

# Ταξινόμηση I

Οι διαφάνειες στηρίζονται στο P.-N. Tan, M. Steinbach, V. Kumar, «Introduction to Data Mining», Addison Wesley, 2006



## Εισαγωγή

**Ταξινόμηση (classification)**  
 Το γενικό πρόβλημα της ανάθεσης ενός αντικειμένου σε μια ή περισσότερες προκαθορισμένες κατηγορίες (κλάσεις)

**Παραδείγματα**

- Εντοπισμός spam emails, με βάση πχ την επικεφαλίδα τους ή το περιεχόμενό τους
- Πρόβλεψη καρκινικών κυττάρων χαρακτηρίζοντας τα ως καλοήθη ή κακοήθη
- Κατηγοριοποίηση συναλλαγών με πιστωτικές κάρτες ως νόμιμες ή προϊόν απάτης
- Κατηγοριοποίηση δευτερευόντων δομών πρωτεϊνών ως alpha-helix, beta-sheet, ή random coil
- Χαρακτηρισμός ειδήσεων ως οικονομικές, αθλητικές, πολιτιστικές, πρόβλεψης καιρού, κλπ

Εφόδια Διδασκίντων: Ακ. Έτος 2007-2008 ΤΑΞΙΝΟΜΗΣΗ 2

## Ορισμός

Είσοδος: συλλογή από εγγραφές  
 Κάθε εγγραφή περιέχει ένα σύνολο από γνωρίσματα (attributes)  
 Ένα από τα γνωρίσματα είναι η κλάση (class)

Βρες (έξοδος) ένα μοντέλο (model) για το γνώρισμα κλάση ως μια συνάρτηση των τιμών των άλλων γνωρισμάτων

Στόχος: νέες εγγραφές θα πρέπει να ανατίθενται σε μία από τις κλάσεις με τη μεγαλύτερη δυνατή ακρίβεια.

Κατηγορία, Κατηγορία, συνεχές, κλάση

Tid	Επιστροφή	Οικογενειακή Κατάσταση	Φορολογητέο Εσοδόμο	Απάτη
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Εφόδια Διδασκίντων: Ακ. Έτος 2007-2008 ΤΑΞΙΝΟΜΗΣΗ 3

## Ορισμός

Είσοδος: συλλογή από εγγραφές  
 Κάθε εγγραφή περιέχει ένα σύνολο από γνωρίσματα (attributes)  
 Ένα από τα γνωρίσματα είναι η κλάση (class)

Βρες ένα μοντέλο (model) για το γνώρισμα κλάση ως μια συνάρτηση των τιμών των άλλων γνωρισμάτων

Στόχος: νέες εγγραφές θα πρέπει να ανατίθενται σε μία κλάση με τη μεγαλύτερη δυνατή ακρίβεια.

Ταξινόμηση είναι η διαδικασία εκμάθησης μιας συνάρτησης στόχου (target function)  $f$  που απεικονίζει κάθε σύνολο γνωρισμάτων  $x$  σε μια από τις προκαθορισμένες ετικέτες κλάσεις  $y$ .

Συνήθως το σύνολο δεδομένων εισόδου χωρίζεται σε:

- ένα σύνολο εκπαίδευσης (training set) και
- ένα σύνολο ελέγχου (test set)

Το σύνολο εκπαίδευσης χρησιμοποιείται για να κατασκευαστεί το μοντέλο και το σύνολο ελέγχου για να το επικυρωθεί.

Σύνολο εγγραφών (x) → Μοντέλο Ταξινόμησης → Ετικέτα κλάσης (y)

Εφόδια Διδασκίντων: Ακ. Έτος 2007-2008 ΤΑΞΙΝΟΜΗΣΗ 4

## Εισαγωγή

Χρησιμοποιείται ως:

- Περιγραφικό μοντέλο (descriptive modeling): ως επεξηγηματικό εργαλείο - πχ ποια χαρακτηριστικά κάνουν ένα ζώο να χαρακτηριστεί ως θηλαστικό
- Μοντέλο πρόβλεψης (predictive modeling): για τη πρόβλεψη της κλάσης άγνωστων εγγραφών - πχ δοσμένων των χαρακτηριστικών κάποιου ζώου να προβλέψουμε αν είναι θηλαστικό, πτηνό, ερπετό ή αμφίβιο

