

# Επιστήμη των Δεδομένων

Δημοκρίτειο Πανεπιστήμιο Θράκης



Αλεξάκης Γεώργιος

58093

Κομνηνός Βασίλειος

58051

Ηλεκτρονικές Μετρήσεις

2 Μαΐου 2023

## Περιεχόμενα

<b>1</b>	<b>Εισαγωγή</b>	<b>3</b>
<b>2</b>	<b>Ταξινόμηση</b>	<b>5</b>
2.1	Εισαγωγή στους ταξινομητές . . . . .	5
2.2	Naive Bayes . . . . .	6
2.3	Support Vector Machine . . . . .	8
2.4	K Nearest Neighbors . . . . .	10
2.5	Decision Trees . . . . .	13
<b>3</b>	<b>Παλινδρόμηση</b>	<b>16</b>
3.1	Γραμμική Παλινδρόμηση . . . . .	16
<b>4</b>	<b>Ομαδοποίηση</b>	<b>19</b>
4.1	Εισαγωγή στην ομαδοποίηση . . . . .	19
4.2	Centroid based clustering . . . . .	20
4.3	Density based clustering . . . . .	21
4.4	Hierarchical clustering . . . . .	22
4.5	Distribution based clustering . . . . .	23
<b>5</b>	<b>Υλοποίηση</b>	<b>24</b>
5.1	CSV Parser . . . . .	24
<b>6</b>	<b>Συμπέρασμα</b>	<b>27</b>
<b>A'</b>	<b>Τύποι για QDA, LDA</b>	<b>30</b>

## Κατάλογος Σχημάτων

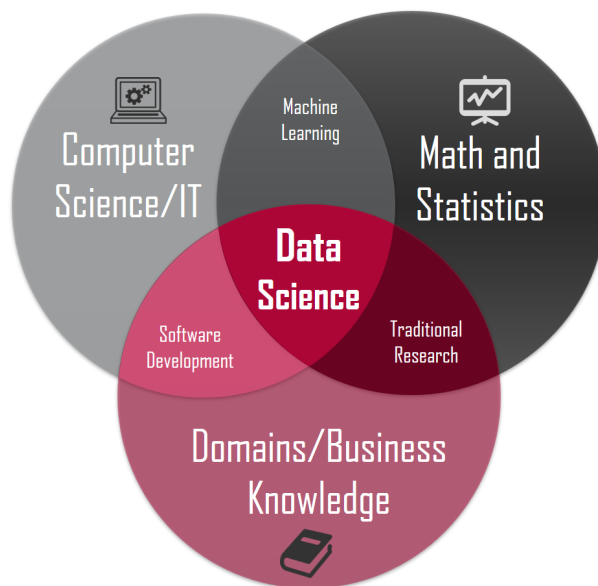
1	Γνώσεις που συνδυάζει η επιστήμη των δεδομένων . . . . .	3
2	Αλγόριθμοι μηχανικής μάθησης . . . . .	4
3	Τύποι ταξινομητών . . . . .	5
4	Παράδειγμα SVM . . . . .	9
5	Επαύξηση διάστασης με SVM . . . . .	9
6	Ταξινόμηση πολλαπλών κλάσεων με SVM . . . . .	10
7	Αλγόριθμος k-NN . . . . .	11
8	Ευκλείδεια απόσταση και Απόσταση Manhattan . . . . .	12
9	Απόσταση Minkowski . . . . .	12
10	Παράδειγμα δέντρου απόφασης . . . . .	13
11	Γραφική αναπαράσταση ομαδοποίησης . . . . .	19
12	Επαναλήψεις του αλγορίθμου k-Means . . . . .	20
13	Αποτελέσματα ομαδοποίησης με DBSCAN vs k-Means . . . . .	21
14	Hierarchical clustering . . . . .	22
15	Γραφική αναπαράσταση Distribution based clustering . . . . .	23

# 1 Εισαγωγή

Η επιστήμη των δεδομένων (Data science) είναι μια επιστήμη που αφορά την εξαγωγή γνώσης από δεδομένα (δομημένα ή αδόμητα) και για να το πετύχει αυτό συνδυάζει γνώσεις από διάφορες άλλες επιστήμες όπως:

- Μαθηματικά
- Στατιστική
- Προγραμματισμός
- Προγνωστική ανάλυση (Predictive analysis)
- Εξόρυξη δεδομένων (Data mining)
- Τεχνητή Νοημοσύνη (Artificial intelligence)
- Μηχανική μάθηση (Machine learning)

Επιπλέον, η επιστήμη των δεδομένων συνδυάζει τα παραπάνω με τη γνώση κάποιου ειδικού τομέα (όπως η ιατρική) με σκοπό να δώσει λύση σε ένα συγκεκριμένο πρόβλημα<sup>[1]</sup>



Σχήμα 1: Γνώσεις που συνδυάζει η επιστήμη των δεδομένων

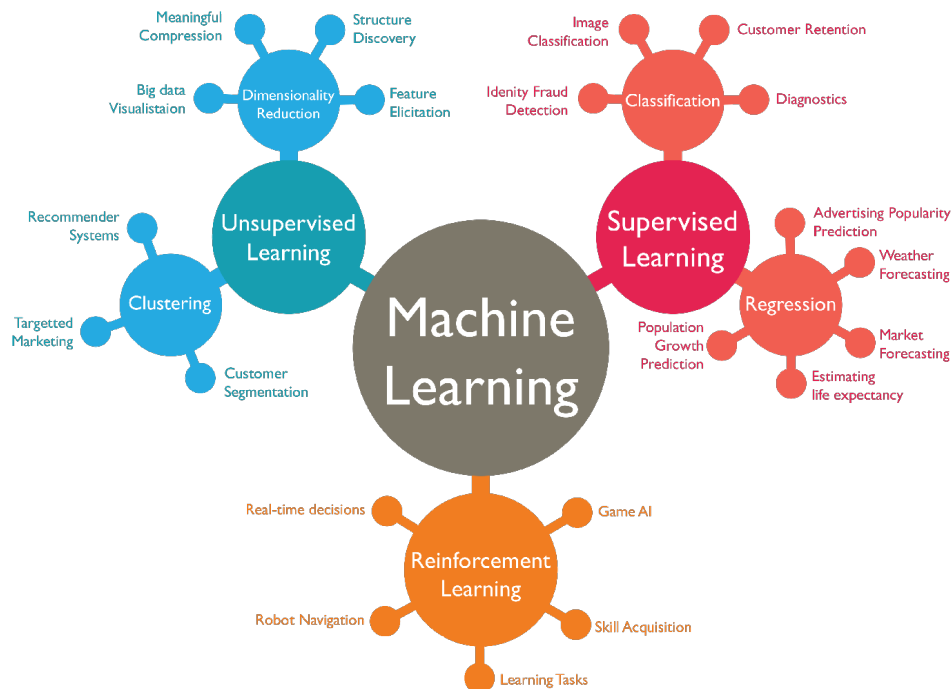
Σε αυτή την εργασία θα ασχοληθούμε κυρίως με την τεχνητή νοημοσύνη και τη μηχανική μάθηση και πιο συγκεκριμένα με τους αλγόριθμους που εφαρμόζονται στην επιστήμη των δεδομένων. Θα δοθεί αναλυτική εξήγηση για πολλούς χρήσιμους αλγόριθμους και θα υλοποιηθούν οι πιο σημαντικοί από αυτούς με τη χρήση της γλώσσας προγραμματισμού C++ η οποία θα μας δώσει τη δυνατότητα να φτιάξουμε γρήγορα προγράμματα.

Οι αλγόριθμοι που θα εξετάσουμε μπορούν να χωριστούν σε πέντε ομάδες με βάση τη χρήση τους και αυτές είναι:

- Αλγόριθμοι ταξινόμησης (Classification)
- Αλγόριθμοι παλινδρόμησης (Regression)
- Αλγόριθμοι ομαδοποίησης (Clustering)

- Αλγόριθμοι μείωσης διαστάσεων (Dimensionality reduction)
- Αλγόριθμοι ενισχυτικής μάθησης (Reinforcement learning)
- Αλγόριθμοι ανίχνευσης ανωμαλίας (Anomaly detection)

Μας ενδιαφέρουν ιδιαίτερα οι πρώτες τρεις ομάδες καθώς είναι αυτές που εμφανίζονται πιο συχνά. Οι παλινδρομητές και οι ταξινομητές είναι πολύ παρόμοιοι αλγόριθμοι οι οποίοι ανήκουν στους αλγόριθμους εποπτευόμενης μάθησης (Supervised learning) και χρησιμοποιούνται για να κάνουν προβλέψεις σύμφωνα με τα δεδομένα που τους έχουμε δώσει. Η διαφορά τους είναι ότι η ταξινόμηση προσπαθεί να προβλέψει την κλάση των καινούριων δεδομένων και να τα κατηγοριοποιήσει σύμφωνα με τις υπάρχουσες κατηγορίες, ενώ η παλινδρόμηση προσπαθεί να προβλέψει την τιμή κάποιου στοιχείου για τα καινούρια δεδομένα σύμφωνα με μια συνάρτηση που έφτιαξε από τα υπάρχοντα δεδομένα. Οι περισσότεροι αλγόριθμοι ταξινόμησης είναι και αλγόριθμοι παλινδρόμησης και το αντίστροφο. Από την άλλη οι αλγόριθμοι ομαδοποίησης είναι μη εποπτευόμενης μάθησης και χρησιμοποιούνται σε αδόμητα δεδομένα για να τα χωρίσουν σε ομάδες.



Σχήμα 2: Αλγόριθμοι μηχανικής μάθησης

## 2 Ταξινόμηση

### 2.1 Εισαγωγή στους ταξινομητές

Σε αυτή την ενότητα θα αναλύσουμε τους αλγορίθμους ταξινόμησης και θα δούμε τη χρήση τους. Για τους αλγορίθμους αυτούς χωρίζουμε το σύνολο των δεδομένων μας σε δυο μέρη:

- Δεδομένα εκπαίδευσης (training dataset)
- Δεδομένα επαλήθευσής (testing dataset)

Τα πρώτα τα χρησιμοποιούμε έτσι ώστε ο αλγόριθμος να βρει κάποιο μοτίβο στα δεδομένα με το οποίο θα μπορεί να κατατάσσει τα καινούρια δεδομένα που δέχεται σε κάποια από τις υπάρχουσες κλάσεις. Αυτή είναι και η διαδικασία εκπαίδευσης του μοντέλου. Αφού η εκπαίδευση τελειώσει τότε θα χρησιμοποιήσουμε τα δεδομένα επαλήθευσής για να επιβεβαιώσουμε την ορθή λειτουργία του μοντέλου. Υπάρχουν τρεις βασικοί τύποι ταξινομητών:

- Δυαδικοί (binary)
- Πολλαπλών κλάσεων (multy-class)
- Πολλαπλών ετικετών (multy-label)

Οι δυαδικοί ταξινομητές χρησιμοποιούνται όταν έχουμε μόνο δύο κλάσεις στις οποίες θέλουμε να εντάξουμε τα δεδομένα, ή όταν η απάντηση που θέλουμε να πάρουμε από το μοντέλο είναι δυαδικής φύσης. Για παράδειγμα, ένα πρόβλημα δυαδικής φύσης θα ήταν να προβλέψουμε εάν ένας ασθενής έχει ή δεν έχει μια ασθένειά σύμφωνα με τις εξετάσεις του.

