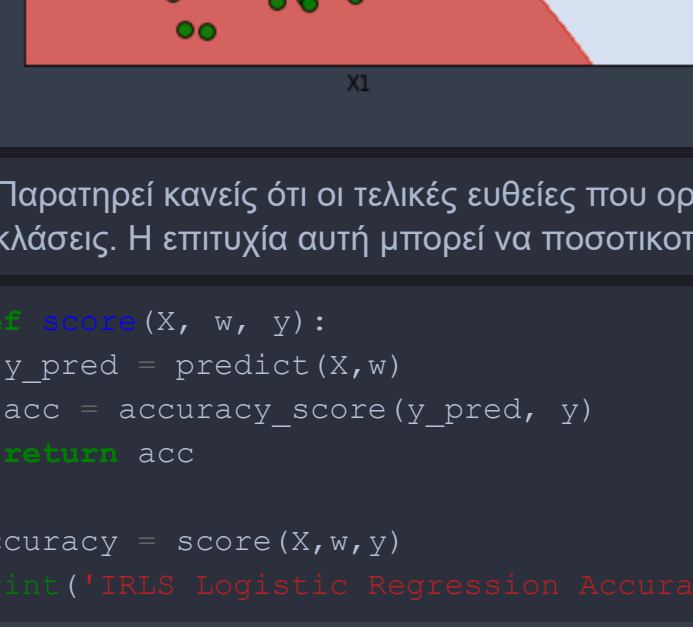
    x_min, x_max = X1.min() - 1, X1.max() + 1
    y_min, y_max = X2.min() - 1, X2.max() + 1
    xx, yy = np.meshgrid(np.arange(x_min, x_max, .05),
                        np.arange(y_min, y_max, .05))

    if clf == 'IRLS':
        Z = predict(np.column_stack((np.ones((len(xx.ravel()))), np.c_[xx.ravel(), yy.ravel()
])),w),
    else:
        Z = np.argmax(np.dot(np.column_stack((np.ones((len(xx.ravel()))), np.c_[xx.ravel(), yy.
ravel()]))),w), axis=1)
        Z = Z.reshape(xx.shape)
        out = ax.contourf(xx, yy, Z, cmap=plt.cm.coolwarm, alpha=0.8)

    zeros = ax.scatter(
        X1[y == 0], X2[y == 0],
        c='blue', label=labels[0],
        s=60, alpha=0.9, edgecolors='k')
    ones = ax.scatter(
        X1[y == 1], X2[y == 1],
        c='red', label=labels[1],
        s=60, alpha=0.9, edgecolors='k')
    twos = ax.scatter(
        X1[y == 2], X2[y == 2],
        c='green', label=labels[2],
        s=60, alpha=0.9, edgecolors='k')

    ax.set_ylabel('X2')
    ax.set_xlabel('X1')
    ax.set_xticks(())
    ax.set_yticks(())
    ax.set_title(title)
    ax.legend()
    plt.show()
```

```
In [ ]: plot_clf(w, X, y, labels=['Class 1', 'Class 2', 'Class 3'])
```



Παρατηρεί κανείς ότι οι τελικές ευθείες που ορίζουν τα βάρη  $w$  χωρίζουν τέλεια τα δείγματα εισόδου στις τρεις κλάσεις. Η επιτυχία αυτή μπορεί να ποσοτικοποιηθεί και με βάση την μετρική accuracy:

```
In [ ]: def score(X, w, y):
    y_pred = predict(X,w)
    acc = accuracy_score(y_pred, y)
    return acc

accuracy = score(X,w,y)
print('IRLS Logistic Regression Accuracy = '+str(100*accuracy)+'%')
```

IRLS Logistic Regression Accuracy = 100.0%

Παρατηρούμε ότι έχουμε ποσοστό ευστοχίας 100%.

## Linear Regression with Sum-of-Squares Loss Function

Στο πρόβλημα γραμμικής πρόβλεψης με συνάρτηση κόστους την Sum-of-Squares ο Newton-Raphson τύπος ανανέωσης βαρών παίρνει σε μόλις ένα βήμα την μορφή:

$$w^{new} = (\Phi^T \Phi)^{-1} \Phi^T t$$

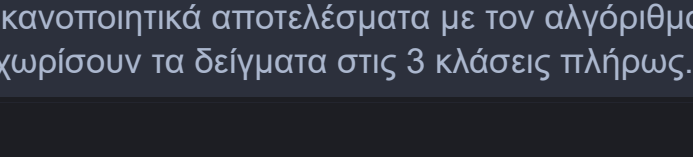
την οποία αναγνωρίζουμε ως τη κλασική λύση ελαχίστων τετραγώνων.

