

# Classifiers !

Ορισμοί:

Έστω ότι έχουμε  $\vec{x}^2$  είσοδος του αλγορίθμου.

- $x_j$  είναι  $j$ -th στοιχείο του  $\vec{x}$
- $x_j$  είναι η τιμή του διανύσματος για  $j = 1, \dots, d$  όπου  $d$  οι διαστάσεις του διανύσματος.

Έστω  $f(\vec{x})$  η τιμή για το target της εισόδου  $\vec{x}$

- η ακριβής μορφή της  $f(x)$  είναι άγνωστη
- χρησιμοποιούμε δεδομένα ενός συνόλου  $D = \{ \vec{x}, f(\vec{x}) \}$  το οποίο είναι γνωστό ώστε να προσπαθίσουμε να βρούμε μία καλή/επαρκή αναπαράσταση της πραγματικής  $f(\vec{x})$

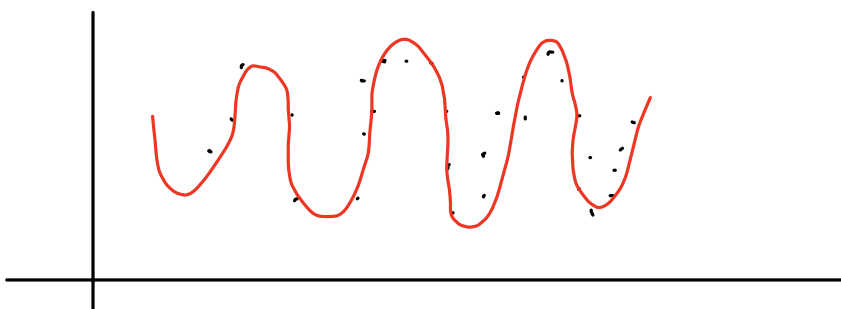
Σκοπός

δημιουργία ενός mapping από το  $\vec{x}$  στο  $f(\vec{x})$

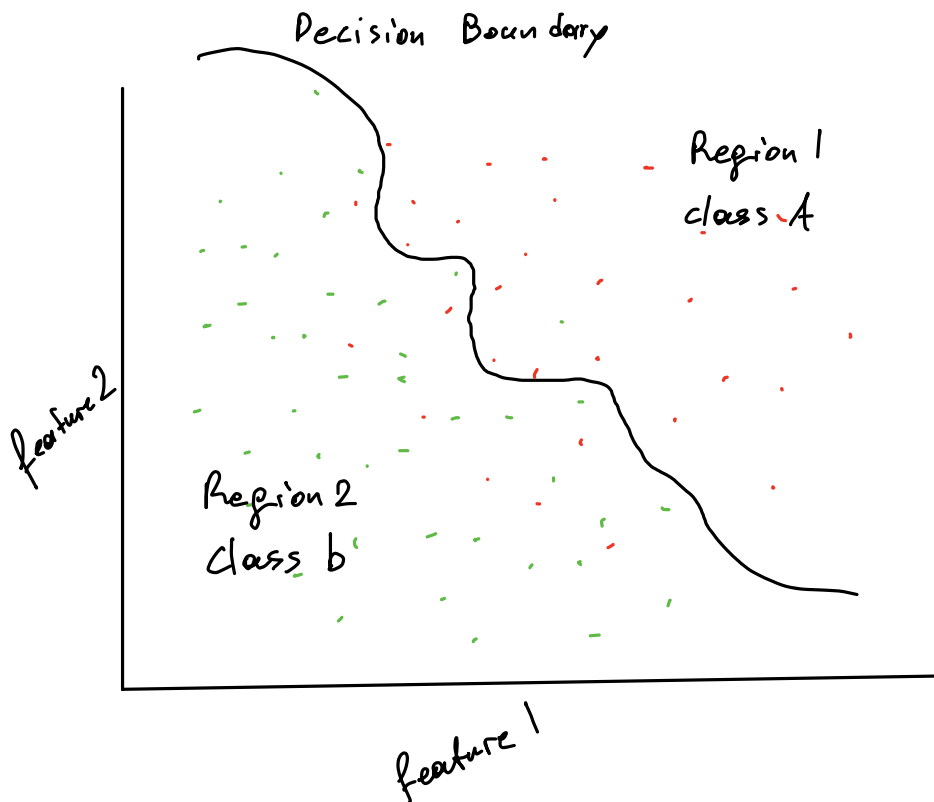
$$h(\vec{x}, \theta) \approx f(\vec{x}) \quad \forall x_j : j \in \{1 \dots d\}$$

$\theta$  είναι οι παράμετροι που έχουμε βρει μέχρι τώρα από τη διαδικασία του classification.

