

Γενικές Οδηγίες: Οι αναλυτικές σειρές ασκήσεων είναι ατομικές, και οι λύσεις που θα δώσετε πρέπει να αντιπροσωπεύουν μόνο την προσωπική σας εργασία. Εξηγήστε επαρκώς την εργασία σας. Αν χρησιμοποιήσετε κάποια άλλη πηγή εκτός των σημειώσεων για την λύση σας, πρέπει να το αναφέρετε. Η παράδοση των λύσεων των αναλυτικών ασκήσεων της σειράς αυτής θα γίνει ηλεκτρονικά στην HELIOS ιστοσελίδα του μαθήματος και θα πρέπει να την υποβάλετε ως ένα ενιαίο αρχείο PDF με το εξής filename format χρησιμοποιώντας μόνο λατινικούς χαρακτήρες: ML22_hwk2_AM_LastnameFirstname.pdf, όπου AM είναι ο 8-ψήφιος αριθμός μητρώου σας. Σκαναρισμένες χειρόγραφες λύσεις επιτρέπονται αρκεί να είναι καθαρογραμμένες και ευανάγνωστες. Επίσης στην 1η σελίδα των λύσεων θα αναγράφεται το ονοματεπώνυμο, Α.Μ., και email address σας. Συμπεριλάβετε και τον κώδικα προγραμμάτων, π.χ. Matlab ή Python, που χρησιμοποιήσατε για αριθμητική επίλυση. Να σημειωθεί ότι η καταληκτική ημερομηνία παράδοσης είναι τελική και δεν θα δοθεί παράταση.

Άσκηση 2.1 (Θεωρία Μηχανικής Μάθησης)

Να δείξετε ότι η κλάση H_{rec}^n των παράλληλων στους άξονες υπερ-παράλληλογράμμων του \mathbb{R}^n είναι PAC εκπαιδευσιμη.

Υπόδειξη. Μπορείτε να συμβολίσετε την κλάση ως $H_{\text{rec}}^n = \{h_{(a_1, a_2, \dots, a_n, b_1, b_2, \dots, b_n)} : a_1 \leq b_1, a_2 \leq b_2, \dots, a_n \leq b_n\}$, όπου $h_{(a_1, a_2, \dots, a_n, b_1, b_2, \dots, b_n)}(x_1, x_2, \dots, x_n) = 1$ για εκείνα τα $(x_1, x_2, \dots, x_n) \in \mathbb{R}^n$ για τα οποία $a_1 \leq x_1 \leq b_1, a_2 \leq x_2 \leq b_2, \dots, a_n \leq x_n \leq b_n$. Για την απόδειξη, μπορείτε να ακολουθήσετε (με τις απαραίτητες επεκτάσεις) τη μεθοδολογία που ακολουθήσαμε για την αντίστοιχη απόδειξη για $n = 2$. Εναλλακτικά μπορείτε να υπολογίσετε τη διάσταση VC της κλάσης H_{rec}^n και να δείξετε ότι είναι πεπερασμένη για κάθε n .

Άσκηση 2.2 (Σύγκριση των Αλγορίθμων k -Means και Fuzzy c-Means)

(α) Όπως αναφέρεται στη σελίδα 38 στη σειρά διαφανειών #9, ο αλγόριθμος fuzzy c-means ελαχιστοποιεί τη συνάρτηση κόστους

$$J(U, \theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q d(x_i, \theta_j), \quad q > 1$$

με τους ακόλουθους περιορισμούς

$$\sum_{j=1}^m u_{ij} = 1, \quad i = 1, 2, \dots, N.$$

Υπενθυμίζεται ότι για τις συναρτήσεις συμμετοχής u_{ij} ισχύει επιπλέον $0 < u_{ij} < 1$. Έστω ότι η απόσταση $d(x_i, \theta_j) = \|x_i - \theta_j\|^2$, δηλαδή ισούται με το τετράγωνο της Ευκλείδειας απόστασης του σημείου x_i από το κέντρο θ_j . Χρησιμοποιώντας τη μέθοδο των πολλαπλασιαστών Lagrange, να δείξετε ότι στο πλαίσιο μιας διαδικασίας εναλλασσόμενης ελαχιστοποίησης της συνάρτησης $J(U, \theta)$, μπορούμε να καταλήξουμε στις ακόλουθες αναδρομικές σχέσεις για τα u_{ij} και τις εκτιμήσεις των κέντρων των κλάσεων θ_j

