

Υπολογιστική σύγκριση αλγορίθμων συσταδοποίησης



Διπλωματική εργασία του Μαντιού Ιωάννη
mai20034

Περιεχόμενα

1. Case study : τί προσπαθούμε να πετύχουμε
2. Αλγόριθμοι clustering
3. Αλγόριθμοι classification
4. Unlabeled Datasets
5. Labeled Datasets
6. Συμπεράσματα

1. Case study : τί προσπαθούμε να πετύχουμε



VS

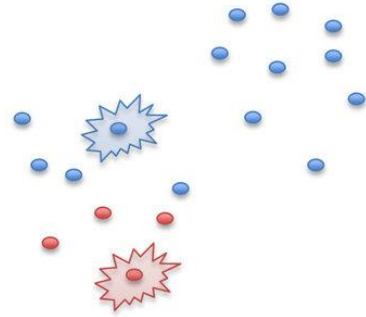


2. Αλγόριθμοι clustering

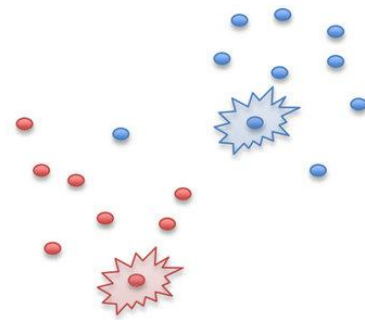
2.1 K-Means

- Κάθε συστάδα σχετίζεται με ένα κεντρικό σημείο (centroid)
- Κάθε σημείο ανατίθεται στην συστάδα με το κοντινότερο κεντρικό σημείο (π.χ. τετραγωνική ευκλείδεια απόσταση)
- Ξαναυπολογίζονται τα κεντρικά σημεία
- Επιλέγουμε εκ των προτέρων τον αριθμό των ομάδων (K)