Οι ταξινομητές πολλαπλών κλάσεων από την άλλη είναι ικανοί να αναγνωρίσουν περισσότερες από δύο κλάσεις και είναι πολύ χρήσιμοι για την αναγνώρισή προτύπων. Συνεχίζοντας με το προηγούμενο παράδειγμα θα θέλαμε χωρίσουμε τους ασθενείς σύμφωνα με την κατάσταση τους σε:

- υγιείς
- ήπια ασθένειά
- σοβαρή ασθένεια

Έτσι οι γιατροί θα μπορούν αν δράσουν ανάλογα.

Οι ταξινομητές πολλαπλών ετικετών δεν έχουν κάποια καινούρια λογική αλλά εφαρμόζουν τις λογικές των προηγούμενων προβλημάτων. Δηλαδή θα μπορούσαμε να θέλουμε να υλοποιήσουμε ένα μοντέλο το οποίο να κάνει και τις δύο προηγούμενες προβλέψεις που συζητήσαμε. Αυτά τα μοντέλα συνήθως δε χρησιμοποιούν ξεχωριστούς αλγορίθμους αλλά συνδυάζουν πολλούς ήδη γνωστούς αλγορίθμους για να φτάσουν στο αποτέλεσμα.



Σχήμα 3: Τύποι ταξινομητών

Στη συνέχεια θα αναλύσουμε τους διασημότερους αλγορίθμους και τη χρήση τους. Οι αλγόριθμοι αυτοί είναι<sup>[2,3,4]</sup>:

- Naive Bayes
- LDA (Linear Discriminant Analysis)
- QDA (Quadratic Discriminant Analysis)
- SVM (Support Vector Machine)
- k-NN (k Nearest Neighbors)
- Decision trees
- Random Forest
- Νευρωνικά δίκτυα (τα οποία θα αναλύσουμε στην παλινδρόμηση)

## 2.2 Naive Bayes

Σύμφωνα με το παρακάτω άρθρο<sup>[5]</sup> ο αλγόριθμος Naive Bayes είναι ένας αλγόριθμος που κάνει την υπόθεση ότι τα χαρακτηριστικά είναι υπό όρους ανεξάρτητα της δεδομένης κλάσης. Αυτή η υπόθεση στην πραγματικότητα δεν ισχύει αλλά ο αλγόριθμος πετυχαίνει πολύ υψηλή ακρίβεια και ταυτόχρονα μεγάλη υπολογιστική απόδοση. Αυτός είναι και ο λόγος που χρησιμοποιείται τόσο συχνά στην επιστήμη των δεδομένων.

Τα πλεονεκτήματα του αλγορίθμου είναι<sup>[6]</sup>:

- Η χρονική πολυπλοκότητα αυξάνεται γραμμικά με το πλήθος των δειγμάτων και των στοιχείων τους
- Χαμηλό variance (η αλλαγή των δεδομένων δεν επηρεάζει πολύ το μοντέλο)
- Μπορούμε εύκολα να προσθέσουμε και άλλα δεδομένα και το μοντέλο θα συνεχίσει την εκπαίδευση χωρίς πρόβλημα
- Δεν είναι επιρρεπής στον θόρυβο
- Δεν επηρεάζεται από έλλειψη τιμών στα δεδομένα

Τα μειονεκτήματά είναι:

- Η υπόθεση ότι τα δεδομένα είναι ανεξάρτητα τον κάνει ακατάλληλο για ορισμένα προβλήματα
- Οι συνδυασμοί στοιχείων που δεν εμφανίζονται στο σύνολο δεδομένων για κάποια πρόβλεψη θα έχει πάντα μηδενική πιθανότητα.
- Οι πιθανότητες που υπολογίζει ενδέχεται να είναι λάθος

Ο αλγόριθμος χρησιμοποιεί τον κανόνα του Bayes:

$$P(y|X) = \frac{P(y) \times P(X|y)}{P(X)}$$

Ο παραπάνω τύπος θα μας δώσει την πιθανότητα ενός δείγματος με στοιχεία  $X$  να ανήκει στην κλάση  $y$ . Το  $X$  είναι διάνυσμα με όλα τα στοιχεία του δείγματος μας. Στο παρακάτω παράδειγμα το  $X$  για το πρώτο δείγμα θα είναι  $\langle 40, 85 \rangle$  και η κλάση θα είναι  $y = 0$ . Σκοπός μας είναι όταν πάρουμε τα

στοιχεία από έναν ασθενή να μπορούμε να προβλέψουμε αν έχει διαβήτη. Θα πρέπει δηλαδή να υπολογίσουμε τις πιθανότητες  $P(y = 0|X = \langle x_1, x_2 \rangle)$  και  $P(y = 1|X = \langle x_1, x_2 \rangle)$  για τον καινούριο ασθενή με μετρήσεις  $x_1, x_2$  και η πρόβλεψη θα είναι αυτή με τη μεγαλύτερη πιθανότητα.

Γλυκόζη	Αρτηριακή πίεση	Διαβήτης
40	85	0
40	92	0
45	63	1
45	80	0
40	73	1
45	82	0
40	85	0
30	63	1
65	65	1
45	82	0
35	73	1
45	90	0
50	68	1
40	93	0
35	80	1
50	70	1

Πίνακας 1: Δείγματα για πρόβλεψη διαβήτη

Μέχρι εδώ ήταν η γνωστή στατιστική. Αλλά όσα περισσότερα στοιχεία έχουμε ανά δείγμα τόσο πιο πολύπλοκος γίνεται ο υπολογισμός του παράγοντα  $P(X|y)$ . Εδώ δίνει λύση στο πρόβλημα ο αλγόριθμος Naive Bayes με την ανεξαρτησία των στοιχείων. Ανεξαρτησία σημαίνει ότι

$$P(\langle x_1, x_2, \dots, x_n \rangle | y) = \prod_{i=1}^n P(x_i | y)$$

Επίσης το  $P(X)$  είναι πάντα σταθερό για τα δεδομένα στοιχεία που έχουμε άρα τελικά πρέπει να υπολογίσουμε το

$$P(y) \times \prod_{i=1}^n P(x_i | y)$$

για κάθε  $y$  και στο τέλος να διαλέξουμε εκείνο το  $y$  που έδωσε τη μεγαλύτερη τιμή (που θα είναι και η κλάση στην οποία ανήκει το δείγμα). Βέβαια, αν τα στοιχεία δεν είναι υπό όρους ανεξάρτητα δεν ισχύει αυτή η απλοποίηση. Παρόλα αυτά μπορεί να χρησιμοποιηθεί σε πολλά σύνολα δεδομένων παρόλο που κάποιες φορές δεν ισχύει και για αυτόν τον λόγο καλείται και αφελής Bayes.

Ο αλγόριθμος είναι ιδανικός για προβλήματα πολλαπλών κλάσεων λόγω του εύκολου υπολογισμού. Πιο συγκεκριμένα ο Naive Bayes χρησιμοποιείται:

- Στην ανίχνευση spam. Μπορεί διαβάζοντας τις λέξεις ενός mail να βρει την πιθανότητα να είναι spam. Κάθε λέξη που διαβάζει θα αυξάνει ή θα μειώνει την πιθανότητα με κάποιον προκαθορισμένο κανόνα μέχρι που στο τέλος θα έχουμε πάρει μια ικανοποιητική απάντηση.
- Στη συναισθηματική ανάλυση. Οι εταιρίες μπορούν να κατατάξουν τα σχόλια των καταναλωτών τους σε θετικά και αρνητικά χρησιμοποιώντας τον αλγόριθμο. Μετά μπορούν να πράξουν ανάλογα έτσι ώστε να ευνοήσουν την εταιρία.



Υπάρχουν πολλές άλλες χρήσεις πέρα από αυτές. Για μια πιο αναλυτική εξήγηση των μαθηματικών του αλγορίθμου μπορείτε να διαβάσετε την εξής έρευνα<sup>[7]</sup>.

Οι αλγόριθμοι QDA και QDA χρησιμοποιούν επίσης τον τύπο του Bayes. Η διαφορά είναι στον τρόπο υπολογισμού του παράγοντα  $P(X|y)$ . Όπως και ο Αλγόριθμος Naive Bayes κάνει μια υπόθεση για να τον υπολογίσει, το ίδιο κάνουν και αυτοί οι αλγόριθμοι. Για να δείτε αυτούς τους τύπους με πλήρη εξήγηση μπορείτε να δείτε το παράρτημα Α'.

## 2.3 Support Vector Machine

Οι μηχανές διανυσμάτων υποστήριξης (SVM) είναι ένα σύνολο μεθόδων εποπτευόμενης μάθησης μέθοδοι που χρησιμοποιούνται για ταξινόμηση και παλινδρόμηση. Ανήκουν σε μια οικογένεια γενικευμένων γραμμικών ταξινομητών. Ένα βασικό χαρακτηριστικό του είναι ότι είναι κάνει προβλέψεις που χρησιμοποιούν τη θεωρία μηχανικής μάθησης για να μεγιστοποιήσει την προγνωστική ακρίβεια ενώ ταυτόχρονα αποφεύγει την υπερβολική προσαρμογή στα δεδομένα (over-fitting)<sup>[8]</sup>.

Τα πλεονεκτήματα του SVM είναι<sup>[9]</sup>:

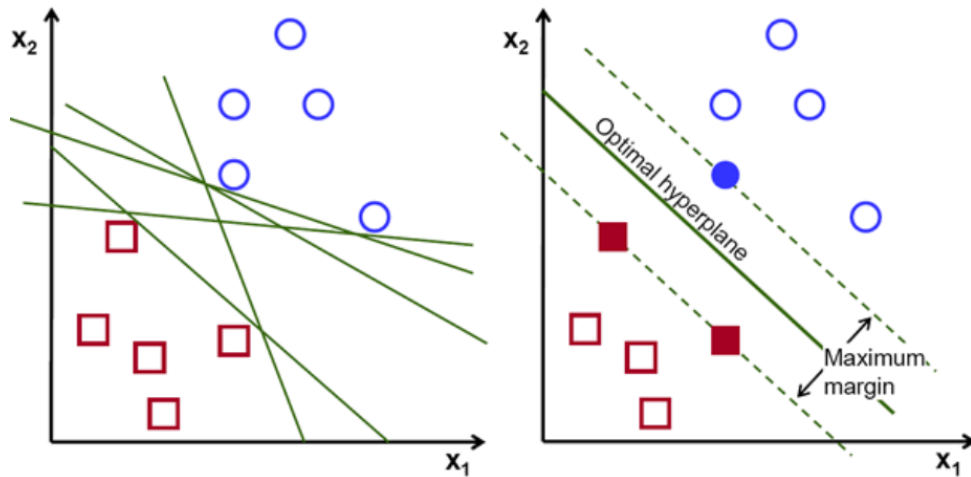
- Δουλεύει πολύ καλά όταν υπάρχει καθαρή διαφοροποίηση στα δεδομένα
- Δουλεύει καλύτερα σε μεγαλύτερες distâncias (πολλά στοιχεία ανά δείγμα)
- Είναι αποδοτικός όταν τα στοιχεία ανά δείγμα είναι περισσότερα από τα δείγματα
- Δεν κάνει μεγάλη χρήση μνήμης

Μειονεκτήματά:

- Δεν είναι κατάλληλος για μεγάλο σύνολο δεδομένων
- Είναι επιρρεπής στον θόρυβο
- Λόγω του τρόπου υλοποίησης του δεν μπορούμε να πάρουμε την πιθανότητα αλλά μόνο την πρόβλεψη

Ο αλγόριθμος αυτός είναι αρκετά απλός και χρησιμοποιείται κυρίως για δυαδική ταξινόμηση. Αρχικά για να τον καταλάβουμε θα πρέπει να αναπαραστήσουμε τα δεδομένα μας σαν σημεία σε έναν χώρο  $N$  διαστάσεων, όπου  $N$  τα στοιχεία που έχουμε ανά δείγμα. Το κάθε στοιχείο θα ανήκει σε μια από τις δύο κατηγορίες τις οποίες θέλουμε να αναγνωρίσουμε. Ο αλγόριθμος λοιπόν προσπαθεί να φέρει ένα υπερεπίπεδο τέτοιο ώστε να χωρίζει τα σημεία ανάλογα με την κατηγορία τους. Έτσι τα σημεία της μιας κατηγορίας θα βρίσκονται από τη μια μεριά και της δεύτερης από την άλλη.

