

2η ΑΝΑΛΥΤΙΚΗ ΑΣΚΗΣΗ ΑΝΑΓΝΩΡΙΣΗ ΠΡΟΤΥΠΩΝ**Άσκηση 2.1: (Hidden Markov Models)**

Ζητείται να χρησιμοποιηθεί ένα HMM για να αποκωδικοποιηθεί μια απλή ακολουθία DNA. Είναι γνωστό ότι μια ακολουθία DNA είναι μια σειρά από στοιχεία του συνόλου $\{A, C, G, T\}$. Ας υποθέσουμε ότι υπάρχει μία κρυμμένη κατάσταση S που ελέγχει τη δημιουργία της ακολουθίας DNA και έχει 2 πιθανές καταστάσεις $\{S_1, S_2\}$. Επίσης, δίνονται οι ακόλουθες πιθανότητες μετάβασης για το HMM λ :

$$P(S_1|S_1) = 0.8 \quad P(S_2|S_1) = 0.2 \quad P(S_1|S_2) = 0.2 \quad P(S_2|S_2) = 0.8$$

οι ακόλουθες πιθανότητες των παρατηρήσεων:

$$P(A|S_1) = 0.4 \quad P(C|S_1) = 0.1 \quad P(G|S_1) = 0.4 \quad P(T|S_1) = 0.1$$

$$P(A|S_2) = 0.1 \quad P(C|S_2) = 0.4 \quad P(G|S_2) = 0.1 \quad P(T|S_2) = 0.4$$

και οι ακόλουθες a-priori πιθανότητες:

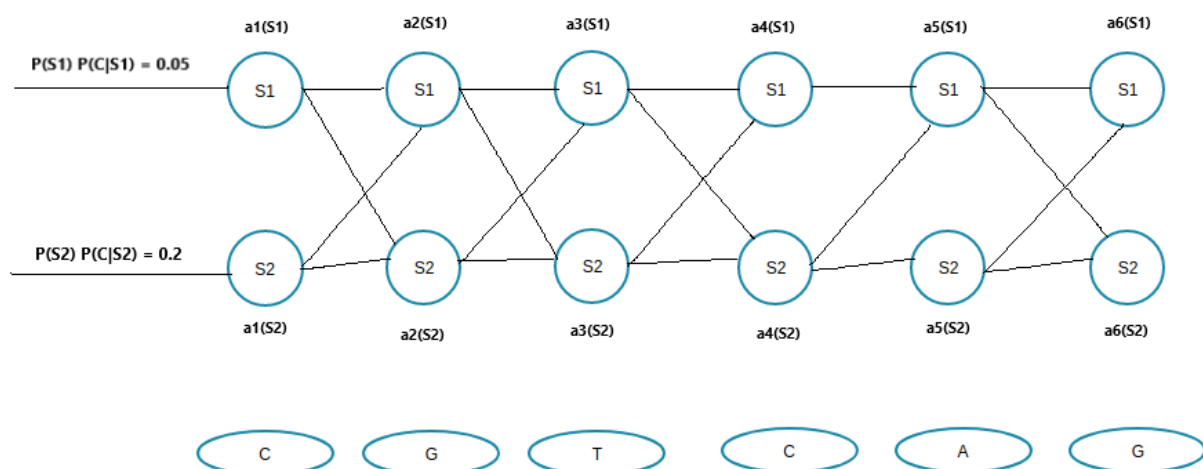
$$P(S_1) = 0.5 \quad P(S_2) = 0.5$$

Έστω ότι η παρατηρούμενη ακολουθία είναι $x = CGTCAG$. Υπολογίστε:

- (α) Την πιθανότητα $P(x|\lambda)$ χρησιμοποιώντας τον forward αλγόριθμο.
- (β) Τις εκ των υστέρων πιθανότητες $P(\pi_i = S_1|x, \lambda)$ για $i = 1, \dots, 6$.
- (γ) Το πιο πιθανό μονοπάτι κρυμμένων καταστάσεων χρησιμοποιώντας τον αλγόριθμο Viterbi.

Λύση**(α) Forward Probabilities**

Θέλουμε να υπολογίσουμε την πιθανότητα $P(x|\lambda)$ κάνοντας χρήση του forward αλγόριθμου, δεδομένου ότι η observed sequence είναι η CGTCAG. Παρακάτω, αναπαριστούμε το HMM model.



Μένει να υπολογίσουμε τις forward probabilities $a_i(S_i)$.

$$\begin{aligned} a_1(s_1) &= 0.05, \\ a_1(s_2) &= 0.2 \end{aligned}$$

$$\begin{aligned} a_2(S_1) &= P(S_1|S_1) \cdot P(G|S_1) \cdot a_1(S_1) + P(S_1|S_2) \cdot P(G|S_1) \cdot a_1(S_2) = 0.032 \\ a_2(S_2) &= P(S_2|S_2) \cdot P(G|S_2) \cdot a_1(S_2) + P(S_2|S_1) \cdot P(G|S_2) \cdot a_1(S_1) = 0.017 \end{aligned}$$

$$\begin{aligned} a_3(S_1) &= P(S_1|S_1) \cdot P(T|S_1) \cdot a_2(S_1) + P(S_1|S_2) \cdot P(T|S_1) \cdot a_2(S_2) = 0.029 \\ a_3(S_2) &= P(S_2|S_2) \cdot P(T|S_2) \cdot a_2(S_2) + P(S_2|S_1) \cdot P(T|S_2) \cdot a_2(S_1) = 0.008 \end{aligned}$$

$$\begin{aligned} a_4(S_1) &= P(S_1|S_1) \cdot P(C|S_1) \cdot a_3(S_1) + P(S_1|S_2) \cdot P(C|S_1) \cdot a_3(S_2) = 0.00248 \\ a_4(S_2) &= P(S_2|S_2) \cdot P(C|S_2) \cdot a_3(S_2) + P(S_2|S_1) \cdot P(C|S_2) \cdot a_3(S_1) = 0.00488 \end{aligned}$$

$$\begin{aligned} a_5(S_1) &= P(S_1|S_1) \cdot P(A|S_1) \cdot a_4(S_1) + P(S_1|S_2) \cdot P(A|S_1) \cdot a_4(S_2) = 0.0046976 \\ a_5(S_2) &= P(S_2|S_2) \cdot P(A|S_2) \cdot a_4(S_2) + P(S_2|S_1) \cdot P(A|S_2) \cdot a_4(S_1) = 0.00044 \end{aligned}$$

$$\begin{aligned} a_6(S_1) &= P(S_1|S_1) \cdot P(G|S_1) \cdot a_5(S_1) + P(S_1|S_2) \cdot P(G|S_1) \cdot a_5(S_2) = 0.001538432 \\ a_6(S_2) &= P(S_2|S_2) \cdot P(G|S_2) \cdot a_5(S_2) + P(S_2|S_1) \cdot P(G|S_2) \cdot a_5(S_1) = 0.000129152 \end{aligned}$$

Έχουμε:

$$P(\mathbf{x}|\lambda) = 0.001538432 + 0.000129152 = 0.001667584, \text{ αποτέλεσμα που προέκυψε με κάποιες στρογγυλοποιήσεις.}$$

(β) γ_i

Υπολογίζουμε πρώτα τα α_i , β_i και ύστερα τα γ_i .

