

Υπολογιστική Εργασία Αναγνώρισης Προτύπων
Ακαδημαϊκό Έτος 2022 – 2023
Γ. Τσιχριντζής Δ. Σωτηρόπουλος

Θέμα: Αλγόριθμοι Σύστασης με Χρήση Τεχνητών Νευρωνικών Δικτύων και Τεχνικών Ομαδοποίησης Δεδομένων

Η συγκεκριμένη εργασία έχει ως στόχο την ανάπτυξη αλγορίθμων σύστασης ταινιών κάνοντας χρήση τεχνητών νευρωνικών δικτύων και τεχνικών συσταδοποίησης. Το σύνολο των δεδομένων επάνω στο οποίο θα εργαστείτε μπορείτε να το κατεβάσετε από την διεύθυνση <https://ieee-dataport.org/open-access/imdb-movie-reviews-dataset> και αφορά σε ένα σύνολο χρηστών $U = \{U_1, U_2, \dots, U_N\}$ οι οποίοι εκφράζουν προτιμήσεις επάνω σε ένα σύνολο αντικειμένων-ταινιών $I = \{I_1, I_2, \dots, I_N\}$. Ο χρήστης $u \in U$ αποδίδει στην ταινία $i \in I$ τον βαθμό προτίμησης $R(u, i) \in R_o = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\} \cup \{0\}$, όπου η τιμή 0 υποδηλώνει την απουσία αξιολόγησης και όχι μηδενικό βαθμό προτίμησης.

Το σύνολο $X = \{x_1, x_2, \dots, x_T\}$ των διαθέσιμων δεδομένων αποτελείται από εγγραφές της μορφής $x_k = [u_k, i_k, r_k, t_k]$ όπου $u_k \in U$ και $i_k \in I$ με $r_k = R(u_k, i_k) \in R_o$. Η τιμή t_k αντιστοιχεί στην ημερομηνία καταχώρησης της συγκεκριμένης εγγραφής.

Προ-επεξεργασία Δεδομένων:

1. Να βρείτε το σύνολο των μοναδικών χρηστών U και το σύνολο των μοναδικών αντικειμένων I .
2. Θεωρείστε την συνάρτηση $\Phi: U \rightarrow P(I)$ (όπου $P(I)$ το δυναμοσύνολο του I) η οποία $\forall u \in U$ επιστρέφει το σύνολο $\varphi(u) \subset I$ των αντικειμένων που αξιολογήθηκαν από τον χρήστη u . Μπορούμε να γράψουμε, επομένως, ότι $\varphi(u) = \{i \in I: R(u, i) > 0\}$. Να περιορίσετε τα σύνολα των μοναδικών χρηστών U και μοναδικών αντικειμένων I στα αντίστοιχα σύνολα \hat{U} και \hat{I} έτσι ώστε $\forall u \in \hat{U}, R_{min} \leq |\varphi(u)| \leq R_{max}$ όπου R_{min} και R_{max} ο ελάχιστος απαιτούμενος και ο μέγιστος επιτρεπτός αριθμός αξιολογήσεων ανά χρήστη. Θεωρήστε προφανώς ότι $|\hat{U}| = n < N$ και $|\hat{I}| = m < M$.
3. Να δημιουργήσετε και να αναπαραστήσετε γραφικά τα ιστογράμματα συχνοτήτων για το πλήθος αλλά και για το χρονικό εύρος των αξιολογήσεων του κάθε χρήστη.
4. Δημιουργήστε μια εναλλακτική αναπαράσταση του συνόλου των δεδομένων ως ένα σύνολο διανυσμάτων προτιμήσεων $R = \{R_1, R_2, \dots, R_n\}$ με $R_j = R(u_j) \in R_o^m, \forall j \in [n]$. Συγκεκριμένα, μπορούμε να γράψουμε ότι:

$$R_j(k) = \begin{cases} R(u_j, i_k), & R(u_j, i_k) > 0; \\ 0, & R(u_j, i_k) = 0. \end{cases} \quad (1)$$

Αλγόριθμοι Ομαδοποίησης Δεδομένων

1. Να οργανώσετε το περιορισμένο σύνολο των χρηστών \hat{U} σε L συστάδες (clusters) της μορφής $\hat{U} = U_{G_1} \cup U_{G_2} \cup \dots \cup U_{G_L}$ έτσι ώστε $U_{G_a} \cap U_{G_b} = \emptyset, \forall a \neq b$ βασιζόμενοι στην διανυσματική αναπαράσταση των προτιμήσεών τους μέσω του συνόλου R . Η διαδικασία ομαδοποίησης των διανυσμάτων προτιμήσεων των χρηστών $R_j = R(u_j) \in R_o^m, \forall j \in [n]$ επάνω στο περιορισμένο σύνολο των αντικειμένων \hat{I} , μπορεί να πραγματοποιηθεί με κατάλληλη παραμετροποίηση του αλγορίθμου **k-means** μεταβάλλοντας την μετρική που