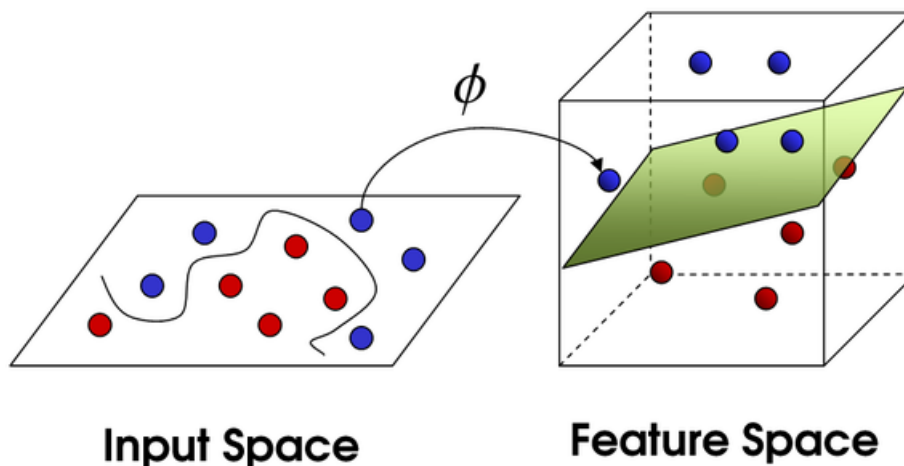
Για να κατανοήσουμε καλύτερα τον αλγόριθμο θα κάνουμε ένα παράδειγμα σε έναν χώρο δύο διαστάσεων που μπορεί εύκολα να αντιληφθεί το ανθρώπινο μάτι.



Σχήμα 4: Παράδειγμα SVM

Στον χώρο δύο διαστάσεων το υπερεπίπεδο είναι απλώς μια ευθεία. Εμείς λοιπόν θέλουμε να βρούμε μια ευθεία που να χωρίζει τα τετράγωνα με τους κύκλους. Όπως μπορούμε να δούμε από το πρώτο διάγραμμα υπάρχουν πολλές ευθείες που μπορεί να χωρίζουν τα σημεία μας αλλά εμείς θέλουμε να διαλέξουμε την ευθεία που να έχει τη μέγιστη απόσταση και από τις δύο μεριές. Επομένως, όταν δώσουμε ένα καινούριο δείγμα στον αλγόριθμο και θέλουμε να το κατατάξει σε μια από τις δύο ομάδες πρέπει απλά να το τοποθετήσει σε αυτόν τον χώρο και μετά να δει από ποια μεριά της ευθείας βρίσκεται. Ανάλογα με τη μεριά αυτή θα αποφασίσουμε και την κλάση του αντικειμένου.

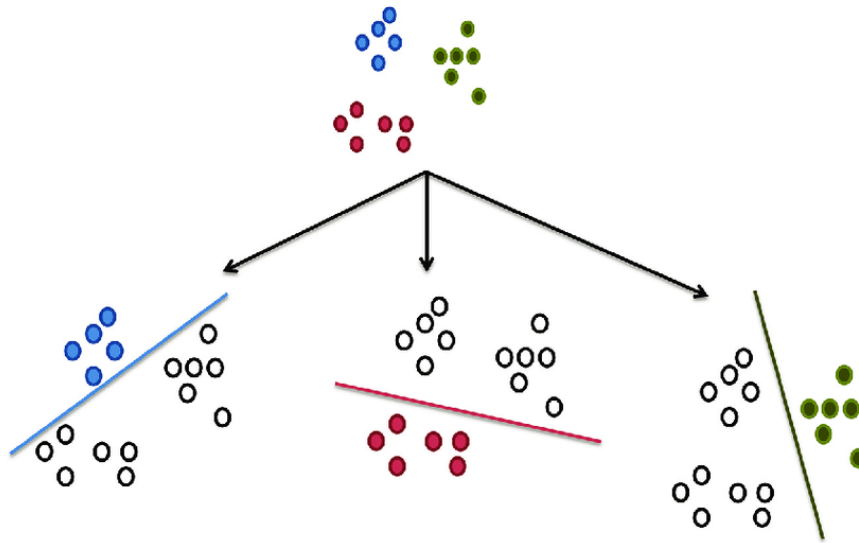
Λύσαμε το πρόβλημα όπου είχαμε πάνω από μία ευθεία αλλά υπάρχει ένα μεγαλύτερο πρόβλημα. Αυτό είναι η περίπτωση που δεν μπορούμε να χωρίσουμε τα στοιχεία μας με μια ευθεία. Εδώ έρχεται το λεγόμενο "kernel trick" για να λύσει το πρόβλημα. Το kernel είναι μια συνάρτηση σύμφωνα με την οποία επαυξάνουμε τη διάσταση του χώρου μας έτσι ώστε με την καινούρια διάσταση που θα δημιουργηθεί τα σημεία μας να είναι διαχωρίσιμα από πλέον ένα επίπεδο. Όταν γυρίσουμε πίσω στην αρχική διάσταση θα φαίνεται σαν τα σημεία μας να είναι διαχωρισμένα με μια καμπύλη.



Σχήμα 5: Επαύξηση διάστασης με SVM

Ένα τελευταίο πρόβλημα που θα μπορούσαμε να λύσουμε θα ήταν να τροποποιήσουμε τον αλγόριθμο έτσι ώστε να δουλεύει για πάνω από δύο κλάσεις. Ένας τρόπος να το κάνουμε αυτό θα ήταν με τη μέθοδο ONE versus ALL. Αυτό που κάνουμε είναι να φέρουμε μια ευθεία για κάθε κλάση που έχουμε η οποία θα χωρίζει τον χώρο στα στοιχεία που ανήκουν σε αυτήν την κλάση και σε όλα τα

υπόλοιπα. Έτσι συνδυάζοντάς την πληροφορία και από τις τρεις ευθείες μπορούμε να πάρουμε τη σωστή απάντηση.



Σχήμα 6: Ταξινόμηση πολλαπλών κλάσεων με SVM

Βέβαια, κάνοντας αυτή τη μέθοδο θα εμφανιστούν σημεία στον χώρο τα οποία μπορεί να ανήκουν σε πάνω από μια κλάση και και άλλα που δεν θα ανήκουν σε καμία. Θα πρέπει λοιπόν να λάβουμε υπόψη μας όχι μόνο από ποια μεριά της κάθε ευθείας βρισκόμαστε αλλά και την απόσταση από αυτές για να μπορέσουμε να καταλήξουμε σε ένα σωστό συμπέρασμα.

Μερικές χρήσεις του αλγορίθμου είναι<sup>[10]</sup>:

- Ανίχνευση προσώπου
- Κατηγοριοποίηση κειμένου
- Ταξινόμηση εικόνων
- Βιοπληροφορική (πχ. ανίχνευση καρκίνου)
- Αναγνώριση γραφικού χαρακτήρα
- Γενικευμένος προγνωστικός έλεγχος (έλεγχος χαοτικών συστημάτων)

## 2.4 K Nearest Neighbors

Ο αλγόριθμος k-NN είναι ένας μη παραμετρικός αλγόριθμος ο οποίος χρησιμοποιεί την εγγύτητα για να κατατάξει τα καινούρια στοιχεία σε μια από τις υπάρχουσες κλάσεις<sup>[11]</sup>. Αυτό σημαίνει ότι το κάθε καινούριο στοιχείο θα παίρνει την κλάση του βλέποντας τις κλάσεις των πιο κοντινών γειτόνων του διαλέγοντας την κλάση της πλειοψηφίας.

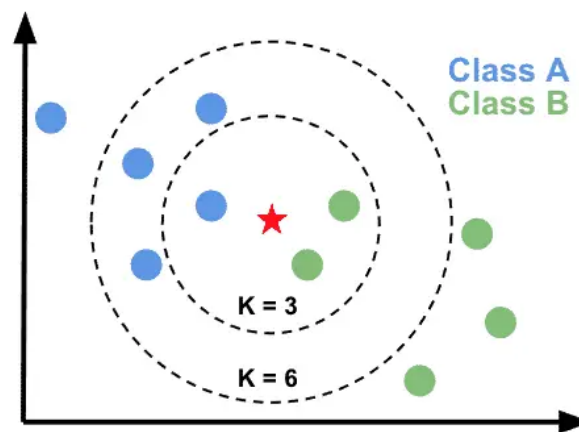
Τα πλεονεκτήματα του k-NN είναι:

- Εύκολη υλοποίηση
- Εύκολη προσαρμογή κατά την εμφάνιση καινούριων δεδομένων
- Δε χρειάζεται παραμέτρους παρά μόνο το k και μια μετρική απόστασης

Κάποια μειονεκτήματά του είναι:

- Δεν έχει καλή επεκτασιμότητα
- Δεν δουλεύει καλά με μεγάλες διαστάσεις
- Μπορεί να γίνει εύκολα over fitting

Για να ξεκινήσουμε τον αλγόριθμο πρέπει πρώτα να ορίσουμε ένα  $k$ . Αυτό θα μας λέει πόσα γειτονικά δείγματα να λάβουμε υπόψη μας για τον προσδιορισμό της κλάσης του καινούριου δείγματος. Έπειτα βρίσκουμε τις αποστάσεις του σημείου από όλα τα σημεία με μια από τις μετρικές που θα αναλύσουμε παρακάτω και διαλέγουμε τα  $k$  σημεία με την κοντινότερη απόσταση στο δικό μας. Τέλος, η κλάση του σημείου μας θα είναι η κλάση που έχει η πλειοψηφία των  $k$  σημείων.