Για το πρώτο state:

$$\alpha_1 = [0.05, 0.65306122, 0.26605505, 0.12311558, 0.60137931, 0.83628137]$$

Για το δεύτερο state:

$$\alpha_2 = [0.2, 0.34693878, 0.73394495, 0.87688442, 0.39862069, 0.16371863]$$

Για το πρώτο state:

$$\beta_1 = [5.12546277, 0.68531469, 0.85191279, 2.29206088, 1.26748252, 1]$$

Για το δεύτερο state:

$$\beta_2 = [3.71863431, 1.59234883, 1.05368161, 0.81859317, 0.59646236, 1]$$

Για το πρώτο state:

$$\gamma_1 = [0.25627314, 0.44755245, 0.2266557, 0.2821884, 0.76223776, 0.83628137]$$

Για το δεύτερο state:

$$\gamma_2 = [0.74372686, 0.55244755, 0.7733443, 0.7178116, 0.23776224, 0.16371863]$$

(γ) Viterbi algorithm

Αρχικοποιούμε $v_1(i) = \pi_i b_i$, $v_1(1) = 1/20$, $v_1(2) = 1/5$.

Επαναλαμβάνουμε την διαδικασία με $v_t(j) = \max[v_{t-1}(i)a_{ij}b_j(x_t)]$. Λαμβάνουμε:

$$v_2(1) = \max(\frac{32}{2000}, \frac{32}{2000})$$

$$v_2(2) = \max(\frac{1}{1000}, \frac{16}{1000})$$

$$v_3(1) = \max(\frac{256}{200000}, \frac{32}{100000})$$

$$v_3(2) = \max(\frac{256}{200000}, \frac{512}{100000})$$

$$v_4(1) = \max(\frac{1024}{10000000}, \frac{1024}{10000000})$$

$$v_4(2) = \max(\frac{1024}{10000000}, \frac{16384}{10000000})$$

$$v_5(1) = \max(\frac{32768}{10^9}, \frac{131072}{10^9})$$

$$v_5(2) = \max(\frac{2048}{10^9}, \frac{131072}{10^9})$$

$$v_6(1) = \max(\frac{4194304}{10^{11}}, \frac{1018576}{10^{11}})$$

$$v_6(2) = \max(\frac{262144}{10^{11}}, \frac{1018576}{10^{11}})$$

Επομένως, η πιο πιθανή ακολουθία με τον αλγόριθμο Viterbi είναι η $2- > 2- > 2- > 2- > 1- > 1$, με πιθανότητα να συμβεί $\frac{4194304}{10^{11}}$.

Άσκηση 2.2: (Principal Component Analysis)

Ζητείται η εφαρμογή της Principal Component Analysis (PCA) πάνω στο ευρέως διαδεδομένο σύνολο δεδομένων κρίνων του Fisher, προκειμένου να μετασχηματιστούν τα δεδομένα σε ένα χώρο χαμηλότερων διαστάσεων. Τα δεδομένα αποτελούνται από 3 κλάσεις (για τους 3 διαφορετικούς τύπους κρίνου), καθεμιά από τις οποίες περιλαμβάνει 50 δείγματα. Τα δεδομένα περιγράφονται από 4 διαφορετικά χαρακτηριστικά:

- μήκος σεπάλων σε εκ.
- πλάτος σεπάλων σε εκ.
- μήκος πετάλων σε εκ.
- πλάτος πετάλων σε εκ.
- τύπος κρίνου (Iris Setosa/Iris Versicolour/Iris Virginica)

(α) Κατεβάστε το σύνολο δεδομένων που έχει ανεβεί στο mycourses (αρχείο PCA.data).

(β) Προεπεξεργαστείτε τα δεδομένα αφαιρώντας τη μέση τιμή και διαιρώντας με την τυπική απόκλιση του κάθε χαρακτηριστικού ξεχωριστά. Τα προκύπτοντα δεδομένα θα πρέπει να έχουν μέση τιμή 0 και διασπορά 1.

(γ) Υπολογίστε τον δειγματικό πίνακα συνδιασπορών $\Sigma = \frac{1}{m-1} \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^T$, όπου \mathbf{x}_i είναι το i -στό δείγμα και m το πλήθος των δειγμάτων.

(δ) Παραγοντοποιήστε τον πίνακα συνδιασπορών κάνοντας χρήση του Singular Value Decomposition (SVD) και βρείτε τις αντίστοιχες ιδιοτιμές και ιδιοδιανύσματα. Προσέξτε εάν η συγκεκριμένη υλοποίηση του SVD δίνει τα αποτελέσματα με φθίνουσα ή αύξουσα σειρά ιδιοτιμών. Ο μετασχηματισμός SVD ενός πίνακα Σ είναι μια παραγοντοποίηση της μορφής $\Sigma = U D V^T$. Για έναν συμμετρικό, θετικά ορισμένο πίνακα Σ , οι $U = V$ περιέχουν τα ιδιοδιανύσματα και D είναι ένας διαγώνιος πίνακας με τις αντίστοιχες ιδιοτιμές.

(ε) Προβάλετε τα δεδομένα πάνω στις δύο πρώτες κύριες συνιστώσες και σχεδιάστε τα αποτελέσματα που προκύπτουν.

(στ) Ποιος είναι ο ελάχιστος απαιτούμενος αριθμός από κύριες συνιστώσες ώστε να “ερμηνεύεται” το 95% της διασποράς των τιμών;

Λύση

(α) Διάβασμα δεδομένων

```
1 import pandas as pd
2 import numpy as np
3
4 # Read Data and take rows and columns
5 df = pd.read_csv('PCA.data', na_values = '?', header = None)
6 print("First 5 records:\n")
7 print(df.head())
8
9 # Drop the last column
10 y = np.array(df[4])
11 del df[4]
```

First 5 records:

	0	1	2	3	4
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa

(β) Προεπεξεργασία δεδομένων

```
: 1 # Preprocessing
2 #  $(X - \mu) / \sigma$ 
3 # Standardization
4 data = np.array(df)
5 for i in range(4):
6     mu = np.mean(data[:,i])
7     sigma = np.sqrt(np.var(data[:,i]))
8     data[:,i] = (data[:,i] - mu) / sigma
```

(γ) Υπολογισμός δειγματικού πίνακα συνδιασπορών

```
: 1 M = len(data)
2 cov = np.zeros((4,4))
3 for i in range(M):
4     cov += np.matmul(np.array([data[i]]).transpose(), np.array([data[i]]))
5
6 cov = cov / (M - 1)
7 print('Δειγματικός πίνακας συνδιασπορών\n')
8 print(cov)
9
10 # print(np.cov(data, rowvar=False)) same result
```

Δειγματικός πίνακας συνδιασπορών

```
[[ 1.00671141 -0.11010327  0.87760486  0.82344326]
 [-0.11010327  1.00671141 -0.42333835 -0.358937 ]
 [ 0.87760486 -0.42333835  1.00671141  0.96921855]
 [ 0.82344326 -0.358937   0.96921855  1.00671141]]
```

(δ) Singular Value Decomposition (SVD)

```
: 1 u, d, vh = np.linalg.svd(cov, full_matrices=True)
2 print(u)
3 print()
4 print(np.diag(d))
5 print()
6 print(vh)
```

```
[[ -0.52237162 -0.37231836  0.72101681  0.26199559]
 [ 0.26335492 -0.92555649 -0.24203288 -0.12413481]
 [-0.58125401 -0.02109478 -0.14089226 -0.80115427]
 [-0.56561105 -0.06541577 -0.6338014  0.52354627]]
```

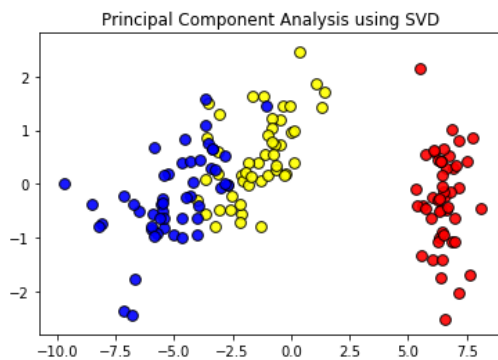
```
[[2.93035378 0.         0.         0.         ]
 [0.         0.92740362 0.         0.         ]
 [0.         0.         0.14834223 0.         ]
 [0.         0.         0.         0.02074601]]
```

```
[[ -0.52237162  0.26335492 -0.58125401 -0.56561105]
 [-0.37231836 -0.92555649 -0.02109478 -0.06541577]
 [ 0.72101681 -0.24203288 -0.14089226 -0.6338014 ]
 [ 0.26199559 -0.12413481 -0.80115427  0.52354627]]
```

Παρατηρούμε πως ο αρχικός δειγματικός πίνακας συνδιασπορών ήταν συμμετρικός και θετικά ορισμένος. Αυτό έχει σαν αποτέλεσμα οι $U = V$ και να περιέχουν τα ιδιοδιανύσματα, ενώ ο D να είναι ένας διαγώνιος πίνακας με τις αντίστοιχες ιδιοτιμές. Στον D οι ιδιοτιμές εμφανίζονται σε φθίνουσα σειρά.