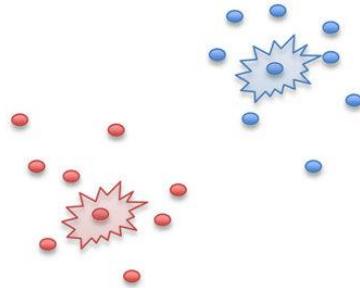
Initial Seeding



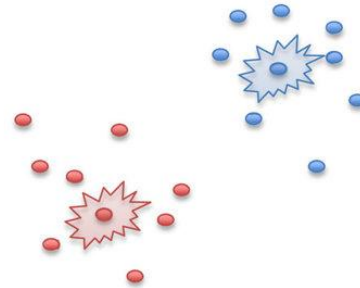
After Round 1



After Round 2



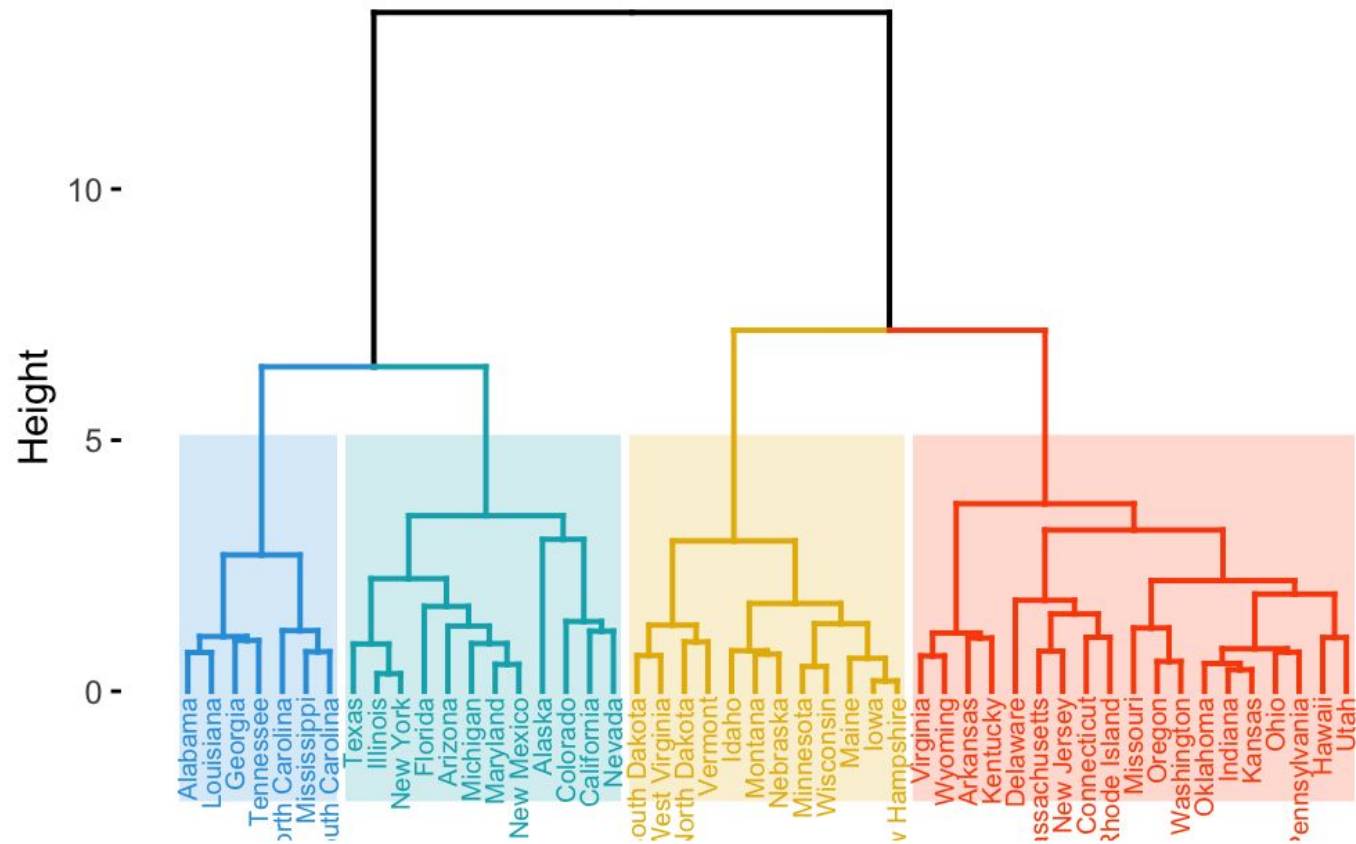
Final



2.2 Agglomerative

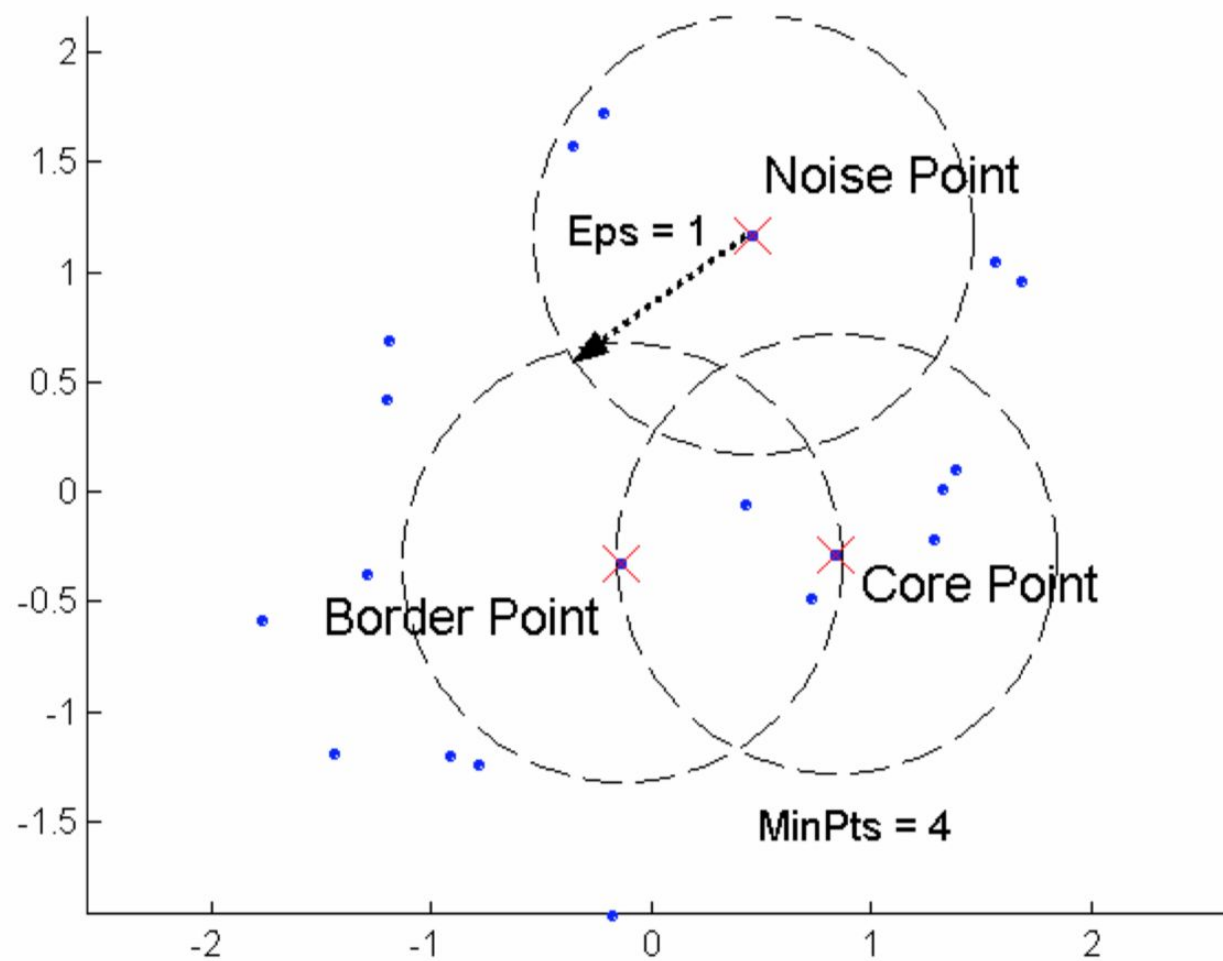
- Δημιουργεί ομάδες με ιεραρχικό τρόπο από κάτω προς τα πάνω
- Στην αρχή όλες οι παρατηρήσεις αποτελούν ατομικές ομάδες οι οποίες συνενώνονται διαδοχικά σε μεγαλύτερες ομάδες.
- Επιλογή και ένωση ομάδων με βάση κάποιο κριτήριο ελάχιστης απόστασης και συγχώνευση των παρατηρήσεων.
- Δενδρόγραμμα

Cluster Dendrogram



2.3 DBSCAN

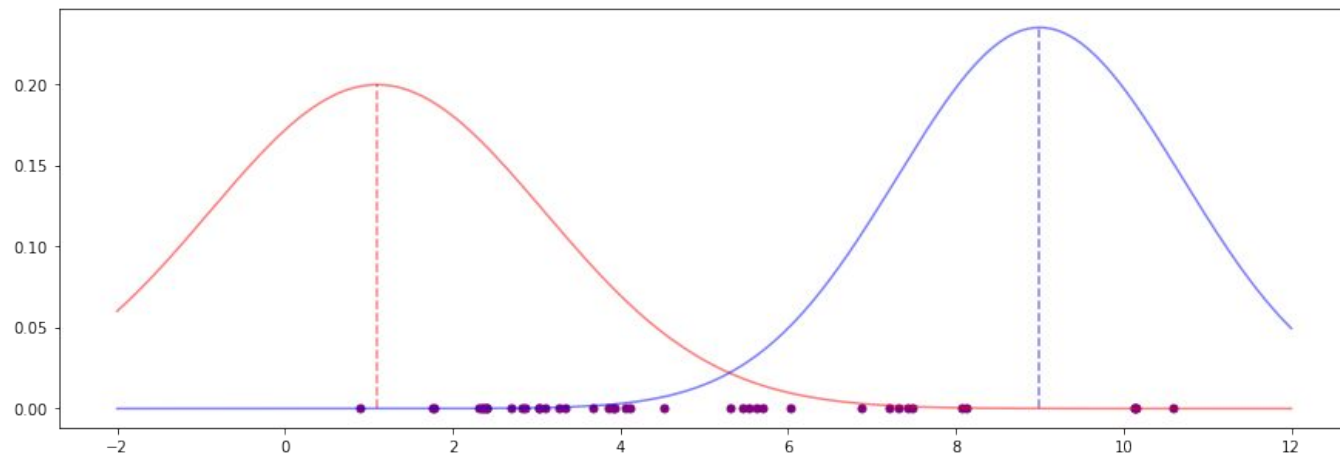
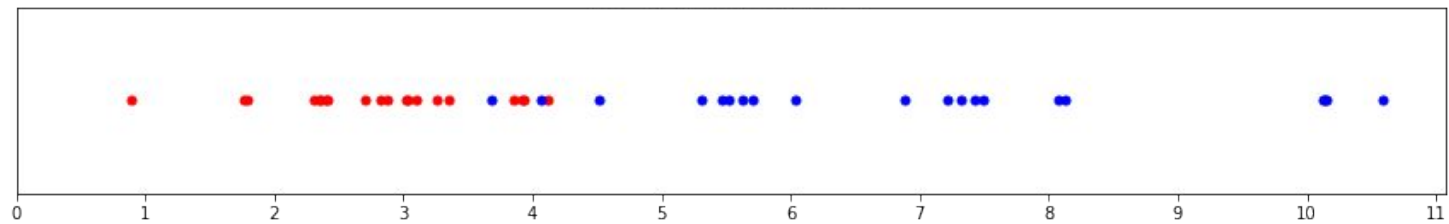
- Βασίζεται στην πυκνότητα των σημείων μέσα σε μια προκαθορισμένη ακτίνα
- Καθορισμός minPts και Eps
- Σημεία: Βασικά, Οριακά, Θορύβου