Σχήμα 7: Αλγόριθμος  $k$ -NN

Για της μετρικές απόστασης γενικότερα στην επιστήμη των δεδομένων και στη μηχανική μάθηση έχουμε τέσσερις βασικές επιλογές:

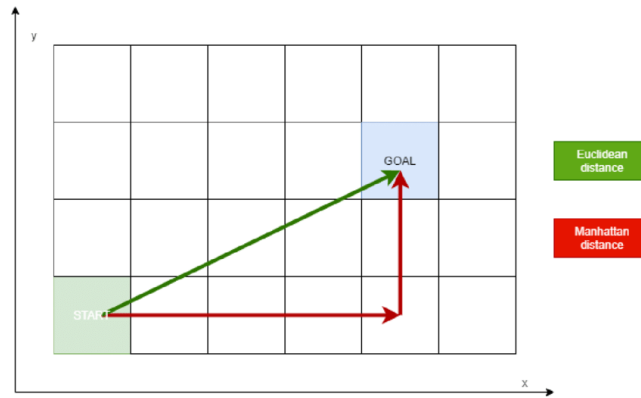
- Ευκλείδεια απόσταση
- Απόσταση Manhattan
- Απόσταση Minkowski
- Απόσταση Hamming

Η Ευκλείδεια απόσταση είναι η πιο γνωστή μετρική και γραφικά είναι το μήκος της ευθείας γραμμής που ενώνει τα δύο σημεία μεταξύ τους. Δίνεται από τον γνωστό τύπο:

$$D = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

Η απόσταση Manhattan είναι επίσης πολύ γνωστή την οποία μπορούμε να παραστήσουμε γραφικά όπως το παρακάτω σχήμα και έχει τύπο:

$$D = \sum_{i=1}^n |x_i - y_i|$$

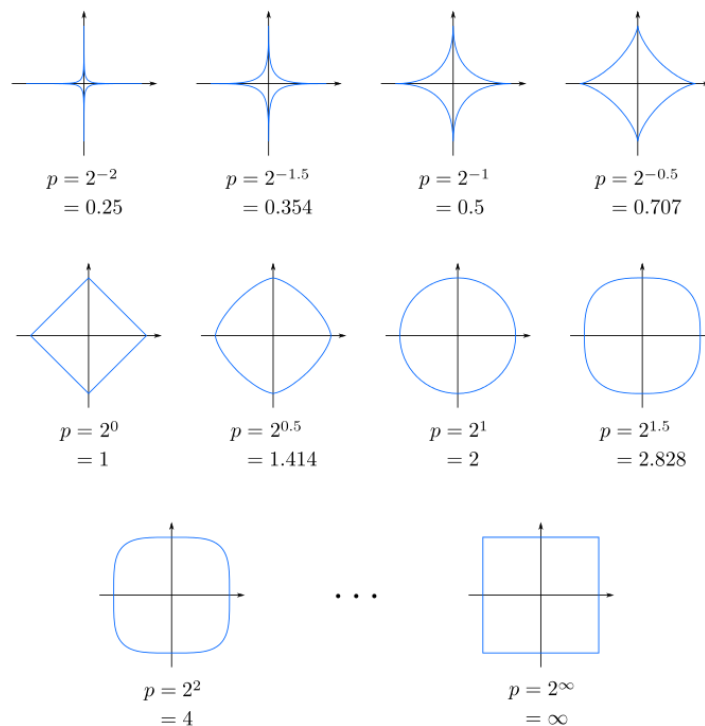


Σχήμα 8: Ευκλείδεια απόσταση και Απόσταση Manhattan

Η απόσταση Minkowski είναι ένα σύνολο αποστάσεων που αποτελεί μια γενίκευση των παραπάνω δύο με τύπο:

$$D = \left( \sum_{i=1}^n |x_i - y_i|^p \right)^{\frac{1}{p}}$$

Όπως μπορούμε αν δούμε πως για  $p = 1$  και  $p = 2$  παίρνουμε την απόσταση Manhattan και την ευκλείδεια απόσταση αντίστοιχα. Δεν μπορούμε να την αναπαραστήσουμε γραφικά αλλά για περαιτέρω κατανόησή μπορούμε να δούμε το παρακάτω σχήμα όπου όλα τα σημεία στην μπλε γραμμή ισαπέχουν από το κέντρο τους με μετρική την απόσταση Minkowski για διαφορετικά  $p$

Unit circles with various values of  $p$  (Minkowski distance)  
OpenGenus

Σχήμα 9: Απόσταση Minkowski

Τέλος, η απόσταση Hamming είναι πολύ απλή και εφαρμόζεται σε δυαδικά διανύσματα και είναι

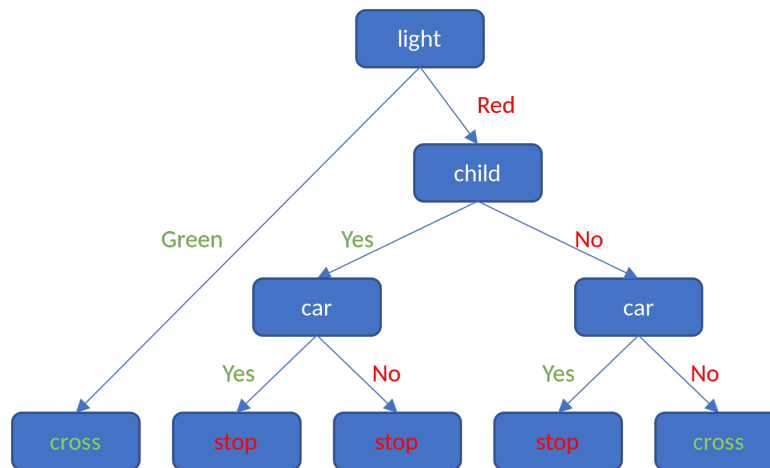
το πλήθος των συνιστωσών που είναι διαφορετικές μεταξύ τους. Για παράδειγμα, αν έχουμε τα διανύσματα  $x = \langle 0, 1, 1, 0, 1 \rangle$ ,  $y = \langle 1, 1, 0, 0, 0 \rangle$  η απόσταση τους θα είναι 3 γιατί οι πρώτες, τρίτες και πέμπτες συνιστώσες είναι διαφορετικές μεταξύ τους.

Οι χρήσεις του k-NN είναι πολύ παρόμοιες με αυτές του SVM. Μια ακόμα χρήση του είναι στην επιλογή προτεινόμενων που υπάρχουν σε πολλές μηχανές αναζήτησης. Μπορεί να βρει τις προτάσεις που θα ταίριαζαν στον κάθε χρήστη ανάλογα με τα "κλικ" που κάνει βρίσκοντας έτσι τα ενδιαφέροντα του.

## 2.5 Decision Trees

Τα δέντρα αποφάσεων είναι μη παραμετρικοί αλγόριθμοι που χρησιμοποιούνται για ταξινόμηση και παλινδρόμηση. Η αναπαράσταση τους γίνεται με τη μορφή δέντρου με κόμβους που χωρίζονται σε κόμβους αποφάσεων και στους τελικούς κόμβους.

Κάθε δέντρο ξεκινάει από έναν κόμβο απόφασης με δύο ή περισσότερους δρόμους που οδηγούν σε άλλους. Ο κάθε κόμβος απόφασης παρουσιάζει μια συνθήκη για τα δεδομένα και ο δρόμος που θα διαλέξουμε εξαρτάτε από το αν η συνθήκη είναι αληθείς ή όχι. Αν συνεχίσουμε αυτή τη διαδικασία κάποια στιγμή θα καταλήξουμε σε έναν τελικό κόμβο ο οποίος θα μας δίνει σαν απάντηση μία από τις κλάσεις που έχουμε ορίσει για το σύνολο δεδομένων.



Σχήμα 10: Παράδειγμα δέντρου απόφασης

Για παράδειγμα στο σχήμα 10 θέλουμε ο υπολογιστής να πάρει την απόφαση για το αν ο άνθρωπος πρέπει να περάσει τον δρόμο ή όχι. Για να το κάνει αυτό περνάει από πολλά στάδια αποφάσεων πριν καταλήξει στην τελική του απόφαση.

Στην πραγματικότητα όμως τα δεδομένα δε θα είναι τόσο απλά. Τα δεδομένα μπορεί να μη διαχωρίζονται με τις συνθήκες που έχουμε επιλέξει για τους κόμβους απόφασης. Αυτό όμως είναι ένα πρόβλημα γιατί ένας τελικός κόμβος θα πρέπει να μας υποδεικνύει μόνο μία κλάση που σημαίνει ότι το σύνολο των συνθηκών που περάσαμε για να φτάσουμε εκεί θα πρέπει να ισχύει μόνο για εκείνη την κλάση.

Για να μπορέσει ο υπολογιστής να βρει τις βέλτιστες συνθήκες για κάθε κόμβο πρέπει να χρησιμοποιήσουμε τη θεωρία της πληροφορίας. Θα κάνουμε ένα παράδειγμα για καλύτερη κατανόηση του αλγορίθμου.

Έστω ότι έχουμε ένα πρόβλημα δυαδικής ταξινόμησης και έναν κόμβο απόφασης που στον οποίο έχουμε 20 στοιχεία και τα σημεία είναι διαχωρισμένα στις 2 κλάσεις στη μέση (50% – 50%). Τώρα θέλουμε να επιλέξουμε ανάμεσα σε 2 συνθήκες:

- Συνθήκη A: Χωρίζει τα σημεία σε 2 κόμβους όπου:

- κόμβος 1: Έχει 14 σημεία με διαχωρισμό 42.8% – 57.2%
- κόμβος 2: Έχει 6 σημεία με διαχωρισμό 33.3% – 66.7%
- Συνθήκη B: Χωρίζει τα σημεία σε 2 κόμβους όπου:
  - κόμβος 1: Έχει 4 σημεία με διαχωρισμό 0% – 100%
  - κόμβος 2: Έχει 16 σημεία με διαχωρισμό 37.5% – 62.5%

Μπορούμε να υπολογίσουμε την εντροπία από τον τύπο:

$$E = \sum_{i=1}^n -p_i \log_2 p_i$$

όπου  $p_i$  η πιθανότητα της κλάσης  $i$ . Άρα έχουμε:

$$E_{parent} = -0.5 \log_2 0.5 - 0.5 \log_2 0.5 = 1$$

$$E_{A1} = -0.428 \log_2 0.428 - 0.572 \log_2 0.572 = 0.99$$

$$E_{A2} = -0.333 \log_2 0.333 - 0.667 \log_2 0.667 = 0.91$$

$$E_{B1} = -0 \log_2 0 - 1 \log_2 1 = 0$$

$$E_{B2} = -0.375 \log_2 0.375 - 0.625 \log_2 0.625 = 0.95$$

Παρατηρούμε ότι όσο καλύτερος είναι ο διαχωρισμός των στοιχείων, τόσο μικρότερη είναι η εντροπία. Το κέρδος πληροφορίας για τις 2 συνθήκες είναι:

$$I_A = E_{parent} - \frac{14}{20} E_{A1} - \frac{6}{20} E_{A2} = 0.034$$

$$I_B = E_{parent} - \frac{4}{20} E_{B1} - \frac{16}{20} E_{B2} = 0.24$$

Ισχύει ότι  $I_B > I_A$  και άρα θα προτιμήσουμε τη συνθήκη B. Αυτός ο έλεγχος θα πρέπει να γίνει για όλες τις πιθανές συνθήκες και για όλους τους κόμβους ξεκινώντας από τον αρχικό κόμβο. Είναι σημαντικό να αναφέρουμε ότι αυτός είναι ένας άπληστος αλγόριθμος διότι επιλέγει την καλύτερη συνθήκη για έναν κόμβο αλλά μπορεί μια χειρότερη επιλογή να οδηγούσε σε καλύτερο αποτέλεσμα αν βλέπαμε τη συνολική εικόνα.

Ένα πρόβλημα με τα δέντρα είναι ότι έχουν μεγάλο variance και γι' αυτό τον λόγο πολλές φορές δεν μπορούν να γενικευτούν για καινούρια δεδομένα. Έτσι σαν μια εξέλιξη του αλγορίθμου δημιουργήθηκε ο αλγόριθμος του τυχαίου δάσους (Random Forest). Ο αλγόριθμος αυτός είναι συνδυασμός μεθόδων και όπως φαίνεται από το όνομα του αποτελείται από πολλά decision trees. Μετά για την τελική απόφαση αν έχουμε πρόβλημα ταξινόμησης θα επιλέξουμε την κλάση που δείχνει η πλειοψηφία, ενώ αν έχουμε πρόβλημα παλινδρόμησης θα πάρουμε τον μέση τιμή των απαντήσεων<sup>[12]</sup>.

