

1. Ανοίγουμε το αρχείο του glass dataset μας και ορίζουμε τις στείλες βάση των στοιχείων που θέλουμε
2. Κάνουμε κανονικοποίηση των τιμών μας και βάζουμε τις στείλες των δεδομένων κάνω και αφαίρεση των αυτλιερς καθώς και αφαίρω τον fe που έχει αδείες γραμμές
3. Κάνω την ελμπουου ανάλυση που μας ζητάει για να βρώ τα απαραίτητα κλάστερς τα οποία είναι 4 για τα δεδομένα που έχουν υποστεί τον καθαρισμό
4. Το επιβεβαιώνω και απο το σιολουέτ ανάλυσις οτι είναι το $K=4$ άρα ο ιδανικός αριθμος κλάστερς μάς είναι 4 για το σετ δεδομένων όπως έχω επεξεργαστεί
5. Εφαρμόζω το Kmeans analysis model για 4 κλάστερ όπως έχουμε βρεί και θέτοντας την παράμετρο random_state σε 0 για να διασφαλιστεί ότι τα αποτελέσματα της ομαδοποίησης είναι αναπαραγώγιμα
6. Κάνω σκάττερ πλότερ απο το dataframe μας δύο στείλες και απο κάτω όλες τις στείλες με το ρεφλέτιβ ιντεξ
7. Χρησιμοποιώ το inconsistency index που είναι ένα μέτρο της αστάθειας της ομαδοποίησης κατά τη χρήση ιεραρχικής ομαδοποίησης με τη μέθοδο μονής σύνδεσης. Ορίζεται ως η μέγιστη διαφορά στις αποστάσεις μεταξύ ζευγών παρατηρήσεων στην ίδια συστάδα σε διαφορετικά επίπεδα της ιεραρχίας. Ένας υψηλότερος συντελεστής inconsistency υποδηλώνει ότι η ομαδοποίηση είναι πιο ασταθής και μπορεί να είναι λιγότερο αξιόπιστη.
8. Βρίσκω και για το σετ δεδομένων μου και για τα κέντρα των κλάστερ οτι η απλή σύνδεση είναι καλύτερα αποτελέσματα ως μέση τιμή του συντελεστή την χαμηλότερη απο όλα
9. Ολές οι μέθοδοι έχουν πολύ μικρό αριθμο inconsistency που συνεπώς σημαίνει οτι όλες παράγουν σωστά κάπως αποτελέσματα με αρκετή ακρίβεια
10. Για να προσδιοριστεί ποιο μοντέλο ομαδοποίησης θα παράγει ένα καλύτερο δενδρογράμμο, μπορούμε να συγκρίνουμε τους μέσους συντελεστές inconsistency για κάθε μοντέλο και να επιλέξουμε το μοντέλο με τον χαμηλότερο μέσο συντελεστή inconsistency. Επομένως, ένα μοντέλο με χαμηλότερο μέσο συντελεστή inconsistency είναι πιθανό να παράγει ένα δενδρογράμμο
11. Κάνοντας την ανάλυση των δενδρογραμμάτων βλέπουμε οτι όλα τα δενδρογράμματα παράγουν κάπως καλά δεδομένα όμως παρατηρούμε οτι το δενδρογράμμο απλής single σύνδεσης έχει μικρότερη απόστασή στο κοψιμό y στις τελευταίες κορυφές μας έχουν διακριτό σχήμα και είναι σχετικά συμπαγείς, παρόμοια αποτελέσματα βλέπουμε και στο average που έχει το αμέσως πιο μικρό inconsistency index μετά το single
12. Συμπαιράινουμε οτι οταν τα δεδομένα καθαριστούν παράγουν καλύτερα σκάττερ πλότερ και έχουμε καλύτερα αποτελέσματα σε όλες τις μετρήσεις και οτι τα κλάστερς μάς είναι 4 και οτι έχουμε καλύτερα αποτελέσματα στα single και average δενδρογράμματα συνεπώς είναι καλύτερες μέθοδοι για τα σέτ δεδομένων μας. αυτο επίσης μπορεί να υποδικνύει οτι τα δεδομένα έχουν ιεραρχική δομή, όπου οι κλάστερς σχηματίζονται από τη συγχώνευση μεμονωμένων δειγμάτων ή μικρών ομάδων δειγμάτων. Αυτό οφείλεται στο γεγονός οτι η μονής σύνδεσης είναι γνωστό οτι σχηματίζει συστάδες που είναι μεγάλες και λεπτές, με τα δείγματα να συγχωνεύονται με διαδοχικό τρόπο.