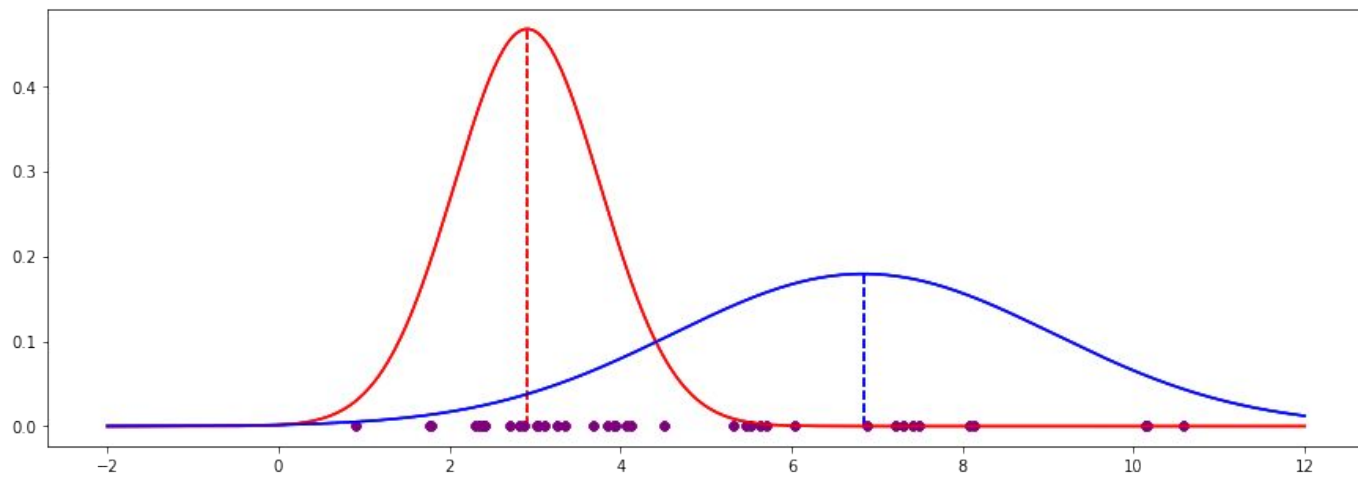
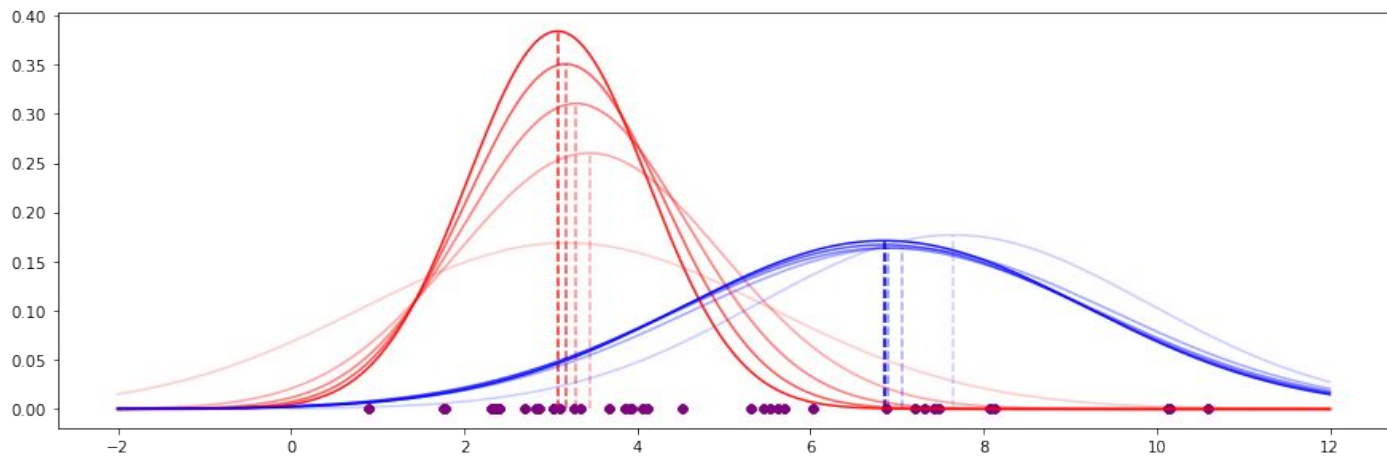


2.4 Expectation-Maximization

- Βασίζεται σε ένα μοντέλο πιθανοτήτων που συνδυάζει διαφορετικές Gaussian κατανομές
- Επιλέγει τυχαίες αρχικές τιμές για το σύνολο Θ των παραμέτρων που προσδιορίζουν την κάθε κατανομή/ ομάδα (π.χ. μέση τιμή, τυπική απόκλιση κτλ.)
- Όσο αλλάζουν οι τιμές των παραμέτρων του Θ : υπολογίζει για κάθε σημείο αν ανήκει σε μια κατανομή (expectation) και για τις πιθανότητες αυτές υπολογίζει το νέο Θ (maximization)

Data (no hidden variables)