(e) Principal Component Analysis (PCA)

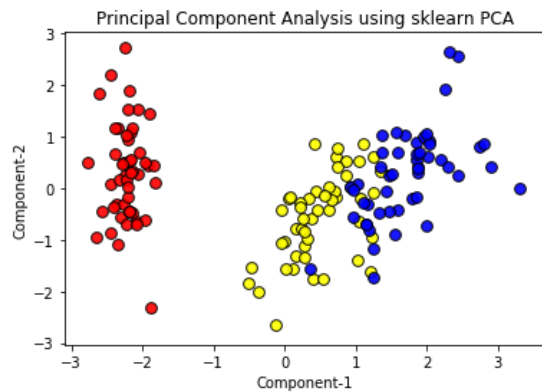
```
1 import matplotlib.pyplot as plt
2
3 u_new = u[:,0:2]
4 d = np.diag(d)
5 d_new = d[0:2,0:2]
6
7 b = np.dot(u_new, d_new)
8 dataPCA = np.dot(data, b)
9
10 # Plot
11 plt.title('Principal Component Analysis using SVD')
12 plt.scatter(dataPCA[0:50,0], dataPCA[0:50,1], color='red', edgecolors='k', s=60, alpha=0.9)
13 plt.scatter(dataPCA[50:100,0], dataPCA[50:100,1], color='yellow', edgecolors='k', s=60, alpha=0.9)
14 plt.scatter(dataPCA[100:150,0], dataPCA[100:150,1], color='blue', edgecolors='k', s=60, alpha=0.9)
15 plt.show()
```



```
1 # unique labels/types
2 print(pd.unique(y))

['Iris-setosa' 'Iris-versicolor' 'Iris-virginica']
```

```
1 from sklearn.decomposition import PCA
2 import matplotlib.pyplot as plt
3
4 pca = PCA(n_components=2)
5 dataPCA = pca.fit_transform(data)
6
7 # Plot
8 X0, X1 = dataPCA[:,0], dataPCA[:,1]
9 colors = {
10     'Iris-setosa': 'red',
11     'Iris-versicolor': 'yellow',
12     'Iris-virginica': 'blue'
13 }
14 fig, ax = plt.subplots()
15 for i, label in enumerate(colors.keys()):
16     ax.scatter(
17         X0[y == label], X1[y == label],
18         c=(colors[label]),
19         s=60, alpha=0.9, edgecolors='k'
20     )
21 ax.set_xlabel('Component-1')
22 ax.set_ylabel('Component-2')
23 ax.set_title('Principal Component Analysis using sklearn PCA')
24 plt.show()
```



(στ) PCA explained variation

```
1 print("Using PCA with SVD\n")
2
3 print('Explained variation for PCA(2): {}'.format(np.sum((d[0] + d[1]) / np.sum(d))))
4 print('Explained variation for PCA(1): {}'.format(d[0,0] / np.sum(d)))
```

Using PCA with SVD

Explained variation for PCA(2): 0.9580097536148198
Explained variation for PCA(1): 0.7277045209380135

```
1 print("Using sklearn PCA\n")
2 var2 = np.sum(pca.explained_variance_ratio_)
3 print('Explained variation for PCA(2): {}'.format(var2))
4
5 pca = PCA(n_components=1)
6 dataPCA = pca.fit_transform(data)
7 var1 = np.sum(pca.explained_variance_ratio_)
8 print('Explained variation for PCA(1): {}'.format(var1))
```

Using sklearn PCA

Explained variation for PCA(2): 0.9580097536148197
Explained variation for PCA(1): 0.7277045209380132

Παρατηρούμε ότι και με χρήση του SVD και με χρήση του sklearn, τα δυο PCA δίνουν ακριβώς τα ίδια αποτελέσματα, όπως και περιμέναμε. Βλέπουμε ακόμα ότι ο ελάχιστος αριθμός κύριων συνιστωσών που ερμηνεύεται τουλάχιστον το 95% της πληροφορίας (variance) των τιμών είναι για $n_components = 2$.

Άσκηση 2.3: (Linear Discriminant Analysis)

Στο μάθημα είδαμε ότι η Linear Discriminant Analysis (LDA) βασίζεται στην ανάστροφη σχέση των μητρών (πινάκων) S_W και S_B :

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

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

όπου το ω_i αναπαριστά μια κλάση με μέση τιμή $\vec{\mu}_i$, $|Classes|$ είναι το πλήθος των κλάσεων και $\vec{\mu}$ είναι η μέση τιμή όλων των δειγμάτων.

(α) Δείξτε ότι στην περίπτωση διαχωρισμού δύο κλάσεων ω_1 και ω_2 , ο πίνακας S_B μπορεί να γραφτεί στη μορφή $S_B = P(\omega_1)P(\omega_2)(\vec{\mu}_2 - \vec{\mu}_1)(\vec{\mu}_2 - \vec{\mu}_1)^T$.

(β) Βασιζόμενοι στο υποερώτημα (α), να βρείτε το ιδιοδιάνυσμα του πίνακα $S_W^{-1}S_B$ και την ιδιοτιμή του.

Λύση

(α)

Για $K = 2$, έχουμε:

$$\mathbf{S}_B = P(\omega_1)(\boldsymbol{\mu}_1 - \boldsymbol{\mu})(\boldsymbol{\mu}_1 - \boldsymbol{\mu})^T + P(\omega_2)(\boldsymbol{\mu}_2 - \boldsymbol{\mu})(\boldsymbol{\mu}_2 - \boldsymbol{\mu})^T$$

Γνωρίζουμε ότι το $\boldsymbol{\mu}$, ως ο μέσος όρος όλων των data points μπορεί να γραφεί:

$$\boldsymbol{\mu} = \frac{1}{N}(N_1\boldsymbol{\mu}_1 + N_2\boldsymbol{\mu}_2) = P(\omega_1)\boldsymbol{\mu}_1 + P(\omega_2)\boldsymbol{\mu}_2$$