Το παραπάνω εξηγεί τη λέξη δάσος όμως δεν εξηγεί τη λέξη τυχαίο. Ο λόγος που είναι τυχαίο είναι επειδή το κάθε δέντρο δε συμπεριλαμβάνει όλα τα δείγματα αλλά ένα τυχαίο τμήμα αυτών και επιπλέον για αυτά τα δείγματα δε θα έχουν όλα τα χαρακτηριστικά τους αλλά θα επιλεχθούν κάποια τυχαία. Παρακάτω βλέπουμε ένα παράδειγμα για το πώς θα χωρίζαμε τα δεδομένα μας για ένα τυχαίο δάσος που θα αποτελείται από τέσσερα δέντρα.

$id$	$x_0$	$x_1$	$x_2$	$x_3$	$x_4$	$y$
0	4.3	4.9	4.1	4.7	5.5	0
1	3.9	6.1	5.9	5.5	5.9	0
2	2.7	4.8	4.1	5.0	5.6	0
3	6.6	4.4	4.5	3.9	5.9	1
4	6.5	2.9	4.7	4.6	6.1	1
5	2.7	6.7	4.2	5.3	4.8	1

Πίνακας 2: Σύνολο δεδομένων για τον αλγόριθμο random forest

$id$	$id$	$id$	$id$
2	2	4	3
0	1	1	3
2	3	3	2
4	1	0	5
5	4	0	1
5	4	2	2

Πίνακας 3: Σύνολο δεδομένων για κάθε δέντρο

Επίσης, όπως είπαμε το κάθε υποσύνολο δεδομένων δε θα περιέχει όλα τα στοιχεία. Για παράδειγμα, το πρώτο θα μπορούσε να έχει μόνο τα στοιχεία  $x_1, x_2$ , το δεύτερο τα  $x_1, x_4$  κ.ο.κ.

Τα πλεονεκτήματα των δέντρων είναι ότι δε χρειάζεται προ επεξεργασία των δεδομένων. Στα δεδομένα πολλές φορές λείπουν στοιχεία που πρέπει ή μπορεί να χρειάζεται κάποια κανονικοποίηση για να εφαρμόσουμε κάποιους αλγόριθμους. Τα δέντρα όμως δεν είναι ένας από αυτούς και αυτό τα κάνει πολύ καλή επιλογή σε τέτοιες περιπτώσεις. Επιπλέον, όπως είπαμε είναι μη παραμετρικοί και επίσης δεν είναι γραμμικοί που είναι κάτι που θέλουμε πολλές φορές για καλύτερο διαχωρισμό των κλάσεων<sup>[13]</sup>. Το βασικό πρόβλημα του αλγορίθμου είναι το overfitting το οποίο λύνει ο αλγόριθμος random forest.



### 3 Παλινδρόμηση

#### 3.1 Γραμμική Παλινδρόμηση

Η πιο απλή και ευρέως χρησιμοποιούμενη τεχνική παλινδρόμησης είναι η γραμμική. Όπως υποδεικνύει το όνομα της η γραμμική παλινδρόμηση συσχετίζει μία ανεξάρτητη με μία εξαρτημένη μεταβλητή με μία γραμμική σχέση προσαρμόζοντας μία ευθεία γραμμή παλινδρόμησης πάνω στα δεδομένα. Η ευθεία δίνεται από την σχέση:

$$y = ax + b + c$$

όπου:

$y$  : εξαρτημένη μεταβλητή (dependent variable)

$x$  : ανεξάρτητη μεταβλητή (Independent / explanatory variable or regressor)

$a$  : κλίση της ευθείας (slope)

$b$  : σημείο τομής με τον άξονα  $y$  ( $y$  - intercept)

$c$  : σφάλμα παλινδρόμησης (regression residual / error)

Η ευθεία αυτή προσαρμόζεται στα δεδομένα με την χρήση της τεχνικής του μικρότερου αθροίσματος των τετραγώνων (των residuals). Τα residuals(σφάλματα) είναι η απόκλιση, ως προς το  $y$ , των δειγμάτων από την γραμμή παλινδρόμησης και βρίσκεται από τον τύπο:

$$c = y_i - \hat{y}_i$$

Παίρνοντας τώρα το άθροισμα των τετραγώνων των σφαλμάτων (Sum of Squared Residuals):

$$SSR = \sum_{i=1}^n e_i^2$$

Ψάχνουμε μία ευθεία γραμμή (παλινδρόμησης) η οποία θα ελαχιστοποιεί το άθροισμα. Στην περίπτωση της απλής γραμμικής παλινδρόμησης το άθροισμα ελαχιστοποιείται για:

$$a = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$b = \bar{y} - a\bar{x}$$

Το τετράγωνο των σφαλμάτων επιλέχτηκε έτσι ώστε να μην αλληλοαναιρούνται τα θετικά με τα αρνητικά σφάλματα. Το πλεονέκτημα έναντι της απόλυτης τιμής είναι ότι το τετράγωνο είναι μία συνεχής και παραγωγίσιμη συνάρτηση. Για ποιο περίπλοκες μορφές γραμμικής παλινδρόμησης χρησιμοποιούνται και επαναληπτικές μέθοδοι προσέγγισης της καλύτερης ευθείας (line of best fit), εξετάζοντας επαναληπτικά τιμές για την κλίση και την μετατόπιση μέχρι να βρούμε το ζευγάρι που ελαχιστοποιεί το άθροισμα των τετραγώνων.

Μια μετρική της απλής γραμμικής παλινδρόμησης είναι το μέσο τετραγωνικό σφάλμα (Mean Square Error), όπου για μία σταθερή διακύμανση του σφάλματος το MSE υπολογίζεται ως:

$$\sigma_\epsilon^2 = \frac{SSR}{n - p - 1}$$

όπου:

$n$  : πλήθος δειγμάτων

$p$  : πλήθος ανεξάρτητων μεταβλητών (regressors)

Για την περίπτωση της απλής γραμμικής παλινδρόμησης που εξετάζουμε τώρα το  $p = 1$

Μια άλλη μετρική που χρησιμοποιείται κατά κόρον είναι το  $R^2$  και στη συνέχεια θα δούμε πως υπολογίζεται. Αρχικά υπολογίζουμε τον μέσο όρο των δειγμάτων (ως προς τον άξονα της εξαρτημένης μεταβλητής). Η τιμή αυτή θα είναι το σημείο τομής του άξονα της εξαρτημένης μεταβλητής και μίας ευθείας παράλληλης ως προς τον άξονα της ανεξάρτητης μεταβλητής. Αυτή η ευθεία είναι μία ευθεία παλινδρόμησης (mean) όταν δεν παίρνουμε υπόψιν μας την εξάρτηση από την ανεξάρτητη μεταβλητή. Υπολογίζουμε το SSR και την διακύμανση γύρω από αυτή την γραμμή παλινδρόμησης. Έπειτα υπολογίζουμε την διακύμανση και γύρω από την γραμμή παλινδρόμησης που εξετάζουμε (fit). Το  $R^2$  υπολογίζεται από τον τύπο:

$$R^2 = \frac{\text{var}(\text{mean}) - \text{var}(\text{fit})}{\text{var}(\text{mean})}$$

Η σχέση αυτή μας δείχνει την μείωση της διακύμανσης όταν παίρνουμε υπόψιν μας την ανεξάρτητη μεταβλητή ή αλλιώς το ποσοστό της μεταβολής της εξαρτημένης μεταβλητής που εξηγεί η μεταβολή της μεταβλητής που εξετάζουμε ή θα μπορούσαμε να πούμε ότι προσδιορίζει το **ποσοστό της βεβαιότητας** που έχουμε όταν κάνουμε μια πρόβλεψη της τιμής της εξαρτημένης για μία συγκεκριμένη τιμή της ανεξάρτητης. Το ίδιο αποτέλεσμα θα είχαμε και εάν υπολογίζαμε το  $R^2$  με τα SSR αντί της διακύμανσης λόγω του ότι απαλείφονται οι παρονομαστές.

Το  $R^2$  είναι προφανές ότι παίρνει τιμές από 0 ως 1 και όσο το μεγαλύτερο τόσο το καλύτερο. Είναι μια μετρική μεγάλης σημασίας για τον προσδιορισμό της ακαταλληλότητας μιας ανεξάρτητης μεταβλητής, καθώς για μεγάλα δείγματα ο προσδιορισμός αυτός και η απόρριψη μεταβλητών μικρότερης σημασίας οδηγεί σε μεγάλη μείωση του χρόνου υπολογισμού και πόρων. Ωστόσο χρειαζόμαστε έναν τρόπο για να προσδιορίζουμε το πόσο σημαντικό είναι στατιστικά (πόσο ακριβές είναι), καθώς όπως φαίνεται και παρακάτω ένα με την προσθήκη έξτρα δεδομένων το  $R^2$  αυξάνεται χωρίς να έχει βελτιωθεί απαραίτητα το μοντέλο μας.

Το **p-value** είναι μια μετρική που κάνει ακριβός αυτή την δουλειά. Το p-value για το  $R^2$  υπολογίζεται από έναν άλλον όρο το  $F$ . Το  $F$  μοιάζει πολύ με το  $R^2$ , η μόνη διαφορά τους βρίσκεται στον παρονομαστή με το  $R^2$  να είναι:

$$R^2 = \frac{\text{διακύμανση του } y \text{ που εξηγείται από το } x}{\text{διακύμανση του } y \text{ όταν δεν παίρνουμε υπόψιν το } x}$$

Ενώ το  $F$ :

$$F = \frac{\text{διακύμανση του } y \text{ που εξηγείται από το } x}{\text{διακύμανση του } y \text{ που δεν εξηγείται από το } x}$$

Ή αλλιώς:

$$F = \frac{SS(\text{mean}) - \frac{SS(\text{fit})}{p_{\text{fit}}} - p_{\text{mean}}}{\frac{SS(\text{fit})}{n} - p_{\text{fit}}}$$

Όπου:

$p_{\text{fit}}$  : παράμετροι (βαθμοί ελευθερίας) για την ευθεία που εξετάζουμε (2 για την απλή γραμμική παλινδρόμηση)

$p_{\text{mean}}$  : παράμετροι (βαθμοί ελευθερίας) για την ευθεία αναφοράς = 1 (αφού η μόνη παράμετρος που εξετάζουμε είναι η μετατόπιση  $b$ )

$n$  : αριθμός των δεδομένων

Ο υπολογισμός τώρα του p-value από το  $F$  γίνεται μέσω της παραγωγής τυχαίων δεδομένων και υπολογίζοντας το  $F$  για αυτό το σετ. Αυτό γίνεται για αρκετές επαναλήψεις (όσο περισσότερες τόσο το καλύτερο) και τα αποτελέσματα ( $F$ ) που παίρνουμε τα τοποθετούμε σε ένα ιστόγραμμα. Έπειτα βρίσκουμε το  $F$  για το σετ δεδομένων που εξετάζουμε και το τοποθετούμε και αυτό στο ιστόγραμμα.