Παράδειγμα



Χρήση σε Διακριτά Προβλήματα :

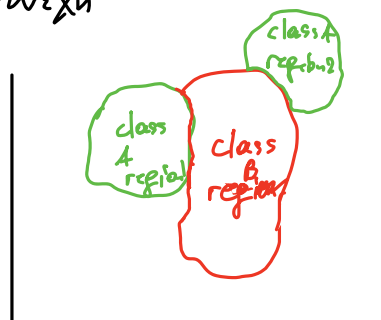


Δεδομένα να βρούμε έναν  
κανόνα / χαρτογράφηση

ο οποίος δίνοντας του  
τις τιμές των feature 1  
και feature 2 αντίστοιχα  
να μπορεί να μας πει  
αν περιλαμβάνει να δώσει  
μόακινη κομμάτια (κλάση  
A) ή πράσινη κομμάτια  
(κλάση B)

## Classification in Euclidean Space

- Η διαδικασία του classification είναι ο χωρισμός ενός  
ευκλείδειου χώρου  $N$ -διάστατων σε επιμέρους χωρία τα οποία:
  - κάθε χωρίο έχει διαφορετικό label (αντιστοιχεί  
δηλαδή σε διαφορετικό class)
  - χωρία τα οποία αντιστοιχούν στο ίδιο label δεν πρέπει  
να είναι συνεχή



- για κάθε καινούρια είσοδο  $\vec{x}_{new}$  θα πρέπει ο classifier να εντοφίσει αυτή το region στο οποίο βρίσκεται το  $\vec{x}_{new}$

- Decision boundaries: τα όρια μεταξύ των διάφορων regions

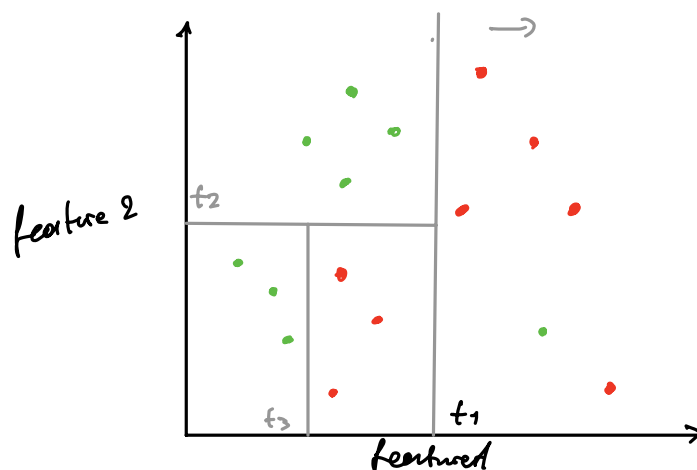
Ο κάθε classifier χαρακτηρίζεται από τις εξισώσεις των decision boundaries που δημιουργεί

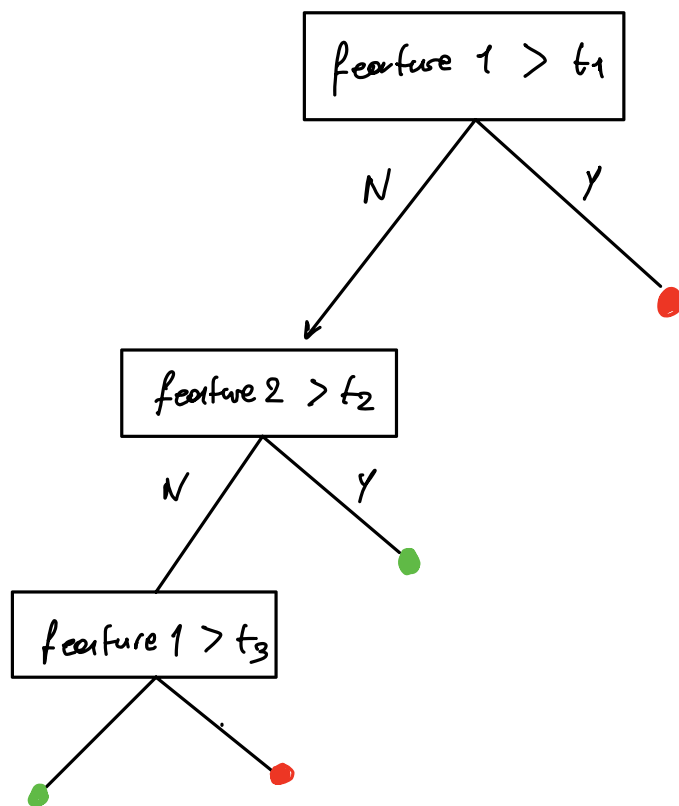
### Αντί Εφαρμογή Classifier: Decision Trees

Το decision tree είναι μία αντί μορφή classifier η οποία χρησιμοποιώντας ευθείες γραμμές παράλληλες στους άξονες  $x, y, z, \dots$  χωρίζει τον  $\mathbb{R}^n$  χώρο που έχουμε σε χωρία, σε κάθε ένα τα των οποίων αναθέτει μία τιμή κλάσης (label)