Κατάλληλη κυρίως για:

- δισαδικές κατηγορίες ή κατηγορίες για τις οποίες δεν υπάρχει διάταξη διακριτές (nominal) vs διατεταγμένες (ordinal)
- για μη ιεραρχικές κατηγορίες

Η τιμή (ετικέτα) της κλάσης -  $y$  - είναι διακριτή τιμή - Διαφορά από regression (οπισθοδρόμηση) όπου το γνώρισμα  $y$  παίρνει συνεχείς τιμές

Εφόδια Διδασκίντων: Ακ. Έτος 2007-2008 ΤΑΞΙΝΟΜΗΣΗ 5

## Τεχνικές Ταξινόμησης

Βήματα Ταξινόμησης

Σύνολο Εκπαίδευσης

Tid	Attrib1	Attrib2	Attrib3	Class
1	No	Small	50K	?
2	No	Medium	80K	?
3	Yes	Large	110K	?
4	No	Small	95K	?
5	No	Large	95K	?
6	No	Medium	60K	?
7	Yes	Large	220K	?
8	No	Small	85K	?
9	No	Medium	75K	?
10	No	Small	90K	?

Επαγωγή Induction → Αλγόριθμος Μάθησης → Κατασκευή Μοντέλου

Αφαίρεση Deduction → Εφαρμογή Μοντέλου

Μοντέλο

Σύνολο Ελέγχου

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	50K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	87K	?

• Ταίριαζει δεδομένα εκπαίδευσης  
 • Προβλέπει την κλάση των δεδομένων ελέγχου  
 • Καλή δυνατότητα γενίκευσης

Εφόδια Διδασκίντων: Ακ. Έτος 2007-2008 ΤΑΞΙΝΟΜΗΣΗ 6

Τεχνικές ταξινόμησης

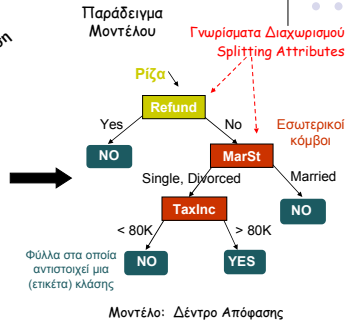
- Τεχνικές βασισμένες σε Δέντρα Απόφασης (decision trees)
- Τεχνικές βασισμένες σε Κανόνες (Rule-based Methods)
- Memory based reasoning
- Νευρωνικά Δίκτυα
- Naive Bayes and Bayesian Belief Networks
- Support Vector Machines

Δέντρα Απόφασης

Δέντρο Απόφασης: Παράδειγμα

Δεδομένα Εκπαίδευσης

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Μοντέλο: Δέντρο Απόφασης

Δέντρο Απόφασης: Παράδειγμα

Μοντέλο = Δέντρο Απόφασης

- Εσωτερικοί κόμβοι αντιστοιχούν σε κάποιο γνώρισμα
- Διαχωρισμός (split) ενός κόμβου σε παιδιά
  - η ετικέτα στην ακμή = συνθήκη/έλεγχος
- Φύλλα αντιστοιχούν σε κλάσεις

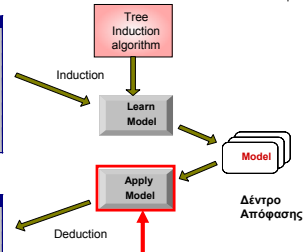
Δέντρο Απόφασης: Βήματα

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	95K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set

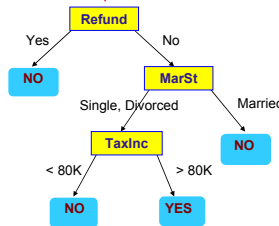


Δέντρο Απόφασης: Εφαρμογή Μοντέλου

Ξεκίνα από τη ρίζα του δέντρου.