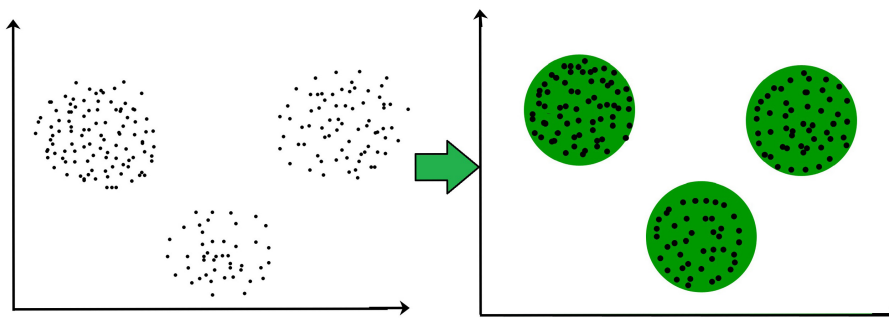
Το p-value υπολογίζεται από τον αριθμό των τιμών που είναι μεγαλύτερες από το  $F$  του σετ που εξετάζουμε διαιρεμένο με τον αριθμό όλων των τιμών. Στην πράξη το γράφημα του ιστογράμματος προσεγγίζεται από μία γραμμή  $F$  κατανομής, έτσι ώστε να μην χρειαστεί να παράξουμε τόσο μεγάλο όγκο τυχαίων δεδομένων. Το σχήμα της γραμμής εξαρτάται από τον αριθμό των δεδομένων και τους βαθμούς ελευθερίας. Το p-value παίρνει μικρότερες τιμές όσο περισσότερα δείγματα έχουμε ανά παράμετρο. Όσο μικρότερο το p-value τόσο καλύτερο το  $R^2$  και κατά συνέπεια η γραμμή παλινδρόμησης

## 4 Ομαδοποίηση

### 4.1 Εισαγωγή στην ομαδοποίηση

Η ενότητα αυτή αφορά τους αλγορίθμους ομαδοποίησης (clustering). Η δουλειά των αλγορίθμων αυτών είναι να χωρίσουν το σύνολο των δεδομένων σε ομάδες τέτοιες ώστε τα δεδομένα μιας ομάδας να εμφανίζουν περισσότερες ομοιότητες μεταξύ τους από ότι με τα δεδομένα των άλλων ομάδων. Αυτός είναι και ο βασικός στόχος της διερευνητικής και στατιστικής ανάλυσης και χρησιμοποιείται για<sup>[14]</sup>:

- Αναγνώριση προτύπων
- Ανάλυση εικόνων
- Εξαγωγή πληροφορίας
- Συμπύεση δεδομένων
- Γραφικά υπολογιστών
- Μηχανική μάθηση



Σχήμα 11: Γραφική αναπαράσταση ομαδοποίησης

Υπάρχουν εκατοντάδες αλγόριθμοι ομαδοποίησης οι οποίοι χρησιμοποιούν διαφορετικές μεθόδους και τεχνικές. Δεν υπάρχει κάποιος που να υπερτερεί πλήρως των υπολοίπων αλλά η καταλληλότητα του Κάθε αλγορίθμου εξαρτάται από το σύνολο δεδομένων. Παρά τις διαφορές τους μπορούμε να τους εντάξουμε σε κατηγορίες συμφωνά με τον τρόπο που κάνουν την ομαδοποίηση. Οι τέσσερις πιο βασικές κατηγορίες είναι:

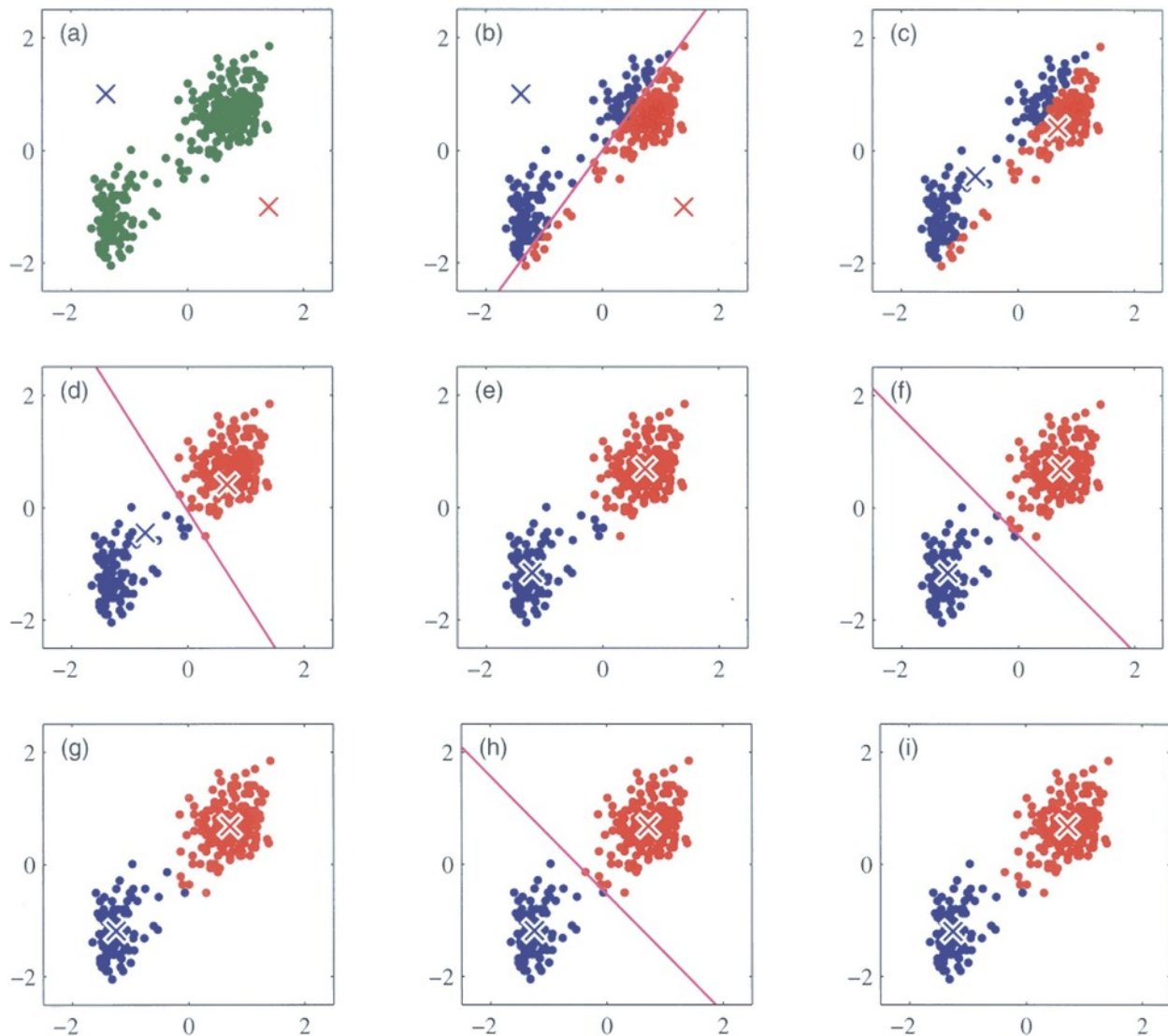
- Centroid based clustering
- Density based clustering
- Distribution based clustering
- Hierarchical clustering

Οι κατηγορίες αυτές θα αναλυθούν περισσότερο στη συνέχεια. Επιπλέον, θα δούμε και τους πιο διάσημους αλγορίθμους που ανήκουν σε αυτές της κατηγορίες και θα αναλύσουμε και τον τρόπο λειτουργίας τους.

## 4.2 Centroid based clustering

Οι αλγόριθμοι τύπου Centroid based clustering όπως μπορούμε να καταλάβουμε κάνουν ομαδοποίηση με βάση τα κέντρα. Αυτό σημαίνει ότι θα έχουμε τόσα κέντρα όσες και οι ομάδες που επιλέξαμε. Έπειτα για να ομαδοποιήσουμε όλα τα σημεία του χώρου θα βλέπουμε σε ποιο κέντρο βρίσκονται πιο κοντά και θα τα εντάσσουμε στην ομάδα αυτού του κέντρου. Το πρόβλημα εδώ είναι πως θα βρούμε τα κατάλληλα κέντρα. Αν δεν επιλέξουμε τα κέντρα σωστά τότε η ομαδοποίηση που θα κάνουμε θα απέχει πολύ από την πραγματικότητα. Γι' αυτό και οι συγκεκριμένοι αλγόριθμοι έχουν ως βασική δουλειά να προσεγγίσουν αυτά τα κέντρα.

Για να καταλάβουμε με ποιον τρόπο συμβαίνει αυτό θα δούμε τον αλγόριθμο k-Means που είναι και ο πιο διάσημος αλγόριθμος ταξινόμησης με βάση τα κέντρα. Το  $k$  είναι το πλήθος των ομάδων που θέλουμε. Ξεκινάμε τον αλγόριθμο επιλέγοντας  $k$  τυχαία κέντρα στον χώρο και χωρίζουμε τα υπόλοιπα σημεία σε ομάδες ανάλογα με την εγγύτητα τους στα τυχαία κέντρα. Έτσι έχουμε καταλήξει σε μια αρχική προσέγγιση. Έπειτα θα βρούμε το πραγματικό κέντρο των ομάδων που δημιουργήθηκαν βρίσκοντας τη μέση τιμή των σημείων που ανήκουν στην ομάδα. Αυτή η διαδικασία επαναλαμβάνεται και έτσι κάποια στιγμή τα τελικά κέντρα θα μπορούν να κάνουν καλή ομαδοποίηση.



Σχήμα 12: Επαναλήψεις του αλγορίθμου k-Means

Ένα από τα πλεονεκτήματα του αλγορίθμου είναι η απλότητά του. Επιπλέον λειτουργεί καλά σε

μεγάλο σύνολο δεδομένων και εξασφαλίζει ότι θα συγκλίνει σε κάποια λύση. Επίσης μπορεί να δημιουργήσει ομάδες διαφορετικού μεγέθους και σχήματος. Παρ' όλα αυτά τα αποτελέσματα τα τελικά αποτελέσματα του αλγορίθμου εξαρτώνται από τις αρχικές συνθηκες. Άρα εφόσον αρχικοποιούμε τα κέντρα τυχαία, αν ξανά τρέξουμε τον αλγόριθμο θα πάρουμε διαφορετικά αποτελέσματα. Επίσης πρέπει να ορίσουμε μόνοι μας το πλήθος των ομάδων που είναι κάτι που μπορεί να μην θέλουμε πάντα<sup>[15]</sup>.

### 4.3 Density based clustering

Ο επόμενος τύπος αλγορίθμου που θα δούμε είναι η ομαδοποίηση με βάση την πυκνότητα. Οι συγκεκριμένοι αλγόριθμοι ανιχνεύουν τις περιοχές με μεγάλη συγκέντρωση σημείων οι οποίες διαχωρίζονται μεταξύ τους από περιοχές που τα σημεία είναι πολύ αραιά. Αυτές οι περιοχές θα είναι και οι ομάδες που θα δημιουργηθούν και τα σημεία στις αραιές περιοχές θα θεωρηθούν θόρυβος<sup>[16]</sup>.

Ας δούμε πως λειτουργεί ο αλγόριθμος DBSCAN (Density-based spatial clustering of applications with noise) ο οποίος είναι από τους πιο διάσημους αλγορίθμους ομαδοποίησης. Για αυτόν τον αλγόριθμο πρέπει να ορίσουμε μια παράμετρο ελάχιστης απόστασης και μία παράμετρο ελάχιστων σημείων. Έπειτα θεωρούμε τις εξής κατηγορίες σημείων:

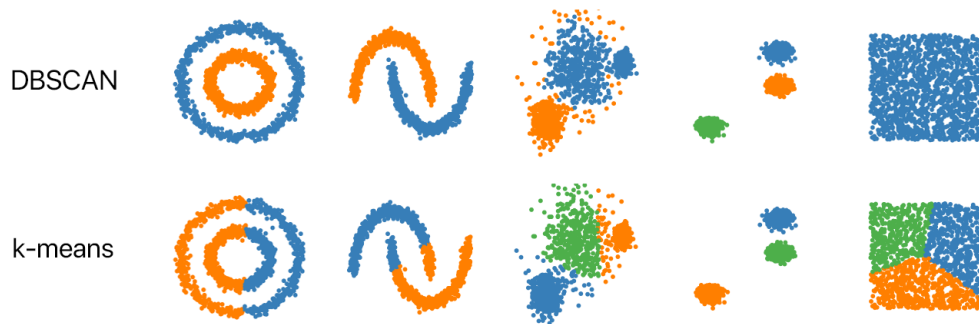
**Κύρια σημεία** Έχει γειτονικά σημεία τουλάχιστον όσα τα ελάχιστα σε απόσταση μικρότερη από την ελάχιστη.