Ακόμη, επειδή έχουμε μόνο δυο κατηγορίες (classes) θα έχουμε ότι $P(\omega_1) = 1 - P(\omega_2)$. Άρα, μπορούμε να γράψουμε:

$$\begin{aligned}\mathbf{S}_B &= P(\omega_1) \left((1 - P(\omega_1))\boldsymbol{\mu}_1 - P(\omega_2)\boldsymbol{\mu}_2 \right) \left((1 - P(\omega_1))\boldsymbol{\mu}_1 - P(\omega_2)\boldsymbol{\mu}_2 \right)^T \\ &\quad + P(\omega_2) \left((1 - P(\omega_2))\boldsymbol{\mu}_2 - P(\omega_1)\boldsymbol{\mu}_1 \right) \left((1 - P(\omega_2))\boldsymbol{\mu}_2 - P(\omega_1)\boldsymbol{\mu}_1 \right)^T \\ &= P(\omega_1) \left(P(\omega_2)\boldsymbol{\mu}_1 - P(\omega_2)\boldsymbol{\mu}_2 \right) \left(P(\omega_2)\boldsymbol{\mu}_1 - P(\omega_2)\boldsymbol{\mu}_2 \right)^T + P(\omega_2) \left(P(\omega_1)\boldsymbol{\mu}_2 - P(\omega_1)\boldsymbol{\mu}_1 \right) \left(P(\omega_1)\boldsymbol{\mu}_2 - P(\omega_1)\boldsymbol{\mu}_1 \right)^T \\ &= P(\omega_1)P(\omega_2)(P(\omega_1) + P(\omega_2)) (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \\ &= P(\omega_1)P(\omega_2) (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^T\end{aligned}$$

(β) Eigenvalue problem

Έχουμε:

$$\mathbf{S}_W^{-1}\mathbf{S}_B\mathbf{x} = \lambda\mathbf{x}$$

Αλλά

$$\mathbf{S}_B\mathbf{x} = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T = a(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

Αυτό σημαίνει πως το παραπάνω διάνυσμα δείχνει πάντα προς την κατεύθυνση του $\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$, αφού ο a είναι scalar. Έτσι, επιλύουμε το πρόβλημα των ιδιοτιμών απευθείας. Πιο αναλυτικά:

$$\mathbf{S}_W^{-1}\mathbf{S}_B\mathbf{x} = \mathbf{S}_W^{-1}[a(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)] = a[\mathbf{S}_W^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)]$$

όπου $\lambda = a$, η ιδιοτιμή, και $\mathbf{x} = \mathbf{S}_W^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$, το ιδιοδιάνυσμα.

Άσκηση 2.4: (Multilayer Perceptron)

Μας δίνονται 10 διανύσματα χαρακτηριστικών που προέρχονται από δύο κλάσεις ω_1 και ω_2 :

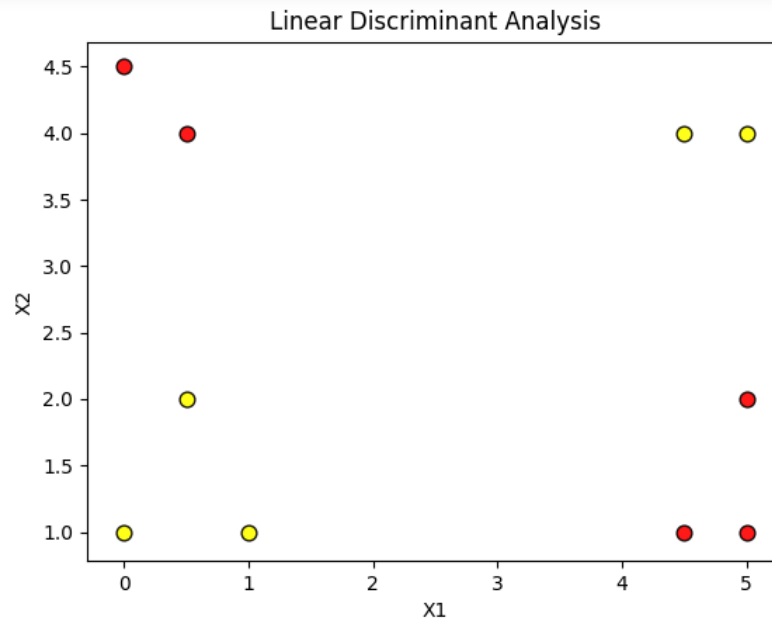
$$\omega_1 : [4.5, 1]^T, [5, 2]^T, [5, 1]^T, [0, 4.5]^T, [0.5, 4]^T$$

$$\omega_2 : [0, 1]^T, [0.5, 2]^T, [5, 4]^T, [4.5, 4]^T, [1, 1]^T$$

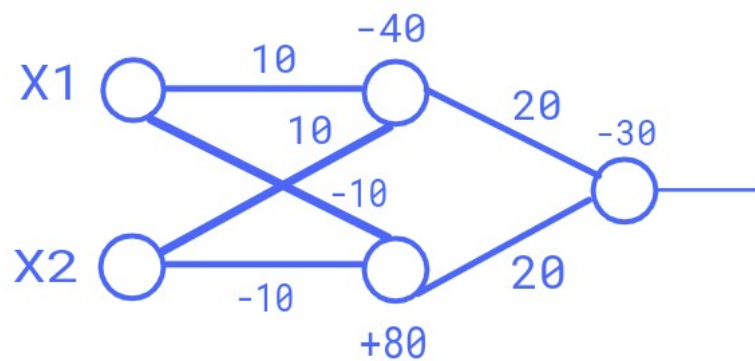
Ελέγξτε αν οι δύο κλάσεις είναι γραμμικά διαχωρίσιμες, και αν όχι, σχεδιάστε ένα κατάλληλο multilayer perceptron με τους κόμβους να έχουν βηματική συνάρτηση ενεργοποίησης (step)

Λύση

Ελέγχουμε αν τα δείγματα είναι γραμμικά διαχωρίσιμα. Από την παρακάτω γραφική παράσταση, βλέπουμε φανερά πως δεν είναι και, από την άλλη, συμπεραίνουμε ότι το πρόβλημά μας είναι παρόμοιο με αυτό της XOR.



Παρακάτω, σχεδιάζουμε το MLP που επιλύει πλήρως το πρόβλημά μας.



Για το παραπάνω MLP, το οποίο χρησιμοποιεί για activation function την step function: $f(x) = 1, x > 0$ αλλιώς -1 , τρέχουμε όλα τα data points που έχουμε δεδομένα και δείχνουμε πως κατηγοριοποιούνται όλα σωστά.

Για το [4.5,1]: $step(4.5(10) + 1(10) - 40) = 1$, $step(4.5(-10) + 1(-10) + 80) = 1$, $step(1(20) + 1(20) - 30) = 1$, αρα ορθώς κατηγοριοποιείται στην πρώτη κατηγορία.

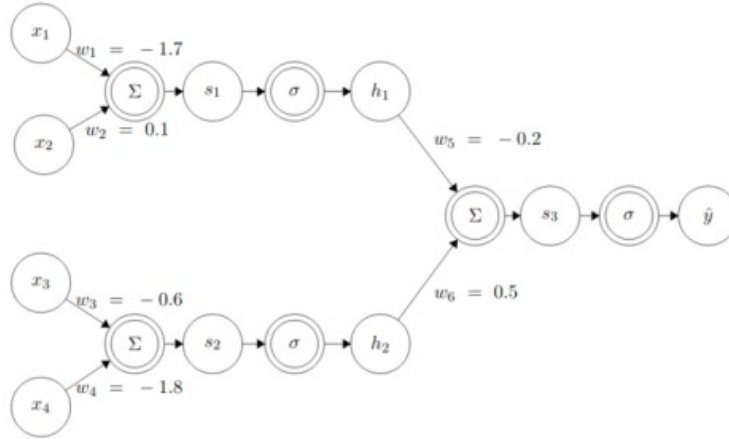
Για το [5,2]: $step(5(10) + 2(10) - 40) = 1$, $step(5(-10) + 2(-10) + 80) = 1$, $step(1(20) + 1(20) - 30) = 1$, αρα ορθώς κατηγοριοποιείται στην πρώτη κατηγορία.

Για το [5,1]: $step(5(10) + 1(10) - 40) = 1$, $step(5(-10) + 1(-10) + 80) = 1$, $step(1(20) + 1(20) - 30) = 1$, αρα ορθώς κατηγοριοποιείται στην πρώτη κατηγορία.

