



ΑΝΑΛΥΤΙΚΗ ΔΕΔΟΜΕΝΩΝ

Απαλλακτική Εργασία

(Ομάδες των 1-3 ατόμων)

Ημερομηνία παράδοσης: *Ημερομηνία εξέτασης του μαθήματος, 23:59μμ*

Σκοπός της εργασίας είναι η εξοικείωση με ένα πραγματικό σύνολο δεδομένων και η εφαρμογή τεχνικών Αναλυτικής Δεδομένων πάνω σε αυτό. Για τον σκοπό αυτό θα επιλέξετε το παρακάτω σύνολο δεδομένων:

- Bank Marketing Data Set (<https://archive.ics.uci.edu/ml/datasets/bank+marketing>). Αποτελείται από περίπου 45 χιλιάδες εγγραφές, και περιέχει πληροφορίες σχετικά με εκστρατείες άμεσου μάρκετινγκ (τηλεφωνικές κλήσεις) ενός Πορτογαλικού τραπεζικού ιδρύματος.

Βήμα 1: Προπαρασκευή δεδομένων (Data preprocessing)

Το πρώτο είναι η εξοικείωση του ερευνητή με τα δεδομένα του. Αφού κατεβάσετε το παραπάνω σύνολο δεδομένων, προχωρήστε σε όποια προπαρασκευαστική εργασία (επιλογή, οπτικοποίηση, καθαρισμό, μετασχηματισμό, δειγματοληψία, κλπ.) θεωρείτε απαραίτητη ώστε: α) να «καθαρίσετε» τα δεδομένα από ελλειπείς ή εσφαλμένες τιμές, εάν υπάρχουν (π.χ., συμπλήρωση κενών πεδίων, απαλοιφή ακραίων τιμών), β) να κανονικοποιήσετε – διακριτοποιήσετε τα δεδομένα (π.χ. για αντιμετώπιση των συνεχών πεδίων τιμών), γ) να μειώσετε τον όγκο των δεδομένων (π.χ. μείωση διαστάσεων). Επίσης θα πρέπει να κάνετε μια απλή στατιστική ανάλυση, σε μορφή ιστογραμμάτων, box plots κλπ., των πιο βασικών (κατά τη γνώμη σας) χαρακτηριστικών του συνόλου δεδομένων.

Βήμα 2: Συσταδοποίηση (Clustering)

Έχοντας εξοικειωθεί με το σύνολο δεδομένων, το επόμενο βήμα της πειραματικής σας διαδικασίας είναι η χρήση τεχνικών συσταδοποίησης, προκειμένου να ανακαλύψετε ιδιότητες του συνόλου και πρότυπα που δεν είναι προφανή με μια απλή στατιστική ανάλυση. Σε αυτό το στάδιο, σημαντικό ρόλο παίζει η μοντελοποίηση του προβλήματος (τι ακριβώς ψάχνετε να

εντοπίσετε). Διαδικαστικά, αφού επιλέξετε (α) τα χαρακτηριστικά του συνόλου δεδομένων τα οποία θα αποφασίσετε να εξετάσετε και (β) μια κατάλληλη μετρική απόστασης/ομοιότητας, χρησιμοποιήστε μέσω του εργαλείου Scikit-Learn δύο διαφορετικές τεχνικές συσταδοποίησης (K-means, DBSCAN), συζητήστε τα αποτελέσματα και την επίπτωση των παραμέτρων των μεθόδων σε αυτά, και συγκρίνετέ τα ως προς την ποιότητα/αποτελεσματικότητα της συσταδοποίησης (π.χ. scatter plots, clustering metrics).

Βήμα 3: Ταξινόμηση (Classification)

Τελευταίο βήμα της πειραματικής σας διαδικασίας είναι η χρήση μοντέλων ταξινόμησης με στόχο την ανάθεση ενός αντικειμένου σε προκαθορισμένες κατηγορίες (κλάσεις). Όπως πριν, και σε αυτό το στάδιο, σημαντικό ρόλο παίζει η μοντελοποίηση του προβλήματος (τι ακριβώς ψάχνετε να εντοπίσετε). Διαδικαστικά, αφού μετασχηματίσετε κατάλληλα το dataset στη μορφή (<Feature(s)>, <Label(s)>), δημιουργήστε μέσω του API Scikit-Learn δύο ταξινομητές (π.χ., Bayesian, LS-SVM, Neural Networks) και - όπως και στο προηγούμενο βήμα - συγκρίνετε τις επιδόσεις τους (π.χ., υπολογίζοντας confusion matrix, ROC-AUC curve). Παρόμοια, κατασκευάστε δυο νευρωνικά δίκτυα μέσω του API TensorFlow/Keras, και δοκιμάστε να εκπαιδεύσετε το ένα ως έχει, και το άλλο μέσω transfer learning και συγκρίνετέ τα με τα προηγούμενα μοντέλα. Τι παρατηρείτε; Βελτιώθηκε η επίδοση των μοντέλων και γιατί;

Βήμα 4: Σύνοψη

Λαμβάνοντας υπόψη τα αποτελέσματα των προηγούμενων βημάτων, καταγράψτε τις παρατηρήσεις σας και 2-3 βασικά συμπεράσματα (“take-home messages”) αναφορικά με το σύνολο δεδομένων με κατάλληλη απεικόνιση/οπτικοποίηση – με άλλα λόγια, “πείτε μια ιστορία” με τα δεδομένα σας (“data story telling”).

Το τελικό παραδοτέο θα αποτελείται από ένα αρχείο zip, το οποίο θα υποβληθεί ηλεκτρονικά στην ενότητα «Εργασίες» του μαθήματος στο e-class και θα περιέχει τα εξής:

- Τεχνική αναφορά (report) με αναλυτική περιγραφή των προσεγγίσεων που ακολουθήσατε σε καθένα από τα βήματα (π.χ. παράμετροι αλγορίθμων, προπαρασκευή δεδομένων, κλπ.) και ερμηνεία των αποτελεσμάτων που προέκυψαν. Στην τεκμηρίωση πρέπει να αναγράφονται τα στοιχεία των φοιτητών της ομάδας.
- Τα αρχεία πηγαίου κώδικα (source code) και τυχόν συμπληρωματικά αρχεία (που είναι απαραίτητα για την εκτέλεση του κώδικα), καθώς και τα αποτελέσματα που παρήχθησαν (π.χ. plots).

Ζητήματα δεοντολογίας

(α) σε περίπτωση αντιγραφής οι εμπλεκόμενες εργασίες μηδενίζονται, (β) σε περίπτωση αμφιβολίας για το κατά πόσο η ομάδα που αναγράφεται ήταν εκείνη που ανέπτυξε την εργασία, ενδέχεται να της ζητηθεί να την παρουσιάσει για τυχόν διευκρινίσεις.