αποτιμά την απόσταση μεταξύ δύο διανυσμάτων προτιμήσεων για ένα ζεύγος χρηστών u και v ως $\mathbf{dist}(\mathbf{R}_u, \mathbf{R}_v)$. Συγκεκριμένα, μπορείτε να χρησιμοποιήσετε τις παρακάτω μετρικές (*):

a. $\mathbf{dist}_{euclidean}(\mathbf{R}_u, \mathbf{R}_v) = \sqrt{\sum_{k=1}^m |\mathbf{R}_u(k) - \mathbf{R}_v(k)|^2 \lambda_u(k) \lambda_v(k)}$ (2) έτσι ώστε:

$$\lambda_j(k) = \begin{cases} 1, & R(u_j, i_k) > 0; \\ 0, & R(u_j, i_k) \leq 0. \end{cases} \quad (3)$$

Η συνάρτηση $\lambda_j()$ αποτιμά το αν ο χρήστης u_j έχει αξιολογήσει ή όχι το αντικείμενο i_k . Προφανώς, μέσω της συγκεκριμένης συνάρτησης μπορούμε αν υπολογίσουμε την τιμή του $|\varphi(u)|$ που αποτιμά το πλήθος των αξιολογήσεων που παρείχε συνολικά ο χρήστης u ως: $|\varphi(u)| = \sum_{k=1}^n \lambda_u(k)$.

b. $\mathbf{dist}_{cosine}(\mathbf{R}_u, \mathbf{R}_v) = 1 - \left| \frac{\sum_{k=1}^m \mathbf{R}_u(k) \mathbf{R}_v(k) \lambda_u(k) \lambda_v(k)}{\sqrt{\sum_{k=1}^m \mathbf{R}_u(k)^2 \lambda_u(k) \lambda_v(k)} \sqrt{\sum_{k=1}^m \mathbf{R}_v(k)^2 \lambda_u(k) \lambda_v(k)}} \right|$ (4)

- c. Να αναπαραστήσετε γραφικά τις συστάδες των χρηστών που αναγνωρίστηκαν από τον αλγόριθμο k-means για κάθε μία από τις παραπάνω μετρικές για διάφορες τιμές της παραμέτρου L .
- d. Να σχολιάσετε την αποτελεσματικότητα των συγκεκριμένων μετρικών στην αποτίμηση της ομοιότητας μεταξύ ενός ζεύγους διανυσμάτων προτιμήσεων χρηστών \mathbf{R}_u και \mathbf{R}_v .

* Οι μετρικές που αναπαρίστανται στις εξισώσεις (2) και (4) υπολογίζονται για κάθε ζεύγος διανυσμάτων \mathbf{R}_u και \mathbf{R}_v μόνο για το υποσύνολο των αντικειμένων που έχουν αξιολογηθεί και από τους δύο χρήστες u και v μέσω της βοηθητικής συνάρτησης $\lambda()$.

Αλγόριθμοι Παραγωγής Συστάσεων με Χρήση Τεχνητών Νευρωνικών Δικτύων

2. Να δημιουργήσετε μια εναλλακτική οργάνωση του περιορισμένου συνόλου των χρηστών σε L συστάδες $\hat{\mathbf{U}} = \mathbf{U}_{G_1} \cup \mathbf{U}_{G_2} \cup \dots \cup \mathbf{U}_{G_L}$ έτσι ώστε $\mathbf{U}_{G_a} \cap \mathbf{U}_{G_b} \neq \emptyset, \forall a, b \in [L]$ με $a \neq b$ κάνοντας χρήση της παρακάτω μετρικής (**):

$$\mathbf{dist}(u, v) = 1 - \frac{|\varphi(u) \cap \varphi(v)|}{|\varphi(u) \cup \varphi(v)|} \quad (5)$$

** Ο αλγόριθμος ομαδοποίησης για την αντιμετώπιση του συγκεκριμένου ερωτήματος είναι της επιλογής σας. Βασική προϋπόθεση, ωστόσο, είναι να μπορεί να λειτουργήσει επάνω στον τετραγωνικό πίνακα των αποστάσεων μεταξύ των χρηστών που περιγράφει η σχέση (5).

- a) Να εξηγήσετε τι εκφράζει η συγκεκριμένη μετρική και να προσδιορίσετε τα μειονεκτήματά της σε σχέση με τις μετρικές που περιγράφονται στις σχέσεις (2) και (4) υπό το πρίσμα της ιδιαίτερης οργάνωσης των χρηστών που μπορεί να επιφέρει.