Ομοίως και τα για τα υπόλοιπα data points, όλα κατηγοριοποιούνται στις σωστές κατηγορίες.

Άσκηση 2.5: (Backpropagation)

Υποθέστε ότι έχουμε το ακόλουθο νευρωνικό δίκτυο. Οι κόμβοι που βρίσκονται μέσα σε μονό κύκλο υποδηλώνουν μεταβλητές (για παράδειγμα η x_1 είναι μια μεταβλητή εισόδου, h_1 είναι μια ενδιάμεση μεταβλητή, και \hat{y} είναι μια μεταβλητή εξόδου). Οι κόμβοι που βρίσκονται μέσα σε διπλό κύκλο υποδηλώνουν συναρτήσεις (για παράδειγμα το Σ υπολογίζει το άθροισμα των εισόδων του και η σ αναπαριστά τη συνάρτηση logistic $\sigma(x) = \frac{1}{1+e^{-x}}$).



Θεωρήστε ότι η συνάρτηση για το L2 loss δίνεται από τη σχέση $L(y, \hat{y}) = \|y - \hat{y}\|_2^2$. Επίσης, μας δίνονται τα δεδομένα ενός δείγματος $(x_1, x_2, x_3, x_4) = (-0.7, 1.2, 1.1, -2)$ με τιμή για το πραγματικό του label ίση με 0.5. Χρησιμοποιήστε τον αλγόριθμο backpropagation για να υπολογίσετε τη μερική παράγωγο $\frac{\partial L}{\partial w_1}$.

Σημείωση: Το gradient για μια συνάρτηση L2 loss είναι ίσο με $2\|y - \hat{y}\|$.

Λύση

Αρχικά, θα εκτελέσουμε forward propagation για να βρούμε τα s_1, s_2, s_3 και το y :

$$s_1 = -1.7(-0.7) + 0.1(1.2) = 1.31$$

$$h_1 = \sigma(s_1) = 0.78751$$

$$s_2 = 1.1(-0.6) + (-2)(-1.8) = 2.94$$

$$h_2 = \sigma(s_2) = 0.94979$$

$$s_3 = h_1(-0.2) + h_2(0.5) = 0.3174$$

$$y = \sigma(s_3) = 0.5787$$

Στην συνέχεια, εκτελούμε το backpropagation για να βρούμε τα deltas, $\delta = \frac{\partial L}{\partial s}$, και τελικά το $\frac{\partial L}{\partial w_1}$:

Για το output unit:

$$\delta_1^{(3)} = 2\|t - y\| \cdot \sigma'(s_3)$$

Για το hidden unit:

$$\delta_1^{(2)} = \sigma'(s_1) \cdot w_5 \cdot \delta_1^{(3)}$$

Τελικά, υπολογίζουμε το ζητούμενο derivative:

$$\frac{\partial L}{\partial w_1} = \delta_1^{(2)} \cdot x_1 = 8.99 \cdot 10^{-4}$$

Άσκηση 2.6: (Support Vector Machine)

Θεωρήστε το πρόβλημα του διαχωρισμού ενός συνόλου από διανύσματα εκπαίδευσης δύο κλάσεων. Τα δεδομένα εκπαίδευσης είναι της μορφής $\{(\mathbf{x}_i, y_i)\}$, όπου τα διανύσματα χαρακτηριστικών $\mathbf{x}_i \in \mathbb{R}^m$ και τα labels των κλάσεων $y_i \in \{-1, 1\}$.

Όπως είναι γνωστό, στην περίπτωση όπου τα δεδομένα εκπαίδευσης δεν είναι γραμμικώς διαχωρίσιμα (για παράδειγμα μέσω ενός κανόνα απόφασης $\text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$ για κάποια \mathbf{w}, b), τότε το πρόβλημα χρειάζεται να διατυπωθεί με χρήση slack variables $\{\xi_i\}$, $1 \leq i \leq n$. Έτσι, ο ταξινομητής SVM με το μεγαλύτερο περιθώριο αποκτάται μέσω της επίλυσης του δυϊκού προβλήματος:

$$L(\mathbf{w}, b, \alpha, \xi, \beta) = \frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i [(\mathbf{w} \cdot \mathbf{x}_i + b) y_i - (1 - \xi_i)] - \sum_{i=1}^n \beta_i \xi_i$$

όπου το C είναι μια σταθερά, $\alpha_i, \beta_i \geq 0$, $\forall i$ είναι οι πολλαπλασιαστές Lagrange, και $\xi_i \geq 0$ είναι οι slack variables.

(α) Υποθέστε ότι $n = 4$ και ότι το \mathbf{x} είναι δύο διαστάσεων $\langle x_1^1, x_1^2 \rangle: \langle 2, 2 \rangle, \langle 2.5, 2.5 \rangle, \langle 5, 5 \rangle, \langle 7, 7 \rangle$. Τώρα το SVM εκπαιδεύεται με βάση την παραπάνω εξίσωση. Να δείξετε ότι οποιαδήποτε labels y και αν έχουν τα τέσσερα δείγματα εκπαίδευσης, το βέλτιστο διάνυσμα

παραμέτρων $\hat{\mathbf{w}} = (\hat{w}^1, \hat{w}^2)$ έχει την ιδιότητα ότι $\hat{w}^1 = \hat{w}^2$.

(β) Θεωρήστε την εκπαίδευση ενός SVM με slack variables, αλλά δίχως την ύπαρξη του bias όρου ($b = 0$). Θα χρησιμοποιήσουμε έναν kernel $\mathbf{K}(\mathbf{u}, \mathbf{v})$ με την ιδιότητα ότι για δύο οποιαδήποτε σημεία \mathbf{u} και \mathbf{v} που ανήκουν στο σύνολο εκπαίδευσης ισχύει ότι $-1 < \mathbf{K}(\mathbf{u}, \mathbf{v}) < 1$. Επιπρόσθετα, $\mathbf{K}(\mathbf{u}, \mathbf{u}) < 1$. Να δείξετε ότι αν υπάρχουν n δείγματα στο σύνολο εκπαίδευσης και η σταθερά C επιλέγεται ούτως ώστε $C < \frac{1}{n-1}$, τότε όλες οι δυϊκές μεταβλητές α_i είναι μη μηδενικές (και άρα όλα τα δείγματα του συνόλου εκπαίδευσης αποτελούν support vectors).

(γ) Θεωρήστε τον εξής kernel:

$$\mathbf{K}(\mathbf{u}, \mathbf{v}) = \mathbf{u} \cdot \mathbf{v} + 4(\mathbf{u} \cdot \mathbf{v})^2$$

όπου τα διανύσματα \mathbf{u} και \mathbf{v} είναι δύο διαστάσεων. Ο kernel αυτός είναι ίσος με το εσωτερικό γινόμενο $\phi(\mathbf{u}) \cdot \phi(\mathbf{v})$ για κάποιο ορισμό της συνάρτησης ϕ . Ποια είναι η συνάρτηση αυτή ϕ ;

Λύση

(α) Example

Έχουμε ότι $\phi(\mathbf{x}_n) = \mathbf{x}_n$, και άρα μπορούμε να γράψουμε:

$$\begin{aligned} \mathbf{w} &= [w_1, w_2] \\ \mathbf{w} &= \sum_{n=1}^N a_n t_n \mathbf{x}_n = [2a_1 t_1 + 2.5a_2 t_2 + 5a_3 t_3 + 7a_4 t_4, 2a_1 t_1 + 2.5a_2 t_2 + 5a_3 t_3 + 7a_4 t_4] \end{aligned}$$

Άρα, λαμβάνουμε ότι $w_1 = w_2$, και άρα μπορούμε να γράψουμε

(β) Support Vectors

Έχουμε ότι:

$$-1 < K(\mathbf{u}, \mathbf{v}) < 1,$$

$$C < \frac{1}{N-1}$$

και θέλουμε να ότι όλες οι δυϊκές μεταβλητές a_n είναι μη μηδενικές.