Δεδομένα Ελέγχου

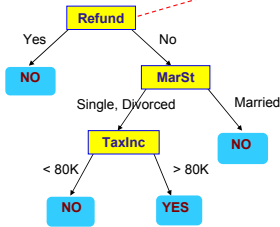
Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



### Δέντρο Απόφασης: Εφαρμογή Μοντέλου

Test Data

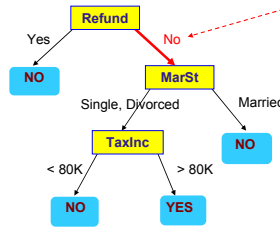
Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



### Δέντρο Απόφασης: Εφαρμογή Μοντέλου

Test Data

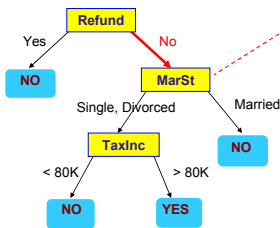
Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



### Δέντρο Απόφασης: Εφαρμογή Μοντέλου

Test Data

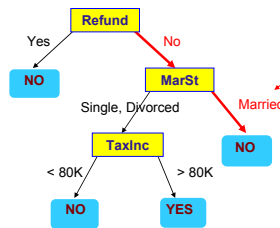
Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



### Δέντρο Απόφασης: Εφαρμογή Μοντέλου

Test Data

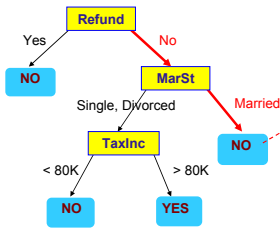
Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



### Δέντρο Απόφασης: Εφαρμογή Μοντέλου

Είσοδος (δεδομένο ελέγχου)

Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?



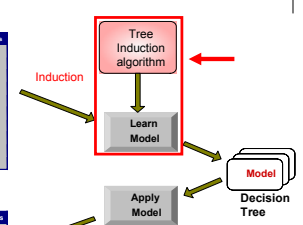
### Δέντρο Απόφασης

Id	Attr01	Attr02	Attr03	Class
1	Yes	Large	120K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	90K	Yes
6	No	Medium	80K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	65K	Yes

Training Set

Id	Attr01	Attr02	Attr03	Class
11	No	Small	85K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	87K	?

Test Set



### Δέντρο Απόφασης: Παράδειγμα

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Για το ίδιο σύνολο εκπαίδευσης υπάρχουν διαφορετικά δέντρα

Εύρωπη Δεδομένων: Ακ. Έτος 2007-2008 ΤΑΞΙΝΟΜΗΣΗ 19

### Δέντρο Απόφασης: Κατασκευή

Ο αριθμός των πιθανών Δέντρων Απόφασης είναι εκθετικός.

Πολλοί αλγόριθμοι για την **επαγωγή (induction)** του δέντρου οι οποίοι ακολουθούν μια **greedy** στρατηγική: για να κτίσουν το δέντρο απόφασης παίρνοντας μια σειρά από **τοπικά βέλτιστες** αποφάσεις

- Hunt's Algorithm (από τους πρώτους)
- CART
- ID3, C4.5
- SLIQ, SPRINT

Εύρωπη Δεδομένων: Ακ. Έτος 2007-2008 ΤΑΞΙΝΟΜΗΣΗ 20

### Δέντρο Απόφασης: Αλγόριθμος του Hunt

Κτίζει το δέντρο **αναδρομικά**

$D_t$ : το σύνολο των εγγραφών εκπαίδευσης που έχουν φτάσει στον κόμβο  $t$

Γενική Διαδικασία:

- Αν το  $D_t$  περιέχει εγγραφές που ανήκουν στην **ίδια** κλάση  $y_t$ , τότε ο κόμβος  $t$  είναι κόμβος φύλλο με ετικέτα  $y_t$
- Αν  $D_t$  είναι το **κενό σύνολο**, τότε ο κόμβος  $t$  είναι κόμβος φύλλο με ετικέτα την default κλάση,  $y_d$
- Αν το  $D_t$  περιέχει εγγραφές που ανήκουν σε περισσότερες από μία κλάσεις, χρησιμοποιήσε έναν **έλεγχο-γνωρίματος** για το διαχωρισμό των δεδομένων σε μικρότερα υποσύνολα. Εφάρμοσε την Διαδικασία αναδρομικά σε κάθε υποσύνολο.