$$u_{ij} = \frac{1}{\sum_{k=1}^m \left(\frac{d(x_i, \theta_j)}{d(x_i, \theta_k)} \right)^{\frac{1}{q-1}}}, \quad \theta_j = \frac{\sum_{i=1}^N u_{ij}^q x_i}{\sum_{i=1}^N u_{ij}^q}.$$

Θα πρέπει να σημειωθεί ότι οι δύο αυτές σχέσεις αποτελούν τα βασικά βήματα του αλγόριθμου fuzzy c-means.

(β) Έστω ένα πρόβλημα ομαδοποίησης δεδομένων σε $m = 2$ κλάσεις, τα σημεία των οποίων ακολουθούν την κανονική κατανομή με μέσα διανύσματα $\mu_1 = [1, 1]^T$, $\mu_2 = [2.5, 2.5]^T$ και πίνακες συμεταβλητότητας

$$\Sigma_1 = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}, \quad \Sigma_2 = \begin{bmatrix} 1 & -0.4 \\ -0.4 & 1 \end{bmatrix}.$$

Δημιουργήστε ένα σύνολο δεδομένων με 300 σημεία από κάθε κλάση. Δώστε μια κατάλληλη γραφική αναπαράσταση του συνόλου των δεδομένων και σχολιάστε.

(γ) Υλοποιήστε (σε Python ή Matlab) τον αλγόριθμο k -means με τυχαία αρχικοποίηση και εφαρμόστε τον στο παραπάνω σύνολο δεδομένων. Ο αλγόριθμος τερματίζει όταν ικανοποιείται το κριτήριο που δίνεται στη σελίδα 32 της 9ης σειράς διαφανειών με $\varepsilon = 10^{-4}$. Δώστε γραφική απεικόνιση της ομαδοποίησης στην οποία καταλήγει ο k -means και των κέντρων των κλάσεων που εκτιμά ο αλγόριθμος.

(δ) Επαναλάβετε τα ζητούμενα του ερωτήματος (γ) για τον αλγόριθμο fuzzy c-means με $q = 2$. Ο αλγόριθμος αρχικοποιείται με τα ίδια διανύσματα $\theta_j, j = 1, 2, \dots, m$ με τα οποία αρχικοποιήθηκε και ο k -means στο προηγούμενο ερώτημα. Μετά τον τερματισμό του fuzzy c-means, το διάνυσμα x_i καταχωρείται στην κλάση j με το μέγιστο $u_{ij}, j = 1, 2, \dots, m$.

(ε) Με βάση τα αποτελέσματα στα δύο προηγούμενα ερωτήματα, συγκρίνετε τους δύο αλγόριθμους ως προς i) τον αριθμό των επαναλήψεων που απαιτούνται μέχρι τον τερματισμό τους, ii) τη μέση απόσταση των εκτιμώμενων κέντρων από τα πραγματικά κέντρα και iii) το ρυθμό επιτυχίας (success rate). Μπορείτε να τρέξετε τους αλγόριθμους για πολλά διαφορετικά σύνολα δεδομένων [παρόμοια με αυτό του ερωτήματος (β)] για να επιβεβαιώσετε τα συμπεράσματά σας.

Άσκηση 2.3 (Ιεραρχική Ομαδοποίηση)

(α) Θεωρήστε δύο σημεία στο \mathbb{R}^l , $\mathbf{x} = [x_1, x_2, \dots, x_l]^T$ και $\mathbf{y} = [y_1, y_2, \dots, y_l]^T$ και έστω $|x_i - y_i| = \max_{j=1,2,\dots,l} \{|x_j - y_j|\}$. Ορίζουμε την απόσταση $d_n(\mathbf{x}, \mathbf{y})$ ως εξής

$$d_n(\mathbf{x}, \mathbf{y}) = |x_i - y_i| + \frac{1}{l/2 + 1} \sum_{j=1, j \neq i}^l |x_j - y_j|.$$

Να αποδείξετε ότι η d_n είναι μετρική. (Μπορεί να δείχτεί ότι η d_n αποτελεί μια προσέγγιση της Ευκλείδειας απόστασης με χαμηλότερη υπολογιστική πολυπλοκότητα.)