Υποθέτουμε πως υπάρχει κάποια δυϊκή μεταβλητή $a_j = 0$ και θέλουμε να οδηγηθούμε σε άτοπο. Λαμβάνοντας υπόψιν ότι $b_n \geq 0$ και ότι $a_n + b_n = C$, παίρνουμε ότι:

$$a_n + b_n < \frac{1}{N-1} \Rightarrow a_n < \frac{1}{N-1}, \text{ για κάθε } n \in [1, N].$$

Ακόμη, επειδή $a_j = 0$ τότε $a_j + b_j = C \Rightarrow b_j \neq 0$. Και επειδή $b_j \xi_j = 0 \Rightarrow \xi_j = 0$. Έτσι, μπορούμε να γράψουμε:

$$\begin{aligned}\xi_j &= |t_j - y(\mathbf{x}_j)| = 0 \\ \Rightarrow \left| t_j - \sum_{n=1}^N a_n t_n K(\mathbf{x}_j, \mathbf{x}_n) \right| &= 0 \\ \Rightarrow t_j &= \sum_{n=1}^N a_n t_n K(\mathbf{x}_j, \mathbf{x}_n) \\ \Rightarrow |t_j| &= \left| \sum_{n=1}^N a_n t_n K(\mathbf{x}_j, \mathbf{x}_n) \right| \\ \Rightarrow 1 &= \left| \sum_{n=1}^N a_n t_n K(\mathbf{x}_j, \mathbf{x}_n) \right|\end{aligned}$$

Αρκεί να δείξουμε ότι η παραπάνω συνθήκη οδηγεί σε άτοπο. Από τριγωνική ανισότητα ισχύει ότι:

$$\left| \sum_{n=1}^N a_n t_n K(\mathbf{x}_j, \mathbf{x}_n) \right| \leq \sum_{n=1}^N |a_n t_n K(\mathbf{x}_j, \mathbf{x}_n)| = \sum_{n=1}^N |a_n| |t_n| |K(\mathbf{x}_j, \mathbf{x}_n)|$$

Έχουμε ακόμη ότι:

$$\begin{aligned}|a_n| &< \frac{1}{N-1} \\ |K(\mathbf{u}, \mathbf{v})| &< 1 \\ |t_n| &= 1\end{aligned}$$

Όμως, το παραπάνω άθροισμα περιέχει όλα τα μη αρνητικά a_n , και άρα αθροίζει συνολικά $N-1$ όρους. Άρα προκύπτει:

$$\left| \sum_{n=1}^N a_n t_n K(\mathbf{x}_j, \mathbf{x}_n) \right| \leq \sum_{n=1}^N |a_n t_n K(\mathbf{x}_j, \mathbf{x}_n)| = \sum_{n=1}^N |a_n| |t_n| |K(\mathbf{x}_j, \mathbf{x}_n)| < \frac{N-1}{N-1} = 1$$

Με αυτόν τον τρόπο, οδηγούμαστε σε άτοπον και άρα δεν υπάρχει κάποια δυϊκή μεταβλητή $a_j = 0$, με αποτέλεσμα όλα τα data points να αποτελούν support vectors.

(γ) Kernel example

Έχουμε ότι:

$$\begin{aligned}K(\mathbf{u}, \mathbf{v}) &= \mathbf{u}\mathbf{v} + 4(\mathbf{u}\mathbf{v})^2 \\ &= x_1 y_1 + x_2 y_2 + 4x_1^2 y_1^2 + 8x_1 x_2 y_1 y_2 + 4x_2^2 y_2^2 \\ &= [x_1, x_2, 2x_1^2, \sqrt{8}x_1 x_2, 2x_2^2] \cdot [y_1, y_2, 2y_1^2, \sqrt{8}y_1 y_2, 2y_2^2] \\ &= \phi(\mathbf{u}) \cdot \phi(\mathbf{v})\end{aligned}$$

Άσκηση 2.7: (Logistic Regression)

Θεωρήστε το πρόβλημα logistic regression για ένα σύνολο δεδομένων $\{\phi_n, t_n\}$, όπου $t_n \in \{0, 1\}$ και $\phi_n = \phi(\mathbf{x}_n)$ είναι οι κατηγορίες και οι συναρτήσεις βάσης, αντίστοιχα, για δείγματα $n = \{1, 2, \dots, N\}$. Η συνάρτηση σφάλματος $E(\mathbf{w})$, η οποία αναφέρεται συνήθως και ως cross-entropy, ορίζεται ως:

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

όπου \mathbf{w} είναι το διάνυσμα βαρών, $y_n = \sigma(\mathbf{w}^T \phi_n)$ η έξοδος του μοντέλου logistic regression στο διάνυσμα εισόδου \mathbf{x}_n , και $\sigma(a) = \frac{1}{1 + \exp(-a)}$ η logistic sigmoid συνάρτηση.

(α) Να δείξετε ότι για ένα γραμμικώς διαχωρίσιμο σύνολο δεδομένων, η λύση μέγιστης πιθανοφάνειας για το μοντέλο logistic regression αντιστοιχεί στην εύρεση ενός διανύσματος \mathbf{w} , για το οποίο η επιφάνεια απόφασης $\mathbf{w}^T \phi(\mathbf{x}) = 0$ διαχωρίζει τις κλάσεις, απειρίζοντας ταυτόχρονα το μέτρο του διανύσματος \mathbf{w} .

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

$$\mathbf{H} = [\Phi^T \mathbf{R} \Phi]$$

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

Να δείξετε ότι η Hessian μήτρα \mathbf{H} είναι θετικά ορισμένη. Ως εκ τούτου, δείξτε ότι η συνάρτηση σφάλματος είναι κυρτή συνάρτηση του \mathbf{w} και ότι έχει μοναδικό ελάχιστο.

(γ) Να γράψετε κώδικα που θα υλοποιεί τον iterative reweighted least squares (IWLS) αλγόριθμο για logistic regression. Χρησιμοποιώντας τον αλγόριθμο αυτό, να υπολογίσετε και να σχεδιάσετε τα διαχωριστικά επίπεδα απόφασης που αντιστοιχούν στο σύνολο δεδομένων του προβλήματος τριών κλάσεων που έχει ανεβεί στο mycourses (αρχείο `MLR.data`). Οι δύο πρώτες στήλες περιλαμβάνουν τα διανύσματα χαρακτηριστικών, ενώ η τρίτη την κλάση. Συγκρίνετε τα αποτελέσματα με εκείνα που θα προέκυπταν εάν εφαρμόζοταν ταξινόμηση με βάση τα ελάχιστα τετράγωνα, σχεδιάζοντας τα αντίστοιχα διαχωριστικά επίπεδα απόφασης.

Λύση

(α) Maximum Likelihood and the cross-entropy error function for a binary classification problem

Για ένα γραμμικά διαχωρίσιμο dataset $\{\phi_n, t_n\}$ όπου $t_n \in \{0, 1\}$ και $\phi_n = \phi(\mathbf{x}_n)$, και ακόμα $p(C_i|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi)$, με $\sigma(\cdot)$ να είναι η logistic sigmoid function, η likelihood function μπορεί να γραφεί:

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

Ο αρνητικός λογάριθμος της παραπάνω συνάρτησης είναι η cross-entropy error function. Άρα η μεγιστοποίηση της likelihood function ανάγεται στην ελαχιστοποίηση της cross-entropy error function. Δεδομένου ότι:

$$\frac{d\sigma}{da} = \sigma(1 - \sigma)$$

όπου a η ανεξάρτητη μεταβλητή της logistic sigmoid.