Εύρωπη Δεδομένων: Ακ. Έτος 2007-2008 ΤΑΞΙΝΟΜΗΣΗ 21

### Δέντρο Απόφασης: Αλγόριθμος του Hunt

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Εύρωπη Δεδομένων: Ακ. Έτος 2007-2008 ΤΑΞΙΝΟΜΗΣΗ 22

### Δέντρο Απόφασης: Αλγόριθμος του Hunt

Γενική Διαδικασία (πιο αναλυτικά):

- Αν το  $D_t$  περιέχει εγγραφές που ανήκουν στην **ίδια** κλάση  $y_t$ , τότε ο κόμβος  $t$  είναι κόμβος φύλλο με ετικέτα  $y_t$
- Αν  $D_t$  είναι το **κενό σύνολο**, αυτό σημαίνει ότι δεν υπάρχει εγγραφή στο σύνολο εκπαίδευσης με αυτό το συνδυασμό τιμών, τότε  $D_t$  γίνεται φύλλο με κλάση αυτή της **πλειοψηφίας** των εγγραφών εκπαίδευσης ή ανάθεση κάποιας default κλάσης
- Αν το  $D_t$  περιέχει εγγραφές που ανήκουν σε περισσότερες από μία κλάσεις, χρησιμοποιήσε έναν έλεγχο-γνωρίματος για το διαχωρισμό των δεδομένων σε μικρότερα υποσύνολα. Εφάρμοσε την Διαδικασία αναδρομικά σε κάθε υποσύνολο.

Το παραπάνω δεν είναι δυνατόν αν όλες οι εγγραφές έχουν τις ίδιες τιμές σε όλα τα γνωρίσματα (δηλαδή, ο ίδιος συνδυασμός αντιστοιχεί σε περισσότερες από μία κλάσεις) τότε φύλλο με κλάση αυτής της πλειοψηφίας των εγγραφών εκπαίδευσης

Εύρωπη Δεδομένων: Ακ. Έτος 2007-2008 ΤΑΞΙΝΟΜΗΣΗ 23

### Δέντρο Απόφασης: Κατασκευή Δέντρου

*Πως θα γίνει η διάσπαση του κόμβου:*

**Greedy** στρατηγική.

- Διάσπαση εγγραφών με βάση έναν έλεγχο γνωρίματος που βελτιστοποιεί ένα συγκεκριμένο **κριτήριο**
- **Θέματα**
  - Καθορισμός του τρόπου διαχωρισμού των εγγραφών
    - Καθορισμός του ελέγχου γνωρίματος
    - Καθορισμός του βέλτιστου διαχωρισμού
  - Πότε θα σταματήσει ο διαχωρισμός (συνθήκη τερματισμού)

Εύρωπη Δεδομένων: Ακ. Έτος 2007-2008 ΤΑΞΙΝΟΜΗΣΗ 24

## Δέντρο Απόφασης: Κατασκευή Δέντρου

Καθορισμός των συνθηκών του ελέγχου για τα γνωρίσματα

- Εξαρτάται από τον τύπο των γνωρισμάτων
  - Διακριτές - Nominal
  - Διατεταγμένες - Ordinal
  - Συνεχείς - Continuous
- Εξαρτάται από τον αριθμό των διαφορετικών τρόπων διάσπασης
  - 2-αδική διάσπαση - 2-way split
  - Πολλαπλή διάσπαση - Multi-way split

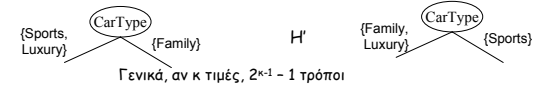
## Δέντρο Απόφασης: Κατασκευή Δέντρου

Διαχωρισμός βασισμένος σε διακριτές τιμές

- Πολλαπλός διαχωρισμός:** Χρησιμοποιήσει τόσες διασπάσεις όσες οι διαφορετικές τιμές

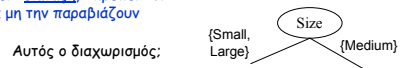


- Διαδικός Διαχωρισμός:** Χωρίζει τις τιμές σε δύο υποσύνολα. Πρέπει να βρει το βέλτιστο διαχωρισμό (partitioning).