**Άμεσα προσιτά σημεία** Για ένα κύριο σημείο ένα άλλο σημείο είναι προσιτό αν είναι μέσα στην ελάχιστη απόσταση.

**Προσιτά σημεία** Ένα σημείο είναι προσιτό αν υπάρχει ένα μονοπάτι από σημεία που το κάθε ένα είναι άμεσα προσιτό με το διπλανό του. Το σημείο εκκίνησης πρέπει να είναι κύριο σημείο.

**Θόρυβος** Όλα τα υπόλοιπα σημεία που περίσσεψαν.

Έπειτα κάθε κύριο σημείο δημιουργεί μία ομάδα με όλα τα προσιτά σημεία σε αυτό<sup>[17]</sup>.



Σχήμα 13: Αποτελέσματα ομαδοποίησης με DBSCAN vs k-Means

Από το παραπάνω σχήμα μπορεί να καταλάβει κανείς τις περιπτώσεις που ο αλγόριθμος DBSCAN υπερσχύει του k-Means. Όπως φαίνεται ο DBSCAN είναι κατάλληλος όταν οι ομάδες δεν έχουν απλά σχήματα και επιπλέον μπορεί να διαλέξει μόνος του τον κατάλληλο αριθμό ομάδων που πρέπει να δημιουργήσει κάθε φορά. Συνοψίζοντας τα πλεονεκτήματα του αλγορίθμου είναι<sup>[18]</sup>:

- Δεν χρειάζεται να ορίσουμε τον αριθμό των ομάδων
- Δουλεύει καλά σε ομάδες με αυθαίρετα σχήματα

- Μπορεί να καταλάβει ποια σημεία αποτελούν θόρυβο και να μην τα λάβει υπόψη του

Παρ' όλα αυτά υπάρχουν και μειονεκτήματα:

- Δεν μπορεί να ομαδοποιήσει δεδομένα που έχουν μεγάλες διαφορές μεταξύ τους γιατί δεν μπορούμε να επιλέξουμε τις παραμέτρους του συστήματος ώστε να ευνοούν όλες τις ομάδες
- Επίσης δεν μπορούμε να διαλέξουμε μια καλή τιμή για την ελάχιστη απόσταση αν δεν έχουμε καλή κατανόηση των δεδομένων
- Δίνει διαφορετικά αποτελέσματα κάθε φορά διότι ξεκινάει την ανάλυση από τυχαία σημεία

## 4.4 Hierarchical clustering

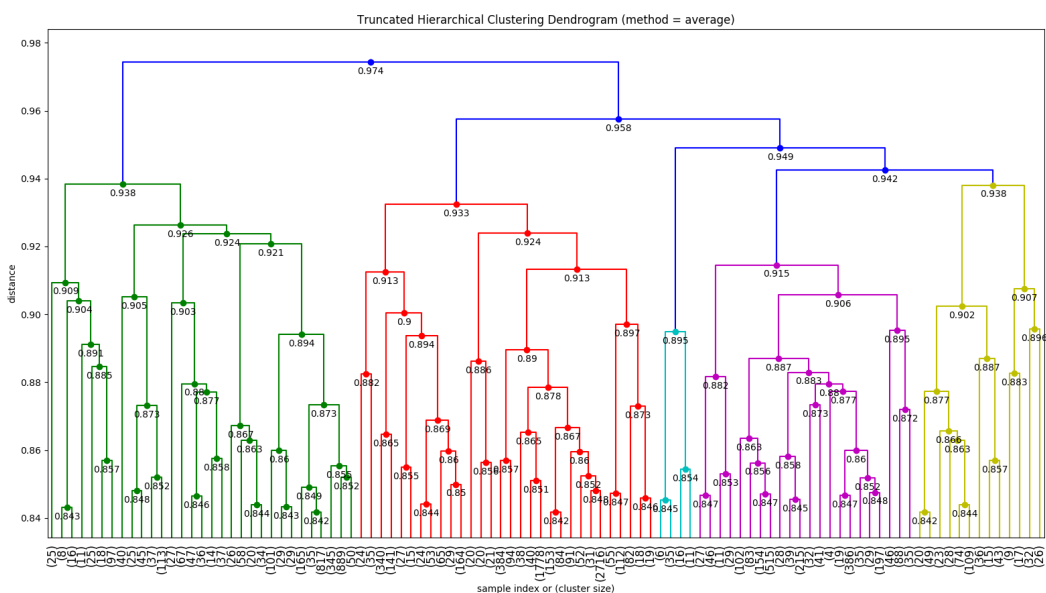
Το Hierarchical clustering είναι μια μέθοδος ομαδοποίησης η οποία προσπαθεί να δημιουργήσει μια ιεραρχία από ομάδες. Για αυτούς τους αλγόριθμους μπορούμε να διακρίνουμε δύο μεγάλες κατηγορίες<sup>[19]</sup>:

**Συσσωματωτική** Προσεγγίζει το πρόβλημα από κάτω προς τα πάνω. Το κάθε δείγμα φτιάχνει τη δικιά του ομάδα και όταν ανεβαίνουμε στην ιεραρχία οι ομάδες αυτές συγχωνεύονται σε μεγαλύτερες

**Διαιρετική** Προσεγγίζει το πρόβλημα από κάτω προς τα πάνω. Ξεκινάμε από μια μεγάλη ομάδα η οποία διαχωρίζεται σε μικρότερες αναδρομικά καθώς κατεβαίνουμε στην ιεραρχία.

Για να αποφασίσουμε πότε δύο ομάδες πρέπει να ενωθούν ή να διαχωριστούν πρέπει να βρούμε μια μετρική ανομοιότητας μεταξύ των σημείων. Αυτή η μετρική τις περισσότερες φορές είναι η απόσταση η οποία μπορεί να είναι η ευκλείδεια ή κάποια άλλη που έχουμε ήδη αναλύσει.

Τα αποτελέσματα του αλγορίθμου μπορούν να αναπαρασταθούν με τη μορφή δέντρου



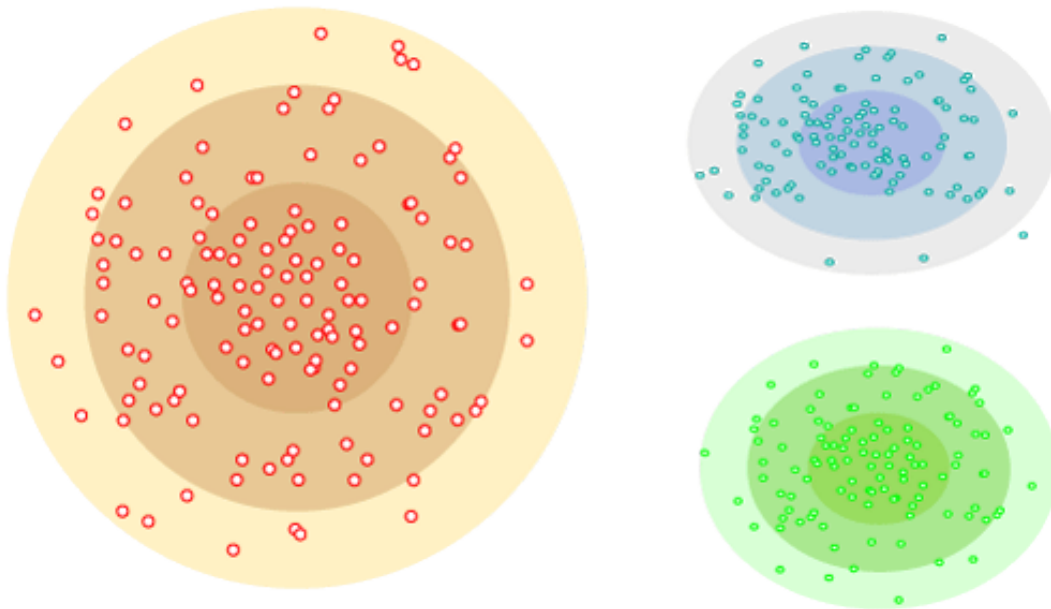
Σχήμα 14: Hierarchical clustering

## 4.5 Distribution based clustering

Το Distribution based clustering είναι μια υποκατηγορία του Hierarchical clustering και χρησιμοποιείται στις περιπτώσεις που τα δεδομένα έχουν κάποια γνωστή κατανομή όπως είναι η κανονική ή η Γκαουσιανή. Το Distribution based clustering μας δίνει τη δυνατότητα να διαχειριστούμε μεγάλα σύνολα δεδομένων και να βρούμε μη γραμμικές σχέσεις ανάμεσα στα στοιχεία<sup>[20]</sup>.

Από τους πιο διάσημους αλγορίθμους είναι ο αλγόριθμος Gaussian mixture models. Όπως μπορούμε να καταλάβουμε από το όνομά του ο αλγόριθμος αυτός κάνει την υπόθεση ότι τα δεδομένα ακολουθούν Γκαουσιανή κατανομή. Στην στατιστική όταν μιλάμε για Γκαουσιανή κατανομή ξέρουμε ότι τα δεδομένα μας θα έχουν μέσο όρο και απόκλιση. Τώρα που τα δεδομένα μας όμως είναι πολλών διαστάσεων θα πρέπει να γενικεύσουμε αυτές τις παραμέτρους. Έτσι κάθε ομάδα που θα δημιουργείται από τον αλγόριθμο θα έχει ένα κέντρο και έναν πίνακα συνδιακύμανσης. Αυτούς τους όρους τους έχουμε ξαναδεί σε προηγούμενη ενότητα και οι τύποι βρίσκονται στο παράρτημα Α'. Μια επιπλέον αναγκαία παράμετρος του συστήματος είναι και η πιθανότητα της κάθε ομάδας. Έτσι με αυτές τις τρεις παραμέτρους μπορούμε να ομαδοποιήσουμε τα δεδομένα.

Η γραφική απεικόνιση των ομάδων στις δύο διαστάσεις θα έχουν τη μορφή ελλείψεων οι οποίες θα περιγράφονται από τις παραμέτρους που αναφέραμε προηγουμένως.



Σχήμα 15: Γραφική αναπαράσταση Distribution based clustering



## 5 Υλοποίηση

Όπως έχουμε ήδη αναφέρει, για την υλοποίηση αλγορίθμων θα χρησιμοποιήσουμε την γλώσσα C++ και σκοπός μας είναι τα μοντέλα που θα δημιουργήσουμε να είναι εύχρηστα και επαναχρησιμοποιήσιμα και από άλλους χρήστες. Θα χρειαστεί επομένως να δημιουργήσουμε κάποια εργαλεία με γενική χρήση.