Γνωρίζουμε για τη logistic sigmoid ότι:

- αν $a > 0$ τότε $\sigma(a) > 0.5$,
- αν $a < 0$ τότε $\sigma(a) < 0.5$,
- αν $a = 0$ τότε $\sigma(a) = 0.5$

Οπότε η επιλογή της κλάσης C_i έναντι της C_j , με $i, j = \{0, 1\}$, $i \neq j$, κατά τη διαδικασία της πρόβλεψης, θα είναι όταν:

$$p(C_i|\phi) > p(C_j|\phi)$$

Έτσι, η μία θα επιλέγεται όταν $\sigma(\mathbf{w}^T \phi) > 0.5$ και η άλλη όταν το $\sigma(\mathbf{w}^T \phi) < 0.5$, πάντοτε με την αρχική μας προϋπόθεση ότι $p(C_i|\phi) = 1 - p(C_j|\phi)$. Άρα, το decision boundary θα είναι όταν:

$$\begin{aligned}\sigma(\mathbf{w}^T \phi) &= 0.5 \\ \Rightarrow \mathbf{w}^T \phi &= 0\end{aligned}$$

Η likelihood προκειμένου να μεγιστοποιηθεί, δηλαδή να ελαχιστοποιηθεί η cross-entropy error function, πρέπει για τα δείγματα i που $t_i = 0$ το $\mathbf{w}^T \phi$ να τείνει στο $-\infty$, ενώ για τα δείγματα j που $t_j = 1$ το $\mathbf{w}^T \phi$ να τείνει στο $+\infty$, αφού πρέπει στην πρώτη περίπτωση να ελαχιστοποιείται το $1 - \ln(y_i)$, ενώ στην δεύτερη να μεγιστοποιείται το $\ln(y_j)$. Δεδομένου ότι μπορούμε να γράψουμε, αφού έχουμε να κάνουμε με Euclidean vectors:

$$\mathbf{w}^T \phi = \|\mathbf{w}^T\| \cdot \|\phi\| \cdot \cos\theta,$$

όπου θ , η γωνία των δυο παραπάνω διανυσμάτων, και το $\|\phi\|$ σταθερό, τότε και στις δυο περιπτώσεις το $\|\mathbf{w}^T\|$ πρέπει να τείνει στο $+\infty$.

Άμεσο συμπέρασμα αυτών είναι ότι αυτή η μέθοδος του maximum likelihood είναι αρκετά ευάλωτη στο overfitting για γραμμικώς διαχωρίσιμα δείγματα, κάτι που μπορεί να αποφευχθεί με χρήση regularization ή με εισαγωγή priors και μέθοδο MAP.

(β) Logistic Regression and Hessian Matrix

Θεωρούμε αυθαίρετο vector $\mathbf{u} \neq \mathbf{0}$ μεγέθους M , όπου M ο αριθμός των χαρακτηριστικών. Προκειμένου ο \mathbf{H} να είναι θετικά ορισμένος (positive definite) αρκεί να δείξουμε ότι:

$$\mathbf{u}^T \mathbf{H} \mathbf{u} > 0$$

Έχουμε:

$$\begin{aligned}\mathbf{u}^T \mathbf{H} \mathbf{u} &= \mathbf{u}^T \left[\sum_{n=1}^N y_n(1 - y_n) \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T \right] \mathbf{u} \\ &= \sum_{n=1}^N y_n(1 - y_n) (\boldsymbol{\phi}_n^T \mathbf{u})^T (\boldsymbol{\phi}_n^T \mathbf{u}) \\ &= \sum_{n=1}^N y_n(1 - y_n) c_n^2 > 0\end{aligned}$$

αφού $c_n^2 > 0$, διότι θα έχει οπωσδήποτε τουλάχιστον μια μη μηδενική τιμή, και $y_n(1 - y_n) > 0$, διότι το y_i αποτελεί πιθανότητα. Άρα, ο \mathbf{H} είναι positive definite και άρα η error function είναι convex και έχει global minimum.

(γ) Iterative Reweighted Least Squares (IRLS) with Logistic Regression

```
1 # Read in the file
2 with open('MLR.data', 'r') as file :
3     filedata = file.read()
4
5 # Replace the "" string with ", "
6 filedata = filedata.replace(' ', ',')
7
8 # Write the file out again
9 with open('MLR.data', 'w') as file:
10     file.write(filedata)
```

```

1 import pandas as pd
2 import numpy as np
3
4 # Read Data and take rows and columns
5 df = pd.read_csv('MLR.data', na_values = '?', header = None)
6 print("First 5 records:\n")
7 print(df.head())
8
9 # Drop the last column
10 t = np.array(df[2])
11 del df[2]
12 X = np.array(df)
13 K = np.unique(t)

```

First 5 records:

```

      0      1  2
0 -3.605388  3.012493  0
1  8.941356  2.108067  1
2  4.161469 -2.959369  2
3  0.047052  4.166693  0
4  6.749173  1.589336  1

```

```

1 from sklearn.base import BaseEstimator, ClassifierMixin
2 from numpy import dot, multiply, squeeze, zeros, ones, reshape, array, identity
3 import numpy as np
4
5 class LogisticRegression(BaseEstimator, ClassifierMixin):
6     # Iterative Reweighted Least Squares (IRLS) for Multiclass Logistic Regression
7     def __init__(self, K, maxiter=10):
8         """
9             K: int, classes
10            maxiter: int
11        """
12
13        self.W = None
14        self.K = K
15        self.maxiter = maxiter
16
17
18    @staticmethod
19    def stable_softmax(x):
20        """
21            x: list of scalars
22        """
23        z = x - max(x)
24        numerator = np.exp(z)
25        denominator = np.sum(numerator)
26        softmax = numerator / denominator
27
28        return softmax
29
30
31    @staticmethod
32    def softmax(W, X, K):
33        """
34            W: K x M
35            X: N x M
36            K: int
37        """
38
39        N = X.shape[0]
40        M = X.shape[1]
41        softmax_all = zeros((N,K))
42
43        for i,xn in enumerate(X):
44            xn = reshape(xn, (M,1)) # M x 1
45
46            A = []
47            for k in range(K):
48                a_nk = squeeze(dot( reshape(W[k], (1,M)), xn )) # scalar
49                A.append(a_nk)
50
51            softmax_all[i,:] = reshape(array(stable_softmax(A)), (1,K))
52
53        return softmax_all # N x K
54
55
56

```



```

57 def fit(self, X, t):
58     """
59     X: N x M
60     t: list of targets
61     """
62
63     K = self.K
64     X_t = X.transpose()
65
66     N = X.shape[0] # n_samples
67     M = X.shape[1] # n_features
68     I = np.identity(K) # K x K
69
70     # 1-of-K coding scheme for target values
71     T = zeros((N,K)) # N x K
72     for i,k in enumerate(t):
73         T[i][k] = 1
74
75     # Initialize w, grad, hessian
76     self.W = zeros(shape=(K,M)) # K x M
77     grad = ones(shape=(K,M)) # K x M
78     H = ones(shape=(K*M,K*M)) # K*M x K*M
79
80     for iteration in range(self.maxiter):
81
82         # Y = softmax for all n,k
83         Y = softmax(W, X, K) # N x K
84
85         # gradient of cross-entropy error function
86         grad = dot( X_t, Y - T ).transpose() # K x M
87
88         # Hessian Matrix: K*M x K*M
89         for k in range(K):
90             for j in range(K):
91                 Ikj = I[k][j]
92                 # Diagonal R
93                 R = zeros((N,N))
94                 for n in range(N):
95                     R[n][n] = Y[n][k] * (Ikj - Y[n][j])
96
97                 H[k*M:k*M+M, j*M:j*M+M] = dot( dot(X_t, R), X ) # M x M
98
99         # Newton-Raphson update
100         self.W = reshape(self.W, (K*M,1))
101         grad = reshape(grad, (K*M,1))
102         Wnew = self.W - dot( np.linalg.pinv(H), grad )
103
104         # keep preferred dimensions
105         grad = reshape(grad, (K,M))
106         self.W = reshape(Wnew, (K,M))
107
108     continue
109
110     return self
111
112
113 def predict(self, X):
114     # Logistic Regression predict on given test data
115     solve = softmax(W, X, self.K)
116     y_predict = np.argmax(solve, axis=1)
117
118     return y_predict
119
120
121 def score(self, X, y_truth):
122     # Logistic Regression accuracy score on given test data
123     y_predict = self.predict(X)
124
125     return ( np.sum([y_predict == y_truth]) ) / ( len(y_predict) )