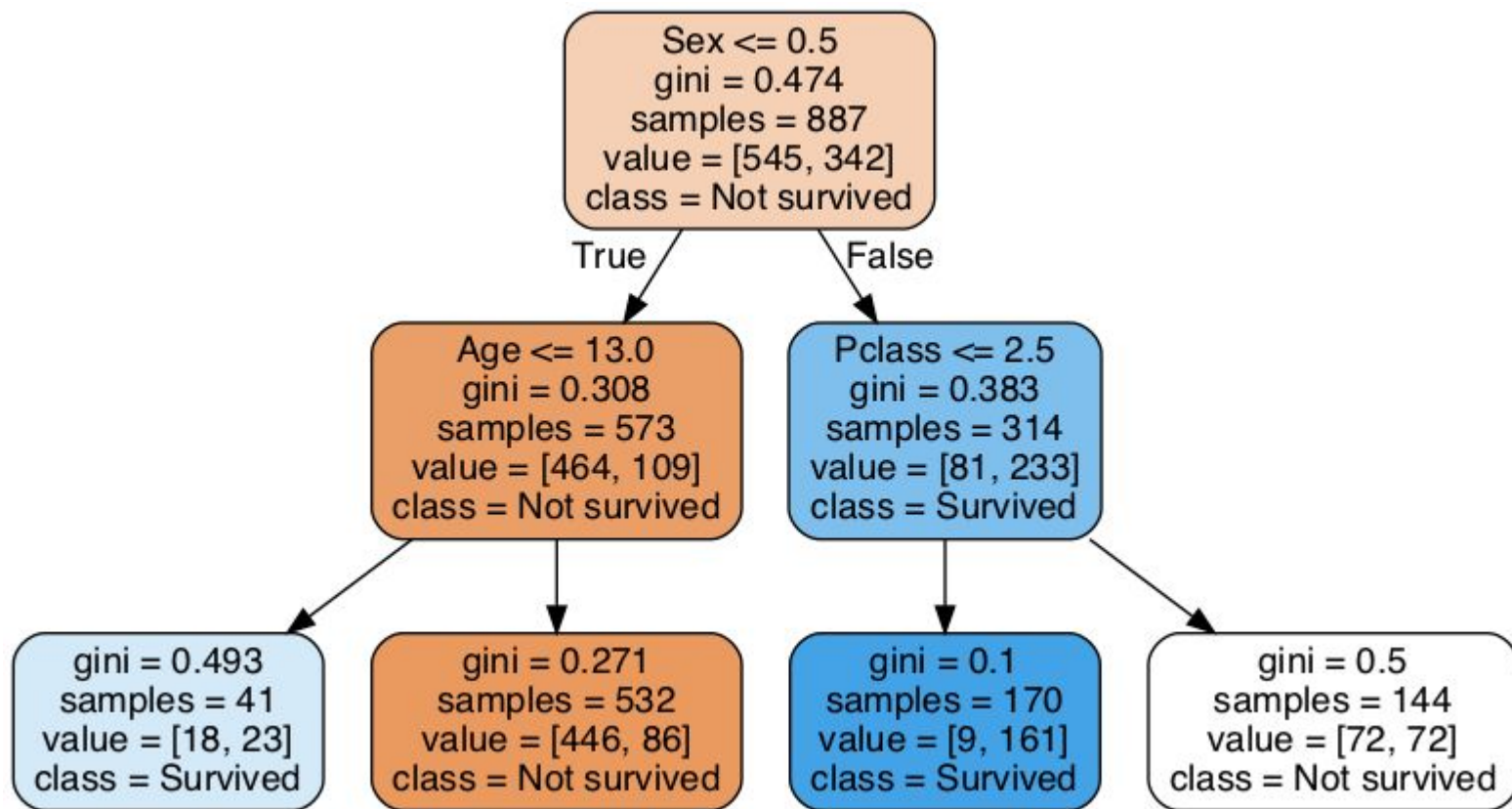




3. Αλγόριθμοι classification

3.1 Decision Trees (CART)

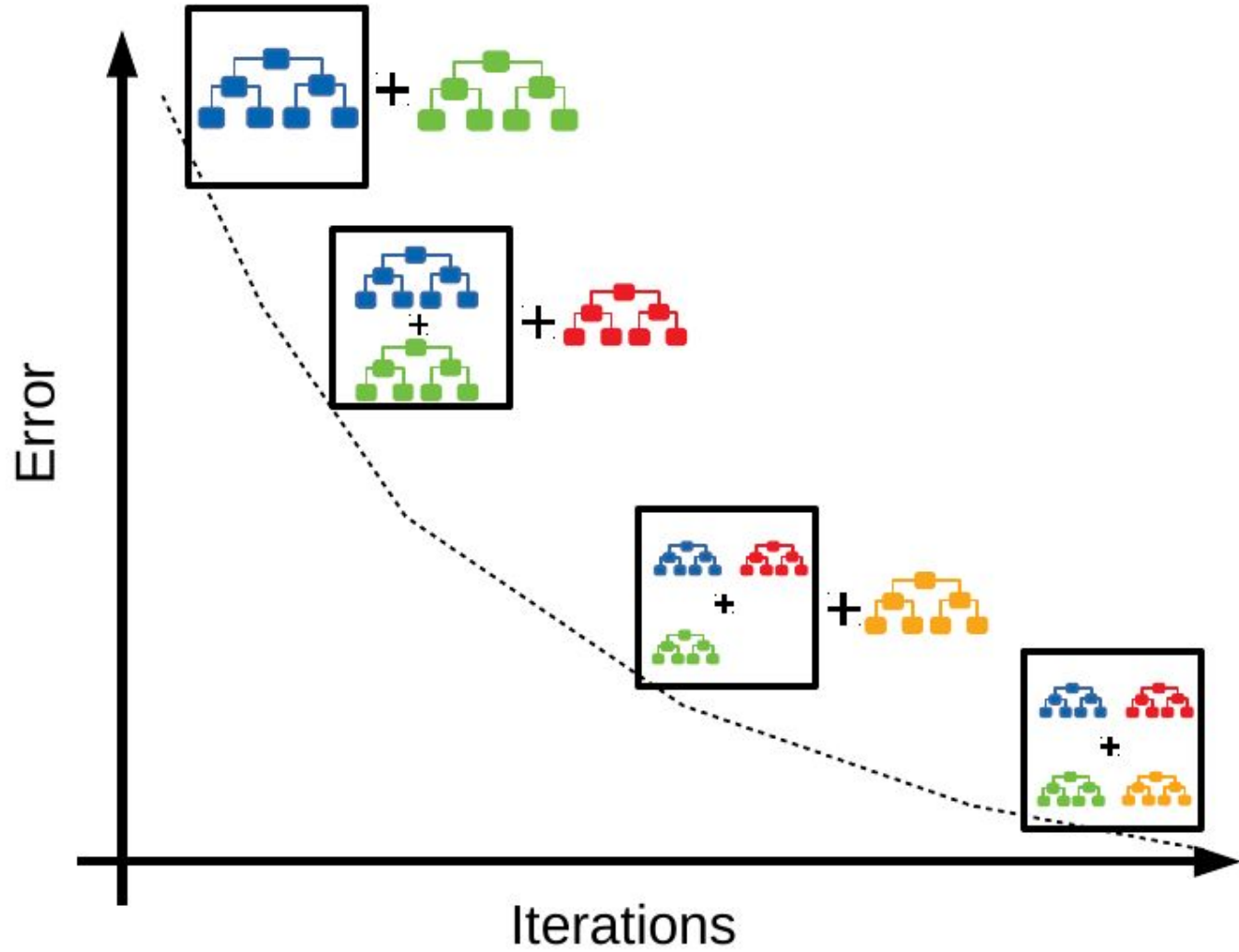
- Επιλέγει το χαρακτηριστικό που πετυχαίνει τον καλύτερο διαχωρισμό μεταξύ των κατηγοριών.
- Χωρίζει τα δεδομένα σε υποσύνολα με βάση της τιμές του χαρακτηριστικού αυτού.
- Για κάθε υποσύνολο που περιέχει περισσότερες από μία κατηγορίες, επαναλαμβάνει τη διαδικασία.
- Σταματάει εφόσον δεν υπάρχουν υποσύνολα που περιέχουν περισσότερες από μία κατηγορίες ή έχουν χρησιμοποιηθεί όλα τα χαρακτηριστικά.