Γενικά, αν κ τιμές,  $2^{k-1} - 1$  τρόποι

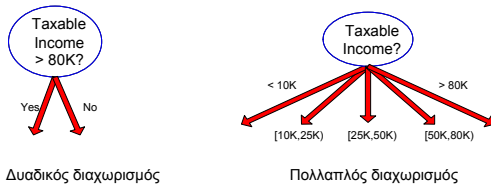
Όταν υπάρχει διάταξη, πρέπει οι διασπάσεις να μη την παραβιάζουν



Αυτός ο διαχωρισμός:

## Δέντρο Απόφασης: Κατασκευή Δέντρου

Διαχωρισμός βασισμένος σε συνεχείς τιμές



Διαδικός διαχωρισμός

Πολλαπλός διαχωρισμός

## Δέντρο Απόφασης: Κατασκευή Δέντρου

Διαχωρισμός βασισμένος σε συνεχείς τιμές

Τρόποι χειρισμού

- Discretization (διακριτοποίηση)** ώστε να προκύψει ένα διατεταγμένο κατηγορικό γνώρισμα
  - Ταξινόμηση των τιμών και χωρισμός τους σε περιοχές καθορίζοντας  $n - 1$  σημεία διαχωρισμού, απεικόνιση όλων των τιμών μιας περιοχής στην ίδια κατηγορική τιμή
  - Στατικό** - μια φορά στην αρχή
  - Δυναμικό** - εύρεση των περιοχών πχ έτσι ώστε οι περιοχές να έχουν το ίδιο διάστημα ή τις ίδιες συχνότητες εμφάνισης ή με χρήση συσταδοποίησης

- Διαδική Απόφαση:** ( $A < v$ ) or ( $A \geq v$ ) εξετάζει όλους τους δυνατούς διαχωρισμούς (τιμές του  $v$ ) και επιλέγει τον καλύτερο - υπολογιστικά βαρύ

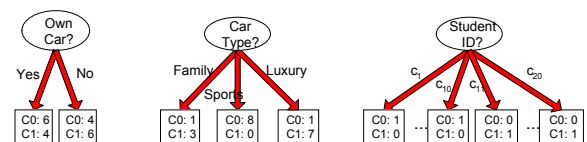
## Δέντρο Απόφασης: Κατασκευή Δέντρου

- Greedy** στρατηγική.
  - Διάσπαση εγγραφών με βάση έναν έλεγχο γνωρίματος που βελτιστοποιεί ένα συγκεκριμένο κριτήριο
- Θέματα**
  - Καθορισμός του τρόπου διαχωρισμού των εγγραφών
    - Καθορισμός του ελέγχου γνωρίματος
    - Καθορισμός του βέλτιστου διαχωρισμού**
- Πότε θα σταματήσει ο διαχωρισμός (συνθήκη τερματισμού)

## Δέντρο Απόφασης: Κατασκευή Δέντρου

Βέλτιστος Διαχωρισμός

Πριν το διαχωρισμό: 10 εγγραφές της κλάσης 0, 10 εγγραφές της κλάσης 1



Ποια από τις 3 διασπάσεις να προτιμήσουμε;

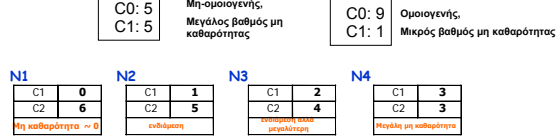
Ποια συνθήκη ελέγχου είναι καλύτερη;

## Δέντρο Απόφασης: Κατασκευή Δέντρου

### Βέλτιστος Διαχωρισμός

Greedy προσέγγιση: προτιμούνται οι κόμβοι με ομοιογενείς κατανομές κλάσεων (homogeneous class distribution)

Χρειαζόμαστε μία μέτρηση της μη καθαρότητας ενός κόμβου (node impurity)



$$I(N1) < I(N2) < I(N3) < I(N4)$$

## Δέντρο Απόφασης: Κατασκευή Δέντρου

Πως θα χρησιμοποιήσουμε τη μέτρηση καθαρότητας:

Κριτήριο για διάσπαση - Το τι κερδίζουμε από την διάσπαση:

Έστω ότι έχουμε ένα μέτρο για τη μέτρηση αυτής της καθαρότητας ενός κόμβου  $n$ :  $I(n)$   
Κοιτάμε την καθαρότητα του γονέα (πριν τη διάσπαση) και των παιδιών του (μετά τη διάσπαση)

Βάρος (εξαρτάται από τον αριθμό εγγραφών)

$$\Delta = I(\text{parent}) - \sum_{i=1}^k \frac{N(u_i)}{N} I(u_i)$$

$N$  είναι ο αριθμός των εγγραφών στο γονέα και  $N(u_i)$  του  $i$ -οστού παιδιού

Διαλέγουμε την «καλύτερη» διάσπαση (μεγαλύτερο  $\Delta$ )

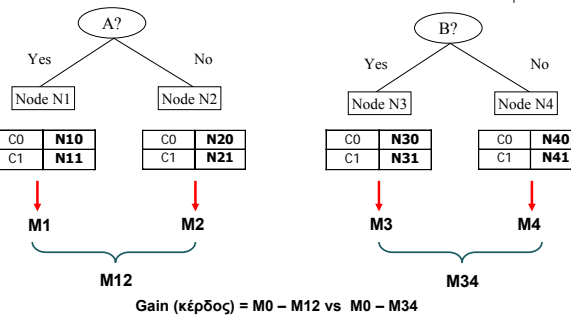
Παράδειγμα

## Δέντρο Απόφασης: Κατασκευή

Πριν τη διάσπαση:

C0	N00
C1	N01

→ M0



## Δέντρο Απόφασης: Κατασκευή Δέντρου

### Μέτρα μη Καθαρότητας

1. Ευρετήριο Gini - Gini Index
2. Εντροπία - Entropy
3. Λάθος ταξινομήσεις - Misclassification error

## Δέντρο Απόφασης: GINI

Ευρετήριο Gini για τον κόμβο  $t$ :

$$GINI(t) = 1 - \sum_{j=1}^c [p(j|t)]^2$$

$p(j|t)$  σχετική συχνότητα της κλάσης  $j$  στον κόμβο  $t$  (ποσοστό εγγραφών της κλάσης  $j$  στον κόμβο  $t$ )  
 $c$  αριθμός κλάσεων

Παράδειγματα:

C1	0	Gini=0.000
C2	6	

C1	1	Gini=0.278
C2	5	

C1	2	Gini=0.444
C2	4	

C1	3	Gini=0.500
C2	3	

## Δέντρο Απόφασης: GINI

Ευρετήριο Gini για τον κόμβο  $t$ :

$$GINI(t) = 1 - \sum_{j=1}^c [p(j|t)]^2$$

$p(j|t)$  σχετική συχνότητα της κλάσης  $j$  στον κόμβο  $t$   
 $c$  αριθμός κλάσεων

Ελάχιστη τιμή (0.0) όταν όλες οι εγγραφές ανήκουν σε μία κλάση (που σημαίνει την πιο ενδιαφέρουσα πληροφορία)

Μέγιστη τιμή ( $1 - 1/c$ ) όταν όλες οι εγγραφές είναι ομοιόμορφα κατανομημένες στις κλάσεις (που σημαίνει τη λιγότερο ενδιαφέρουσα πληροφορία)  
εξαρτάται από τον αριθμό των κλάσεων

## Δέντρο Απόφασης: GINI

Χρήση του στην κατασκευή του δέντρου απόφασης

- Χρησιμοποιείται στα CART, SLIQ, SPRINT.

Όταν ένας κόμβος  $p$  διασπάται σε  $k$  κόμβους (παιδιά), (που σημαίνει ότι το σύνολο των εγγραφών του κόμβου χωρίζεται σε  $k$  υποσύνολα), η ποιότητα του διαχωρισμού υπολογίζεται ως:

$$GINI_{split} = \sum_{i=1}^k \frac{n_i}{n} GINI(i)$$

όπου,  $n_i$  = αριθμός εγγραφών του παιδιού  $i$ ,  
 $n$  = αριθμός εγγραφών του κόμβου  $p$ .