(β) Δίνεται το σύνολο προτύπων στο \mathbb{R}^2 ,

$$X = \{(0, 3), (1.4, 2.6), (-1.5, 3.4), (-0.2, -0.4), (1, -1), (2, -1.5), (2.6, -1.8), (3, -2)\}.$$

Δώστε τον πίνακα προτύπων $D(X)$ και προσδιορίστε τον πίνακα εγγύτητας $P(X)$ με βάση την μετρική d_n .

(γ) Με βάση τον πίνακα εγγύτητας $P(X)$ που υπολογίσατε στο ερώτημα (β), δώστε τις διαδοχικές ομαδοποιήσεις που θα προκύψουν από την εφαρμογή του ιεραρχικού αλγόριθμου απλού δεσμού, καθώς και το αντίστοιχο δενδρόγραμμα εγγύτητας.

(δ) Επαναλάβετε τα ζητούμενα στο ερώτημα (γ) για τον ιεραρχικό αλγόριθμο πλήρους δεσμού.

(ε) Συγκρίνετε και σχολιάστε τα αποτελέσματα που πήρατε στα δύο προηγούμενα ερωτήματα και προσδιορίστε τη βέλτιστη ομαδοποίηση σε κάθε περίπτωση από τα αντίστοιχα δενδρογράμματα.

Άσκηση 2.4 (Θεωρία PCA και SVD)

Μας δίνεται μια ακολουθία δεδομένων (τυχαία διανύσματα με μηδενικό μέσο) $\mathbf{x}_n \in \mathbb{R}^d, n = 1, \dots, N$, και θέλουμε να βρούμε μια κατεύθυνση (μοναδιαίο διάνυσμα) $\mathbf{e} \in \mathbb{R}^d$ και σταθερές a_n έτσι ώστε, αν προσεγγίσουμε κάθε δεδομένο μας (διάνυσμα στήλης) \mathbf{x}_n με ένα διάνυσμα $a_n \mathbf{e}$, το συνολικό μέσο τετραγωνικό λάθος J να είναι ελάχιστο:

$$J(a_1, \dots, a_n, \mathbf{e}) = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - a_n \mathbf{e}\|^2, \quad \|\cdot\| = \text{Euclidean norm}$$

(α) Θεωρήστε γνωστό το \mathbf{e} και αποδείξτε ότι τα βέλτιστα a_n που ελαχιστοποιούν το J είναι

$$a_n = \langle \mathbf{x}_n, \mathbf{e} \rangle = \mathbf{x}_n^T \mathbf{e}$$

(β) Αντικαθιστώντας τα βέλτιστα a_n στο J , αποδείξτε ότι αποκτούμε ένα λάθος

$$J_1(\mathbf{e}) = -\mathbf{e}^T \mathbf{R}_x \mathbf{e} + (1/N) \sum_{n=1}^N \|\mathbf{x}_n\|^2, \quad \mathbf{R}_x = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T,$$

όπου \mathbf{R}_x είναι ο εμπειρικός πίνακας αυτοσυσχέτισης των δεδομένων.

(γ) Αποδείξτε ότι το βέλτιστο \mathbf{e} που ελαχιστοποιεί το J_1 είναι το ιδιοδιάνυσμα του \mathbf{R}_x που αντιστοιχεί στην μέγιστη ιδιοτιμή λ_1 . (Υπόδειξη: ελαχιστοποιήστε το J_1 εκμεταλλευόμενοι τον περιορισμό $\|\mathbf{e}\| = 1$ με Lagrange πολλαπλασιαστή.)

(δ) Τι σχέση έχει η ανωτέρω λύση με PCA (Principal Component Analysis)?