3.2 Gradient Boosting

- Η βασική ιδέα του Boosting είναι ο συνδυασμός αλγορίθμων των οποίων το σφάλμα είναι ελαφρώς καλύτερο από την τυχαία επιλογή.
- Ο αλγόριθμος Gradient Boosting προσθέτει νέα μοντέλα με στόχο την διόρθωση σφαλμάτων που έγιναν από τα υπάρχοντα μοντέλα. Τα μοντέλα προστίθενται διαδοχικά μέχρις ότου δεν μπορούν να γίνουν περαιτέρω βελτιώσεις.

- Ένα ενισχυμένο δέντρο χρησιμοποιεί μικρότερα δέντρα που εξηγούν μόνο ένα κομμάτι των χαρακτηριστικών ανά δέντρο.
- Ο αλγόριθμος δημιουργεί ένα νέο δέντρο που αναλύει τα προηγούμενα δένδρα. Στη συνέχεια δημιουργεί ένα νέο δέντρο που επιχειρεί να διορθώσει τα σφάλματα στα προηγούμενα κ.ο.κ.



3.3 Extreme Gradient Boosting

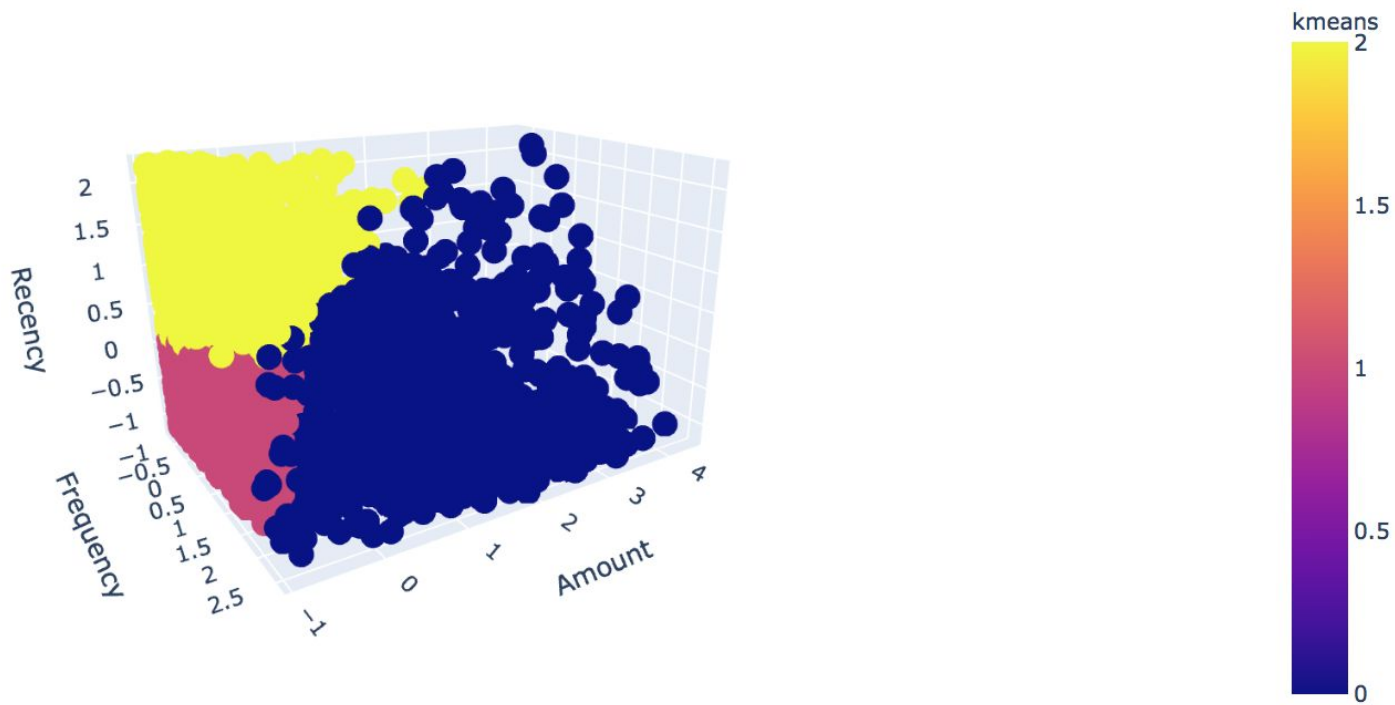
- Η state of the art έκδοση του gradient boosting.
- Ο XGBoost επικεντρώνεται στην υπολογιστική ταχύτητα και την απόδοση του μοντέλου.
- Καινοτομίες του XGBoost: αυτόματος χειρισμός των δεδομένων που λείπουν και δομή block για τον παραλληλισμό της δομής των δέντρων

4. Unlabeled Datasets

4.1 Online Retail II Dataset

- Συναλλαγές 2 ετών ενός online καταστήματος
- Χαρακτηριστικά του δείγματος: Αριθμός τιμολογίου, κωδικός προϊόντος, περιγραφή, ποσότητα, ημερομηνία αγοράς, τιμή μονάδας, κωδικός πελάτη, χώρα πελάτη.
- Ανάλυση RFM (Recency, Frequency, Monetary) για την αξιολόγηση κάθε πελάτη.