### 5.1 CSV Parser

Πριν μπορέσουμε να ξεκινήσουμε την υλοποίηση αλγορίθμων θα χρειαστούμε έναν τρόπο για να μεταφέρουμε τα δεδομένα μας μέσα στο πρόγραμμα μας. Ο πιο συνηθισμένος τύπος αρχείου για δεδομένα είναι αυτός με την κατάληξη .csv (comma separated values). Όπως φαίνεται και από το όνομα είναι ένα αρχείο στο οποίο τα δεδομένα είναι διαχωρισμένα με κόμμα. Ένα παράδειγμα είναι το παρακάτω:

```
1 glucose , bloodpressure , diabetes
2 40 , 85 , 0
3 40 , 92 , 0
4 45 , 63 , 1
5 45 , 80 , 0
6 40 , 73 , 1
7 45 , 82 , 0
8 40 , 85 , 0
9 30 , 63 , 1
10 65 , 65 , 1
11 45 , 82 , 0
```

Listing 1: diabetes.csv

Το παραπάνω αρχείο μας δίνει δείγματα τα οποία έχουν ως στοιχεία τη γλυκόζη, την πίεση και τα αν έχουν ή όχι διαβήτη, όπου κάθε καινούρια σειρά είναι και ένα διαφορετικό δείγμα. Θα φτιάξουμε μια συνάρτηση τώρα ώστε να εισάγουμε αυτά τα δεδομένα στο πρόγραμμα μας. Επιπλέον θα φτιάξουμε και μια συνάρτηση για να εμφανίζουμε αυτά τα δεδομένα ώστε να επιβεβαιώσουμε ότι η συνάρτηση δουλεύει σωστά. Για να το κάνουμε αυτό θα χρησιμοποιήσουμε κάποια templates που μας παρέχει η C++ τα οποία υπάρχουν στους παρακάτω headers:

```
1 #pragma once
2
3 #include <iostream>
4 #include <iomanip>
5 #include <fstream>
6 #include <vector>
7 #include <unordered_map>
```

Listing 2: parser includes

Θα χρησιμοποιήσουμε τα πρώτα δύο includes για να εμφανίσουμε τα δεδομένα μορφοποιημένα στο command line, ενώ τα επόμενα τρία θα τα χρειαστούμε για να εισάγουμε τα δεδομένα.

```
1 template <class T>
2 std::unordered_map<std::string, std::vector<T>> parse_csv (
3     std::string filepath, size_t features)
4 {
5     // The data will be held in an unordered map
6     std::unordered_map<std::string, std::vector<float>> data, tempData;
```

```

7 // Csv file to read from
8 std::ifstream file;
9 T temp;
10 std::string key;
11 file.open(filepath);
12 for (size_t i = 0; i < features - 1; i++) {
13     std::getline(file, key, ',');
14     tempData.insert({key, {}});
15 }
16 std::getline(file, key, '\n');
17 tempData.insert({key, {}});
18 for (const auto &pair : tempData) data.insert(pair);
19 bool eof = false;
20 while (true) {
21     for (auto &pair : data) {
22         if (file >> temp) {
23             pair.second.push_back(temp);
24             if (file.peek() == ',' ||
25                 file.peek() == '\n' ||
26                 file.peek() == '\r') file.ignore();
27         }
28         else {
29             eof = true;
30             break;
31         }
32     }
33     if (eof) break;
34 }
35 file.close();
36 return data;
37 }

```

Listing 3: parse\_csv function

Το map είναι μια δομή δεδομένων η οποία αντιστοιχίζει έναν τύπο μεταβλητής σε μία άλλη. Στην συγκεκριμένη περίπτωση αντιστοιχίζουμε strings και συγκεκριμένα τα ονόματα των χαρακτηριστικών των δεδομένων μας σε vectors δηλαδή πίνακες. Μπορούμε να χρησιμοποιήσουμε τώρα το data για να πάρουμε όποια τιμή θέλουμε. Για παράδειγμα το data["glucose"] θα μας επιστρέψει έναν πίνακα που θα περιέχει όλη την πρώτη στήλη από το diabetes.csv και αν θέλουμε να πάρουμε ένα συγκεκριμένο στοιχείο από αυτή τη στήλη μπορούμε απλά να το πάρουμε γνωρίζοντας τη θέση του. Το data["glucose"][2] θα μας επιστρέψει το τρίτο στοιχείο της στήλης.

Τώρα θα ορίσουμε τη συνάρτηση που θα εμφανίζει τα αποτελέσματα στην οθόνη.

```

1 template <class T>
2 void print_map (
3     const std::unordered_map<std::string, std::vector<T>>& data)
4 {
5     size_t size;
6     for (const auto &pair : data) {
7         std::cout << pair.first << '\t';
8         size = pair.second.size();
9     }
10    std::cout << '\n';
11    for (size_t i = 0; i < size; i++) {

```

```
12     for (const auto &pair : data)
13         std::cout << std::setw(pair.first.size())
14                     << std::left
15                     << pair.second[i]
16                     << '\t';
17     std::cout << '\n';
18 }
19 }
```

Listing 4: print\_map function

Για το παράδειγμά μας θα πάρουμε την εξής έξοδο:

```
1 glucose  bloodpressure  diabetes
2 40       85             0
3 40       92             0
4 45       63             1
5 45       80             0
6 40       73             1
7 45       82             0
8 40       85             0
9 30       63             1
10 65       65             1
11 45       82             0
```

Listing 5: output for diabetes.csv

## 6 Συμπέρασμα

## Αναφορές

- [1] Data science. [https://en.wikipedia.org/wiki/Data\\_science](https://en.wikipedia.org/wiki/Data_science). Accessed: 2023-03-13.
- [2] Rohit Garg. 7 types of classification algorithms. <https://analyticsindiamag.com/7-types-classification-algorithms/>. Accessed: 2023-03-15.
- [3] Machine learning classification - 8 algorithms for data science aspirants. <https://data-flair.training/blogs/machine-learning-classification-algorithms/>. Accessed: 2023-03-15.
- [4] Ekin Keserer. 8 types of machine learning classification algorithms. <https://www.akkio.com/post/5-types-of-machine-learning-classification-algorithms>. Accessed: 2023-03-15.
- [5] Webb, Geoffrey I and Keogh, Eamonn and Miikkulainen, Risto. Naïve Bayes. *Encyclopedia of machine learning*, 15:713–714, 2010.
- [6] Pavan Vadapalli. Naive Bayes explained: Function, Advantages and disadvantages applications in 2023. <https://www.upgrad.com/blog/naive-bayes-explained/>. Accessed: 2023-03-18.
- [7] Rish, Irina and others. An empirical study of the naive Bayes classifier. 3(22):41–46, 2001.
- [8] Jakkula, Vikramaditya. Tutorial on support vector machine (svm). *School of EECS, Washington State University*, 37(2.5):3, 2006.
- [9] Dhiraj K. Top 4 advantages and disadvantages of Support Vector Machine or SVM. <https://dhirajkumarblog.medium.com/top-4-advantages-and-disadvantages-of-support-vector-machine-or-svm-a3c06a2b107>. Accessed: 2023-03-18.
- [10] Real-Life Applications of SVM (Support Vector Machines). <https://data-flair.training/blogs/applications-of-svm/>. Accessed: 2023-03-22.
- [11] What is the k-nearest neighbors algorithm? <https://www.ibm.com/topics/knn>. Accessed: 2023-03-24.
- [12] Random Forest. [https://en.wikipedia.org/wiki/Random\\_forest](https://en.wikipedia.org/wiki/Random_forest). Accessed: 2023-04-08.
- [13] Advantages and disadvantages of decision tree in machine learning. <https://www.analytixlabs.co.in/blog/decision-tree-algorithm>. Accessed: 2023-04-08.
- [14] Cluster analysis. [https://en.wikipedia.org/wiki/Cluster\\_analysis](https://en.wikipedia.org/wiki/Cluster_analysis). Accessed: 2023-04-10.
- [15] k-Means Advantages and Disadvantages Machine Learning Google Developers. <https://developers.google.com/machine-learning/clustering/algorithm/advantages-disadvantages>. Accessed: 2023-04-16.
- [16] How Density-based Clustering works. <https://pro.arcgis.com/en/pro-app/latest/tool-reference/spatial-statistics/how-density-based-clustering-works.htm>. Accessed: 2023-04-19.
- [17] DBSCAN. <https://en.wikipedia.org/wiki/DBSCAN>. Accessed: 2023-04-19.

- 
- [18] Density-Based Spatial Clustering of Applications with Noise (DBSCAN). <https://ml-explained.com/blog/dbscan-explained>. Accessed: 2023-04-19.
  - [19] Hierarchical clustering. [https://en.wikipedia.org/wiki/Hierarchical\\_clustering](https://en.wikipedia.org/wiki/Hierarchical_clustering), 2023. Accessed: 2023-04-22.
  - [20] Different types of Clustering in Machine Learning. <https://vitalflux.com/different-types-clustering-algorithm-machine-learning/>, 2022. Accessed: 2023-04-22.

## Α' Τύποι για QDA, LDA

Στον αλγόριθμο QDA ο παράγοντας  $P(X|y)$  δίνεται από τον τύπο:

$$P(X|y = k) = \frac{1}{(2\pi)^{\frac{n}{2}} |\sum_k|^{\frac{1}{2}}} e^{-\frac{1}{2}(X-M_k)^T \sum_k^{-1} (X-M_k)}$$

όπου:

$k$  η κλάση που εξετάζουμε

$X$  το διάνυσμα χαρακτηριστικών του δείγματος

$M_k$  το κέντρο των σημείων της κλάσης  $k$

$n$  η διάσταση του  $X$

$\sum_k$  ο πίνακας συνδιακύμανσης

και

$$X = \langle x_1, x_2, \dots, x_n \rangle$$

αν τα σημεία που ανήκουν στην κλάση  $k$  είναι  $K_1, K_2, \dots, K_m$  με:

$$K_i = \langle k_{i_1}, k_{i_2}, \dots, k_{i_n} \rangle$$

τότε το  $M$  για αυτή την κλάση θα είναι:

$$M = \langle m_1, m_2, \dots, m_n \rangle = \left\langle \frac{1}{m} \sum_{i=1}^m k_{i_1}, \frac{1}{m} \sum_{i=1}^m k_{i_2}, \dots, \frac{1}{m} \sum_{i=1}^m k_{i_n} \right\rangle$$

Η διακύμανση είναι μας δείχνει τη διασπορά των σημείων της κλάσης από το κέντρο της κλάσης για μια συγκεκριμένη διάσταση και υπολογίζεται για κάθε διάσταση ξεχωριστά. Η διακύμανση για μια διάσταση  $u$  από τις  $n$  διαστάσεις θα είναι:

$$V_u = \frac{1}{m} \sum_{i=1}^m (k_{i_u} - m_u)^2$$

Πολύ παρόμοια είναι και η συνδιακύμανση η οποία όμως υπολογίζεται για δύο διαστάσεις  $u$  και  $z$ :

$$C_{uz} = \frac{1}{m} \sum_{i=1}^m (k_{i_u} - m_u) (k_{i_z} - m_z)$$

Ο πίνακας συνδιακύμανσης είναι ένας συμμετρικός πίνακας με διαστάσεις  $n \times n$  και είναι:

$$\sum = \begin{bmatrix} V_1 & C_{12} & \dots & C_{1n} \\ C_{12} & V_2 & \dots & C_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ C_{1n} & C_{2n} & \dots & V_n \end{bmatrix}$$

Ο αλγόριθμος LDA είναι μια απλοποίηση του QDA όπου θεωρούμε ότι όλες οι κλάσεις έχουν τον ίδιο πίνακα συνδιακύμανσης το οποίο απλοποιεί πάρα πολύ τις πράξεις.