```



```

1 lr = LogisticRegression(K=3, maxiter=5)
2 lr.fit(X, t)
3 print("IRLS LogisticRegression Score:\t", lr.score(X, t) * 100, "%")
4 print("Final W:\n", lr.W)

```

IRLS LogisticRegression Score: 100.0 %

Final W:

```

[[-0.35798298  1.48913729]
 [ 0.39032115  0.39589378]
 [-0.03233817 -1.88503107]]

```

```

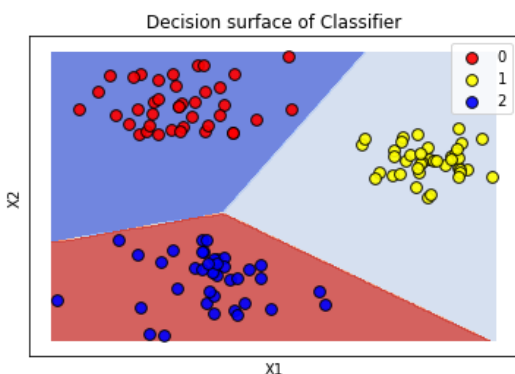
1 import warnings
2 warnings.filterwarnings("ignore",category=DeprecationWarning)
3 warnings.filterwarnings("ignore",category=FutureWarning)
4
5 import numpy as np
6 import matplotlib.pyplot as plt
7
8 def plt_decision_boundaries(clf, X, y, labels, plot_labels):
9     # Plot Decision boundaries for a 10 labels classification problem, supposed that labers are in range[0,3]
10     ...
11     clf:          classifier that has fit, predict and score methods
12     X:            X_train
13     y:            y_train
14     labels:       an numpy array of labels
15     plot_labels:  a list of the axis titles on the plot
16     ...
17     fig, ax = plt.subplots()
18     # title for the plots
19     title = ('Decision surface of Classifier')
20
21     # Initializations
22     X0, X1 = X[:,0], X[:,1]
23
24     x_min, x_max = X0.min() - 0.2, X0.max() + 0.2
25     y_min, y_max = X1.min() - 0.2, X1.max() + 0.2
26     xx, yy = np.meshgrid(np.arange(x_min, x_max, .05),
27                           np.arange(y_min, y_max, .05))
28
29     Z = clf.predict(np.c_[xx.ravel(), yy.ravel()])
30     Z = Z.reshape(xx.shape)
31     out = ax.contourf(xx, yy, Z, cmap=plt.cm.coolwarm, alpha=0.8)
32
33     colors = ['red', 'yellow', 'blue']
34
35     for label in labels:
36         ax.scatter(
37             X0[y == label], X1[y == label],
38             c=(colors[int(label)]), label=int(label),
39             s=60, alpha=0.9, edgecolors='k'
40         )
41
42     ax.set_ylabel(plot_labels[1])
43     ax.set_xlabel(plot_labels[0])
44     ax.set_xticks(())
45     ax.set_yticks(())
46     ax.set_title(title)
47     ax.legend()
48     plt.show()

```

```

1 lr = LogisticRegression(K=3, maxiter=5)
2 lr.fit(X, t)
3 plt_decision_boundaries(lr, X, t, np.arange(3), ["X1", "X2"])

```



Linear Regression and Mean Squared Error cost function (Normal Equations method)

```

1 from sklearn.base import BaseEstimator, ClassifierMixin
2 from numpy import dot, multiply, squeeze, zeros, ones, reshape, array, identity
3 import numpy as np
4
5 class LinearRegression(BaseEstimator, ClassifierMixin):
6     # Normal Equations for Multiclass Logistic Regression
7     def __init__(self, K):
8         """
9         K: int, classes
10        """
11
12        self.Wml = None
13        self.K = K
14
15    def fit(self, X, t):
16        """
17        X: N x M
18        t: list of targets
19        """
20
21        K = self.K
22        N = X.shape[0] # n_samples
23        M = X.shape[1] # n_features
24        X_t = X.transpose()
25
26        # 1-of-K coding scheme for target values
27        T = zeros((N,K)) # N x K
28        for i,k in enumerate(t):
29            T[i][k] = 1
30
31        # Mean Squared Error cost function - Maximum Likelihood solution
32        pseudo_inv = dot( np.linalg.pinv(dot(X_t, X)), X_t ) # pseudo-inverse
33        self.Wml = dot( pseudo_inv, T)
34
35        return self
36
37    def predict(self, X):
38        # Linear Regression predict on given test data
39        Y = dot( X, self.Wml )
40        y_predict = np.argmax(Y, axis=1)
41
42        return y_predict
43
44    def score(self, X, y_truth):
45        # Linear Regression accuracy score on given test data
46        y_predict = self.predict(X)
47
48        return ( np.sum([y_predict == y_truth]) ) / ( len(y_predict) )

```

```

1 lr = LinearRegression(K=3)
2 lr.fit(X, t)
3 print("Normal Equations LinearRegression Score:\t", lr.score(X, t) * 100, "%")
4 print("Wml:\n", lr.Wml)

```

Normal Equations LinearRegression Score: 99.16666666666667 %

Wml:

```

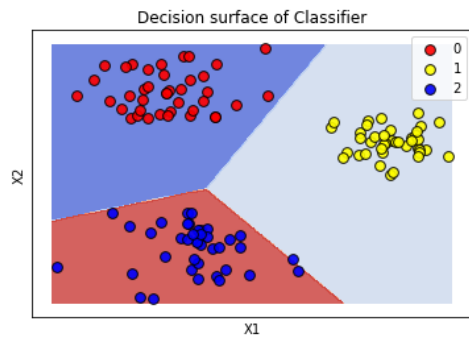
[[-0.03973394  0.09163179  0.00441572]
 [ 0.15567397  0.0324988  -0.08580454]]

```

```

1 lr = LinearRegression(K=3)
2 lr.fit(X, t)
3 plt_decision_boundaries(lr, X, t, np.arange(3), ["X1", "X2"])

```



Παρατηρούμε πως και τα δυο μοντέλα δίνουν σχεδόν άριστα ποσοστά. Το Logistic Regression δίνει 100%, ενώ το Linear Regression δίνει 99.17%. Αυτό οφείλεται στην διαφορετική cost function καθώς και στην activation function που χρησιμοποιούν. Βλέπουμε πως αν και τα δείγματα είναι linearly separable, η χρήση της softmax και της cross-entropy καταφέρνουν να μην κάνουν κανένα λάθος, τοποθετώντας τα δείγματα πιο κεντρικά στην αντίστοιχη περιοχή απόφασης, σε αντίθεση με τη χρήση των normal equations και της mean-squared που κάποια δείγματα τοποθετούνται πολύ κοντά στο decision boundary, με αποτέλεσμα ένα να γίνεται misclassified.