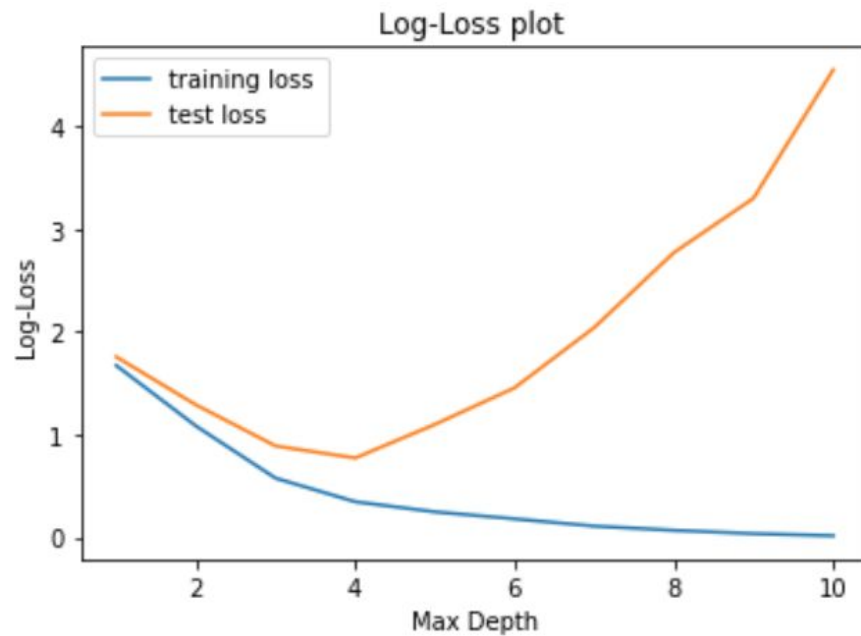
- Μετά από καθαρισμό των κακών δεδομένων και ομαδοποίηση κατά το μοντέλο RFM: πίνακας διαστάσεων $1.000.000 \times 8 \ggg 4.075 \times 3$
- Πριν το clustering: διαγραφή outliers και standardization



Μετά την τεχνητή παραγωγή κλάσεων μέσω clustering:

- Oversampling με αλγόριθμο SMOTE
- Stratified Train/test split: 80%/20%
- Αντιμετώπιση overfitting βάσει log-loss μεταξύ training και test set

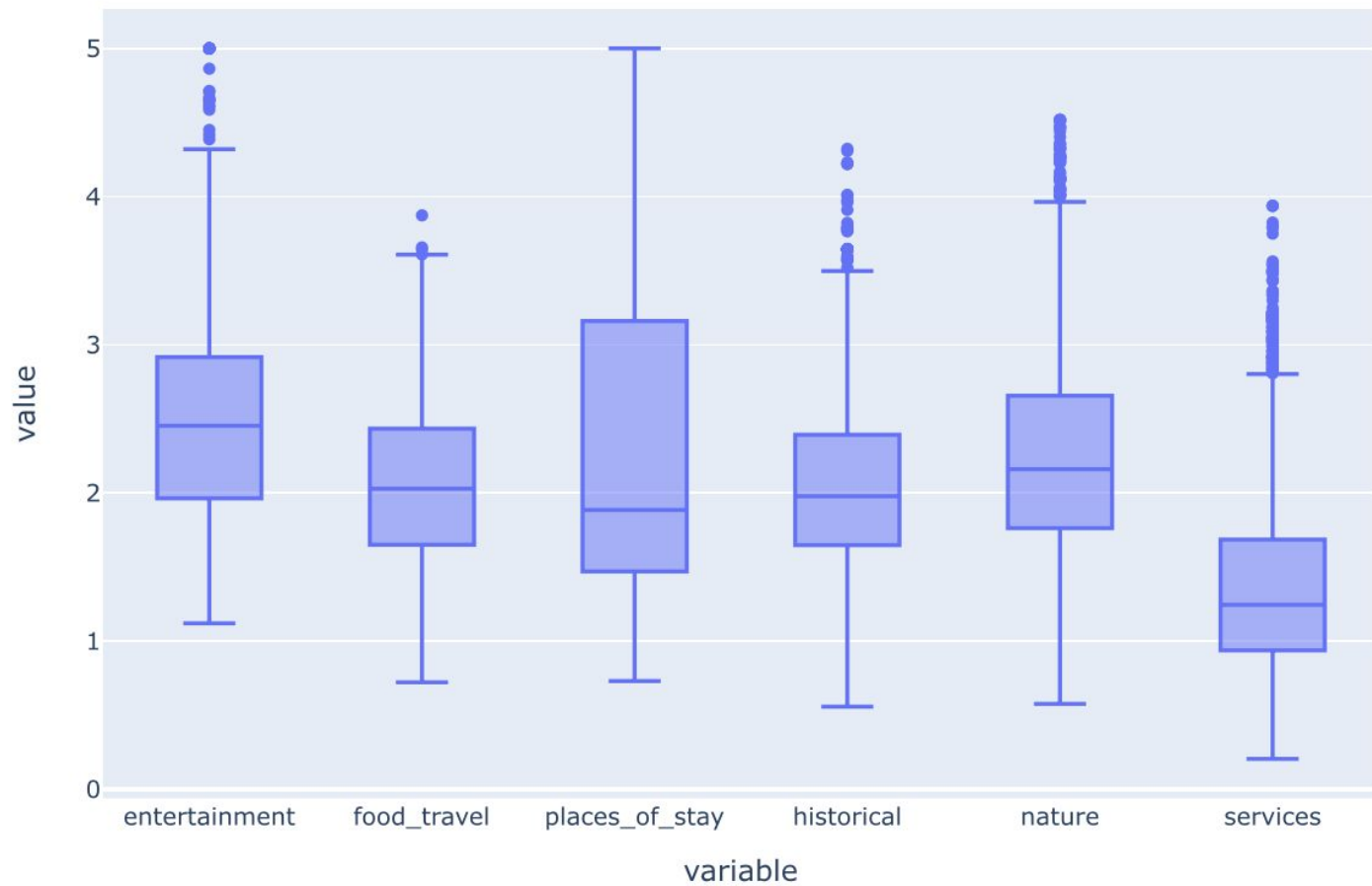
Παράδειγμα log-loss plot



4.2 Travel Review Ratings Dataset

- Κριτικές google σε 24 κατηγορίες χώρων στην Ευρώπη
- Αξιολόγηση με αστέρια από 1 έως 5

- Γκρουπάρισμα ομοειδών τοποθεσιών σε κατηγορίες για καλύτερο clustering (μείωση του dimensionality)
- Δεν χρειάζεται standardization, ίδιο εύρος τιμών μεταξύ των κατηγοριών