```
In [ ]: pseudo_inv = np.linalg.pinv(X)
w_lin = np.dot(pseudo_inv,T)
print('Weights of Linear Regression with S.o.S Function:')
print(w_lin)
```

Weights of Linear Regression with S.o.S Function:  
[[ 0.38015285 0.12967629 0.57017086]  
[ -0.05663662 0.08432925 -0.02769264]  
[ 0.124494785 0.01922407 -0.14417192]]

Παρουσιάζουμε στον χώρο τις τρεις ευθείες διαχωρισμού για τα άνωθι βάρη, καθώς και το ποσοστό ευστοχίας του ταξινομητή και την τιμή της συνάρτησης σφάλματος cross-entropy.

```
In [ ]: plot_clf(w_lin, X, y, labels=['Class 1', 'Class 2', 'Class 3'], clf='Linear Regression')
```



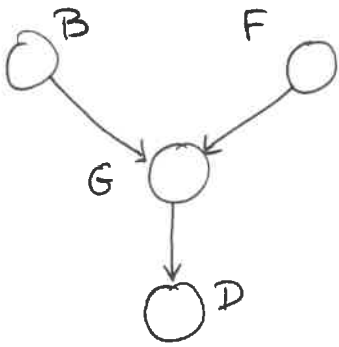
```
In [ ]: y_pred = np.argmax(np.dot(X, w_lin), axis=1)
acc2 = accuracy_score(y_pred, y)
print('Accuracy = '+str(100*acc2), '%')
```

Accuracy = 100.0 %

Παρατηρούμε ότι η λύση της linear regression με συνάρτηση κόστους την Sum of Squares δίνει εξίσου ικανοποιητικά αποτελέσματα με τον αλγόριθμο IRLS (Accuracy = 100%). Οι ευθείες διαχωρισμού καταφέρνουν να χωρίσουν τα δείγματα στις 3 κλάσεις πλήρως.



### Aktion 2.10: (Conditional Independence)



Battery B  
Fuel Tank F  
Gauge G  
Driver D

B	0	1
P(B)	0.05	0.95

F	O	L
P(F)	0.2	0.8

$P(G|BF)$ :

$BF \backslash G$	0	1
00	0.8	0.2
01	0.75	0.25
10	0.7	0.3
11	0.05	0.95

P(D|G):

G \ D	0	1
0	0.8	0.2
1	0.2	0.8

Παρατηρούμε ότι  $D = 0$

1.) Πιθανότητα το ντερόγιο  $F$  να είναι ίσδιο ( $F=0$ ), δεδομένου της παρατήρησης  $D=0$ , δηλαδή:

$$P(F=0 | D=0) = \frac{P(F=0, D=0)}{P(D=0)} = \frac{\sum_{B, G} P(B, G, F=0, D=0)}{\sum_{B, F, G} P(B, F, G, D=0)}$$

Επομένως :

B	G	F	D	$P(F)$	$P(B)$	$P(G BF)$	$P(D G)$	$P(BFGD)$
0	0	0	0	0.2	0.05	0.8	0.8	0.0064
0	1	0	0	0.2	0.05	0.2	0.2	0.0004
1	0	0	0	0.2	0.95	0.7	0.8	0.1064
1	1	0	0	0.2	0.95	0.3	0.2	0.0114
								0.1246

Telikius  $P(F=0|p=0) = \frac{0.1246}{0.3254} = \frac{38.29\%}{\underline{\underline{0.3829\%}}}$

B	F	G	D	$P(B)$	$P(F)$	$P(G BF)$	$P(D G)$	$P(BFGD)$
0	0	0	0	0.05	0.2	0.8	0.8	0.0064
0	0	1	0	0.05	0.2	0.2	0.2	0.0004 (+)
								0.0068

0.0328

Τελικώς:

$$P(F=0|D=0, B=0) = \frac{0.0068}{0.0328} = \underline{0.2073} \text{ ή } \underline{20.73\%} < 38.29\%$$

Παρατηρούμε ότι η άνωθεν πιθανότητα είναι μικρότερη από αυτήν που υπολογίσαμε στο υποερώτημα (1), γεγονός που οφείλεται στο explaining away, δηλαδή το μοτίβο συλλογιστικής σύμφωνα με το οποίο αν προκαθορίσουμε την <sup>αιτία</sup> ενός παρατηρούμενου <sup>πιδανός</sup> γεγονότος, τότε οι πιθανότητες που αφορούν εναλλακτικές αιτίες μειώνονται.

Στην προκειμένη περίπτωση, για το παρατηρούμενο γεγονός  $D=0$  προκαθορίζουμε ως αιτία την άδεια μπαταρία ( $B=0$ ) και ως αποτέλεσμα μειώνεται η πιθανότητα της αιτίας το άδειο ντεπόζιτο ( $F=0$ ).



## References

- [1] <https://towardsdatascience.com/the-mathematics-of-decision-trees-random-forest-and-feature-importance-in-scikit-learn-and-spark-f2861df67e3>
- [2] <https://sites.millersville.edu/bikenaga/linear-algebra/spectral-theorem/spectral-theorem.html>
- [3] C. M. Bishop, Pattern Recognition and Machine Learning, Springer, 2006. Ch. 4.3.2, 4.3.3, 4.3.4