(ε) Να βρεθεί αναλυτικά πως σχετίζεται η ανωτέρω λύση με την SVD (Singular Value Decomposition) του $N \times d$ πίνακα \mathbf{X} που σχηματίζεται στιβάζοντας τα διανύσματα $\mathbf{x}_n, n = 1, \dots, N$, ως γραμμές?

Άσκηση 2.5 (Αριθμητική εφαρμογή του PCA)

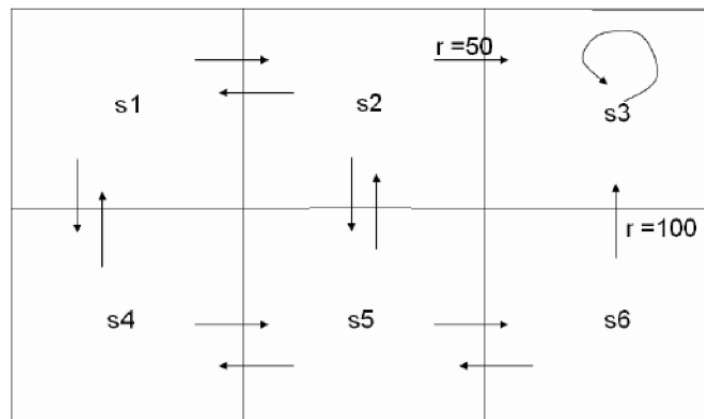
Ζητείται η εφαρμογή της Principal Component Analysis (PCA) πάνω στο ευρέως διαδεδομένο σύνολο δεδομένων κρίνων (Iris dataset) του Fisher, προκειμένου να μετασχηματιστούν τα δεδομένα σε ένα χώρο χαμηλότερων διαστάσεων. Τα δεδομένα αποτελούνται από 3 κλάσεις (για τους 3 διαφορετικούς τύπους κρίνου), καθεμιά 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 στην Helios ιστοσελίδα του μαθήματος (αρχείο PCA_iris_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 (sample covariance matrix) $C = (1/N) \sum_{n=1}^N x_n x_n^T$, όπου x_n είναι το n -στό δείγμα και N 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 C είναι μια παραγοντοποίηση της μορφής $C = U \Sigma V^T$. Για έναν συμμετρικό θετικά-ορισμένο πίνακα, C , οι πίνακες $U = V$ περιέχουν τα ιδιοδιανύσματα και Σ 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;

Άσκηση 2.6 (Ενισχυτική Μάθηση)

Θεωρήστε το πρόβλημα μετακίνησης με ντετερμινιστικό τρόπο ενός ρομπότ στον διδιάστατο κόσμο της εικόνας. Σε όσες μεταβάσεις δεν αναγράφεται ανταμοιβή (r) θεωρήστε μηδενική τιμή ανταμοιβής. Επίσης, θεωρήστε συντελεστή έκπτωσης $\gamma = 0.8$.



- (α) Για κάθε κατάσταση s υπολογίστε την τιμή $v^*(s)$.
- (β) Σημειώστε στην εικόνα τα βέλτη μεταβάσεων καταστάσεων-ενεργειών που αντιστοιχούν σε βέλτιστη πολιτική. Εάν υπάρχει ισοπαλία, να το αναφέρετε.
- (γ) Πόσες πλήρεις επαναλήψεις ως προς την αξία είναι επαρκείς για να είναι εγγυημένη η εύρεση της βέλτιστης πολιτικής; Θεωρήστε ότι οι αρχικές τιμές αξίας είναι μηδενικές και ότι οι καταστάσεις λαμβάνονται με τυχαία σειρά σε κάθε επανάληψη.
- (δ) Είναι δυνατό να μεταβληθούν οι ανταμοιβές, ώστε η v^* να αλλάξει, αλλά η βέλτιστη πολιτική π^* να παραμείνει αμετάβλητη; Εάν ναι, δώστε ένα παράδειγμα και περιγράψτε την αλλαγή που θα συμβεί στην v^* . Εάν όχι, εξηγήστε σύντομα γιατί είναι αδύνατο.