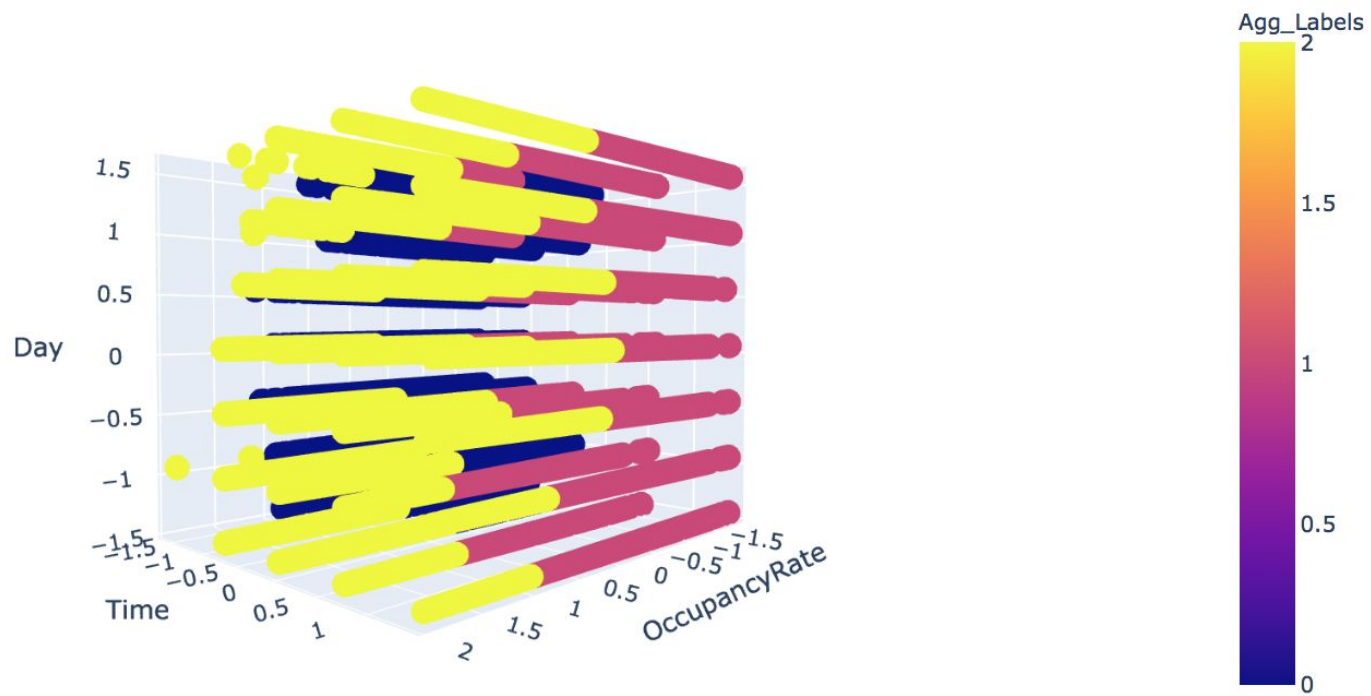
Μετά την τεχνητή παραγωγή κλάσεων μέσω clustering:

- Oversampling με αλγόριθμο SMOTE
- Stratified Train/test split: 80%/20%
- Αντιμετώπιση overfitting βάσει log-loss μεταξύ training και test set

4.3 Birmingham Parking Dataset

- Ημερήσια δεδομένα 2 μηνών από δημοτικά πάρκινγκ στο Birmingham
- Χαρακτηριστικά του δείγματος: κωδικός πάρκινγκ, χωρητικότητα, πληρότητα, στιγμή καταγραφής

- Ομαδοποίηση του χρόνου σε δίωρα και ημέρες
- Δημιουργία στήλης 'Ποσοστό Πληρότητας'
- Τελικές στήλες για clustering : Χωρητικότητα, Ποσοστό Πληρότητας, Ώρα, Ημέρα



Μετά την τεχνητή παραγωγή κλάσεων μέσω clustering:

- Oversampling με αλγόριθμο SMOTE
- Stratified Train/test split: 80%/20%
- Αντιμετώπιση overfitting βάσει log-loss μεταξύ training και test set

5. Labeled Datasets

5.1 Teaching Assistant Evaluation Dataset

- Αξιολογήσεις 151 καθηγητών σε περίοδο 2 χρόνων στο πανεπιστήμιο του Wisconsin
- Βαθμοί αξιολόγησης: low, medium, high
- Χαρακτηριστικά του δείγματος: γλώσσα ομιλίας του διδάσκοντα, διδάσκοντας, μάθημα, εξάμηνο, μέγεθος τάξης, αξιολόγηση

- Δεν χρειάζεται προεπεξεργασία των δεδομένων
- Oversampling με αλγόριθμο SMOTE
- Stratified Train/test split: 80%/20%
- Αντιμετώπιση overfitting βάσει log-loss μεταξύ training και test set

5.2 Car Evaluation Dataset

- Αξιολογήσεις 1728 αυτοκινήτων βάσει 6 χαρακτηριστικών
- Χαρακτηριστικά δείγματος: τιμή αγοράς, κόστος συντήρησης, αριθμός θυρών, χωρητικότητα ατόμων, μέγεθος πορτ μπαγκάζ, εκτιμώμενη ασφάλεια
- Κλάσεις: ακατάλληλο, κατάλληλο, καλό, πολύ καλό

- Μετατροπή κατηγορικών μεταβλητών σε σειριακούς ακεραίους
- Oversampling με αλγόριθμο SMOTE
- Stratified Train/test split: 80%/20%
- Αντιμετώπιση overfitting βάσει log-loss μεταξύ training και test set