- b) Η μετρική που περιγράφεται μέσω της σχέσης (5) μπορεί να χρησιμοποιηθεί προκειμένου να προσδιοριστεί το σύνολο $N_k(u_a) = \{u_a^{(1)}, u_a^{(2)}, \dots, u_a^{(k)}\}$ των k πλησιέστερων γειτόνων ενός χρήστη $u_a \in \mathbf{U}_{G_a}, \forall a \in [L]$. Επομένως, για τον εκάστοτε χρήστη u_a της εκάστοτε συστάδας \mathbf{U}_{G_a} μπορούμε να το διάνυσμα τον προσωπικών του προτιμήσεων R_{u_a} καθώς και τα διανύσματα των k πλησιέστερων γειτόνων του $R_{u_a^{(1)}}, R_{u_a^{(2)}}, \dots, R_{u_a^{(k)}}$ εντός της συστάδας \mathbf{U}_{G_a} . Στόχος του συγκεκριμένου ερωτήματος είναι να αναπτύξετε ένα πολυστρωματικό νευρωνικό δίκτυο για κάθε συστάδα χρηστών \mathbf{U}_{G_a} με $a \in [L]$ το οποίο θα προσεγγίζει τις αξιολογήσεις του κάθε χρήστη εντός αυτής μέσω των αξιολογήσεων των k πλησιέστερων γειτόνων του μέσω μιας συνάρτησης της μορφής:

$$R_{u_a} = f_a(R_{u_a^{(1)}}, R_{u_a^{(2)}}, \dots, R_{u_a^{(k)}}), \forall u_a \in \mathbf{U}_{G_a}, \forall a \in [L] \quad (6)$$

- c) Το σύνολο των χρηστών της κάθε συστάδας μπορεί να διαμεριστεί περεταίρω σε ένα υποσύνολο χρηστών για εκπαίδευση του νευρωνικού δικτύου $\mathbf{U}_{G_a}^{train}$ και σε ένα υποσύνολο χρηστών για τον έλεγχο της επίδοσης του νευρωνικού δικτύου $\mathbf{U}_{G_a}^{test}$ έτσι ώστε $\mathbf{U}_{G_a} = \mathbf{U}_{G_a}^{train} \cup \mathbf{U}_{G_a}^{test}$ με $\mathbf{U}_{G_a}^{train} \cap \mathbf{U}_{G_a}^{test} = \emptyset$. Η καταλληλότερη διαμόρφωση των δεδομένων για την εκπαίδευση των νευρωνικών δικτύων αυτού του ερωτήματος θα μπορούσε να αναπαρασταθεί ως εξής:

Έστω ότι το σύνολο των χρηστών που μετέχουν στην συστάδα \mathbf{U}_{G_a} δίνεται από την σχέση: $\mathbf{U}_{G_a} = \{u_{a,1}, u_{a,2}, \dots, u_{a,n_a}\}$ με $n_a = |\mathbf{U}_{G_a}|$. Τότε το σύνολο των διανυσμάτων χαρακτηριστικών και το αντίστοιχο σύνολο των ετικετών θα μπορούσε να οργανωθεί σε έναν πίνακα της μορφής:

$$\begin{bmatrix} R_{u_{a,1}}^{(1)} & R_{u_{a,1}}^{(2)} & \dots & R_{u_{a,1}}^{(k)} \\ R_{u_{a,2}}^{(1)} & R_{u_{a,2}}^{(2)} & \dots & R_{u_{a,2}}^{(k)} \\ \vdots & \vdots & \vdots & \vdots \\ R_{u_{a,n_a}}^{(1)} & R_{u_{a,n_a}}^{(2)} & \dots & R_{u_{a,n_a}}^{(k)} \end{bmatrix} \text{ και } \begin{bmatrix} R_{u_{a,1}} \\ R_{u_{a,2}} \\ \vdots \\ R_{u_{a,n_a}} \end{bmatrix}$$

- d) Η προσεγγιστική ακρίβεια των παραπάνω νευρωνικών δικτύων μπορεί να μετρηθεί μέσω της μετρικής του μέσου απόλυτου σφάλματος ανάμεσα στις πραγματικές και τις εκτιμώμενες αξιολογήσεις των χρηστών. Να παρουσιάσετε πίνακες των αποτελεσμάτων σας τόσο για την ακρίβεια εκπαίδευσης όσο και για την ακρίβεια ελέγχου για κάθε συστάδα χρηστών.

Μπορείτε να εργαστείτε σε ομάδες 3 φοιτητών το πολύ. Η υλοποίηση της εργασίας μπορεί να πραγματοποιηθεί σε MATLAB ή Python. Ο κώδικας της εργασίας σας θα πρέπει να συνοδεύεται από κείμενο αναλυτικής τεκμηρίωσης.

Καλή Επιτυχία!