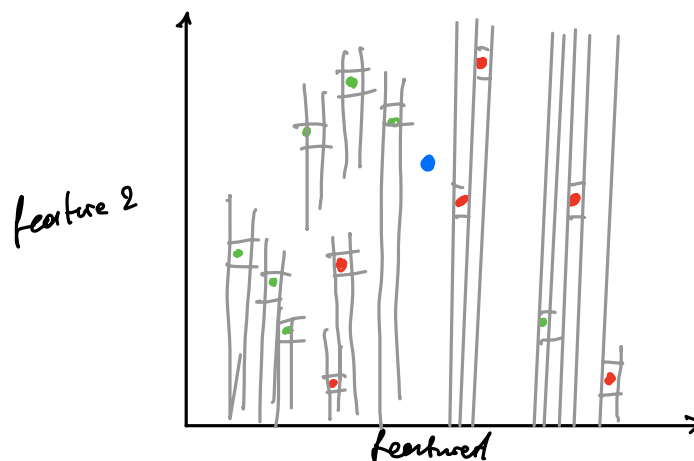
Ο τρόπος αναίρεσης στην περίπτωση του decision tree είναι μιας μορφής:  $x_j > t$ , όπου το  $t$  είναι ένα optimizable arbitrary value

Η διαδικασία αυτή μπορεί να δημιουργήσει μέχρι  $N-1$  χωρία για  $N$  δεδομένα.





Overfitting problem



training accuracy: 100%

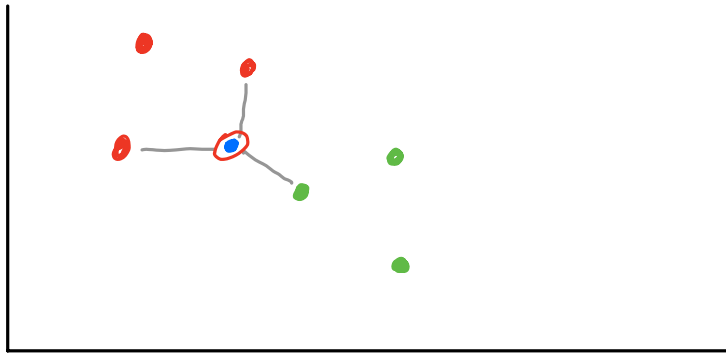
testing accuracy: very low

Thumb rule: τα χυρία δεν πρέπει να είναι πάνω από του  
 να χρησιμοποιείται για τον ~~απο~~ έλεγχο του έργο

# Classifier k-Nearest Neighbors

Υπόθεση: similar inputs give similar outputs

Κανόνας: για ένα test input  $x$ , δίνουμε την πιο κοντινή τιμή  
των neighboring datapoints



Μαθηματικός ορισμός kNN

- test point  $x$
- ορίζεται ένα σύνολο  $k$  στοιχείων τα οποία έχουν την ελάχιστη Ευκλείδεια απόσταση από το  $x$ ,  $S_x$

$$S_x \subset D : |S_x| = k$$

↓  
το σύνολο  
των αρχικών  
datapoint

↳ Δίνει επίσης ότι

δίνεται να χρησιμοποιήσε  $k$  γείτονες

$h > 3$

$\text{min1} = +\infty$

$\text{min2} = +\infty$

$\text{min3} = +\infty$

for  $i \in 1..|D|$

if  $\text{dist}(x_j, x_i) < \text{min1}$

$\text{min1} \leftarrow x_1$

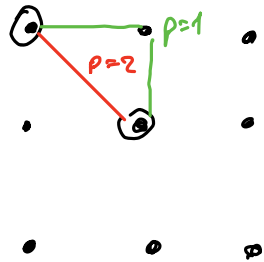
:

### Distance function

$$\text{dist}(x, z) = \left( \sum_{r=1}^d |x_r - z_r|^p \right)^{1/p}$$

$p=2$  : Euclidean distance

$p=1$  : Manhattan distance



# Principal Component Analysis - PCA

$y$	$f_1$	$f_2$	$f_3$	$\dots$	$f_{145}$	$\dots$	$f_d$
$y_1$							
$y_2$							
$\vdots$							
$y_n$							

διδούμε γενικά  $d \leq \sqrt{n}$

1. data matrix  $X$   $n \times d$

2. κανονικοποιήσω την data

$$X'_{ij} = \frac{X_{ij} - \overset{\rightarrow \text{mean}}{\mu_j}}{\underset{\rightarrow \text{std}}{\sigma_j}}$$

3. Covariance Matrix  $\Sigma$

$$\Sigma = \frac{1}{n-1} X'^T X'$$

4. ιδιοτιμές & ιδιοδιανύσματα

$$\Sigma \vec{v} = \lambda \vec{v}$$

$\lambda$ : ιδιοτιμή που αντιστοιχεί στο ιδιοδιάνυσμα  $\vec{v}$

5. Ταξινομήστε τις ιδιοτιμές από τη μεγαλύτερη στη μικρότερη

$y$	PC 1	PC 2	...	PC d
$y_1$	0,8	0,05		0,00004

$\underbrace{\hspace{10em}}_{+1}$