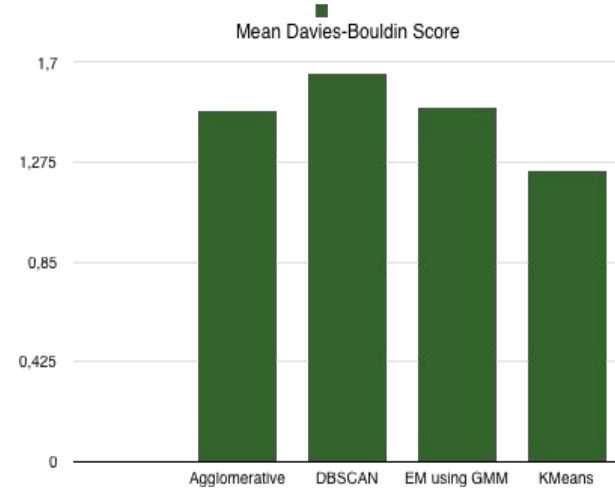
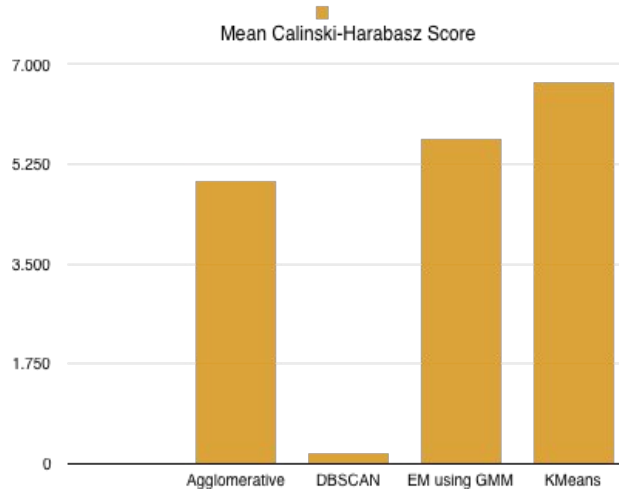
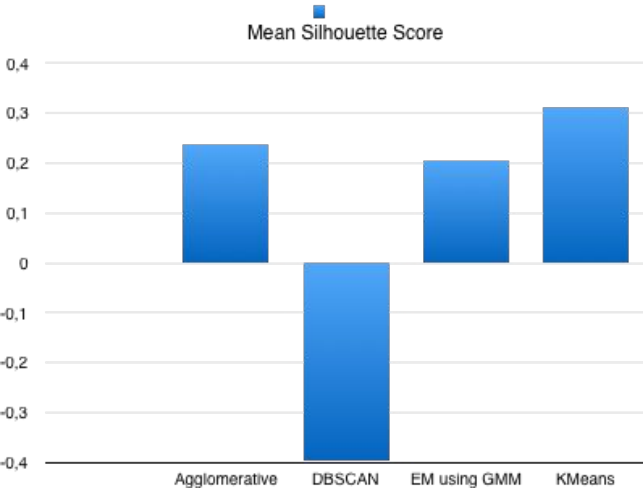
5.3 Cardiotocography Dataset

- Μετρήσεις εμβρυικού καρδιακού ρυθμού σε 2126 βρέφη (καρδιοτοκογράφημα)
- 22 χαρακτηριστικά όπως χτύποι το λεπτό, εμβρυικές κινήσεις το δευτερόλεπτο κ.ο.κ
- 10 κλάσεις που αντιπροσωπεύουν 10 διαφορετικά μοτίβα καρδιακών ρυθμών

- Oversampling με αλγόριθμο SMOTE
- Stratified Train/test split: 80%/20%
- Αντιμετώπιση overfitting βάσει log-loss μεταξύ training και test set

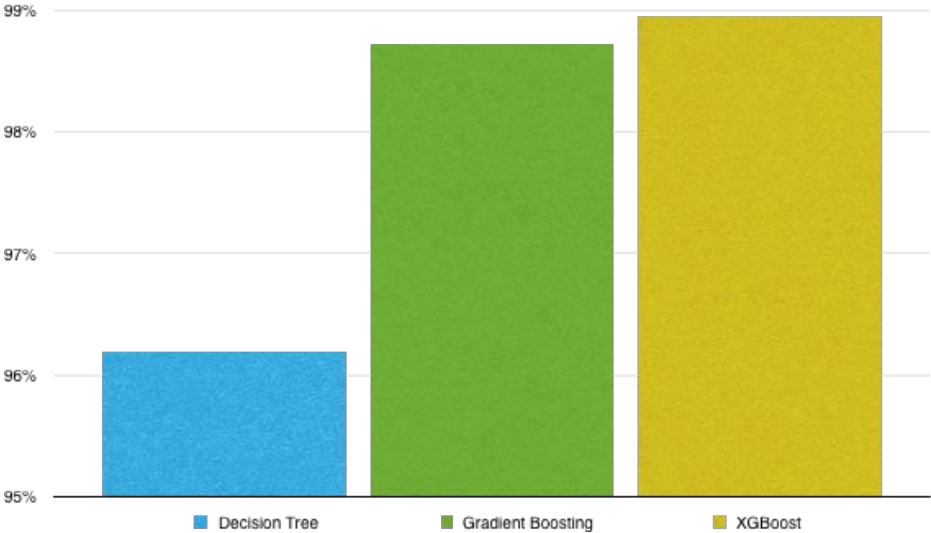
6. Συμπεράσματα

Μέσοι Όροι Clustering Metrics

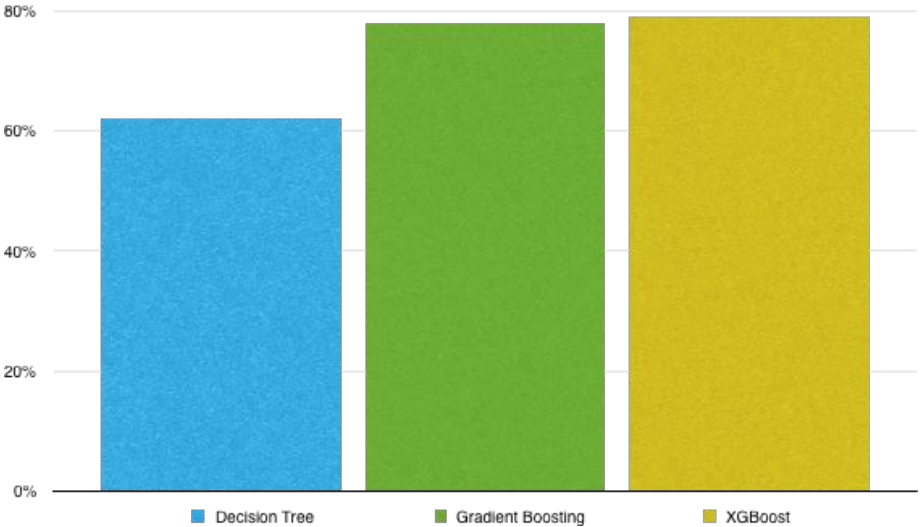


Μέσοι Όροι Accuracy

Clustered Datasets Average Accuracy



Labeled Datasets Average Accuracy



- Οι αλγόριθμοι κατηγοριοποίησης έχουν πολύ υψηλή απόδοση όταν οι κλάσεις έχουν παραχθεί με clustering
- Ένας απλός αλγόριθμος classification σε συνδυασμό με clustering έχει αρκετά υψηλότερη ακρίβεια πρόβλεψης σε σχέση με έναν πολύπλοκο αλγόριθμο classification χωρίς να έχει προηγηθεί clustering.

- Ο αριθμός των παραγόμενων clusters εξαρτάται σε τεράστιο βαθμό από το hyperparameter tuning
- *Η σύγκριση μοντέλων πάνω σε διαφορετικά σεντ δεδομένων χρειάζεται επιφύλαξη λόγω της διαφορετικότητας των δεδομένων από dataset σε dataset

Ευχαριστώ