Ψάχνουμε για:

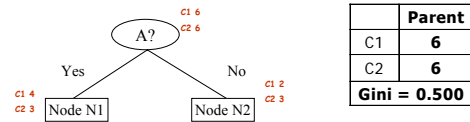
- Πιο καθαρές
- Πιο μεγάλες (σε αριθμό) μικρές διασπάσεις

## Δέντρο Απόφασης: GINI

Παράδειγμα Εφαρμογής

Περίπτωση 1: Διαδικά Γνωρίσματα

Αρχικός κόμβος



	N1	N2
C1	4	2
C2	2	3
Gini	0.486	

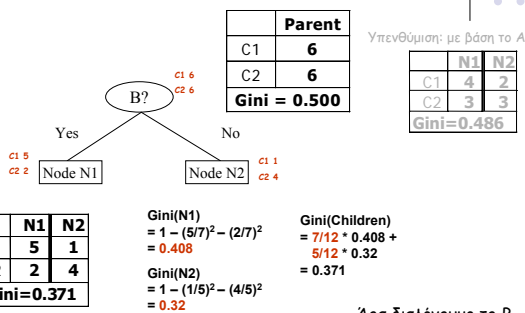
$$Gini(N1) = 1 - (4/7)^2 - (3/7)^2 = 0.49$$

$$Gini(N2) = 1 - (2/5)^2 - (3/5)^2 = 0.48$$

$$Gini(Children) = 7/12 * 0.49 + 5/12 * 0.48 = 0.486$$

## Δέντρο Απόφασης: GINI

Παράδειγμα Εφαρμογής (συνέχεια)



	N1	N2
C1	5	1
C2	2	4
Gini	0.371	

$$Gini(N1) = 1 - (5/7)^2 - (2/7)^2 = 0.408$$

$$Gini(N2) = 1 - (1/5)^2 - (4/5)^2 = 0.32$$

$$Gini(Children) = 7/12 * 0.408 + 5/12 * 0.32 = 0.371$$

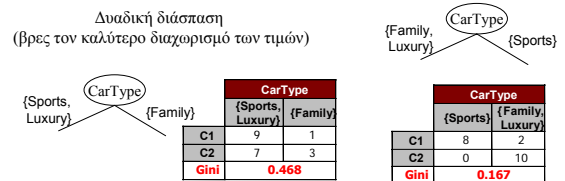
Άρα διαλέγουμε το B

## Δέντρο Απόφασης: GINI

Περίπτωση 2: Κατηγορικά Γνωρίσματα

Για κάθε διαφορετική τιμή, μέτρησε τις τιμές στα δεδομένα που ανήκουν σε κάθε κλάση

Χρησιμοποίησε τον πίνακα με τους μετρητές για να πάρεις την απόφαση



	{Sports, Luxury}	{Family}
C1	9	1
C2	7	3
Gini	0.468	

	{Sports}	{Family, Luxury}
C1	8	2
C2	0	10
Gini	0.167	

## Δέντρο Απόφασης: GINI

Περίπτωση 2: Κατηγορικά Γνωρίσματα



Πολλαπλή Διάσπαση

	CarType		
	Family	Sports	Luxury
C1	1	8	1
C2	3	0	7
Gini	0.163		

## Δέντρο Απόφασης: GINI

Συνεχή Γνωρίσματα

- Χρήση **δυναμικών αποφάσεων** πάνω σε μία τιμή
- Πολλές επιλογές για την τιμή διαχωρισμού
  - Αριθμός πιθανών διαχωρισμών = Αριθμός διαφορετικών τιμών - έστω  $N$
- Κάθε τιμή διαχωρισμού  $v$  συσχετίζεται με έναν πίνακα μετρητών
  - Μετρητές των κλάσεων για κάθε μια από τις δύο διασπάσεις,  $A < v$  and  $A \geq v$
- Απλή μέθοδος για την επιλογή της καλύτερης τιμής  $v$ 
  - Για κάθε  $v$ ,  $scan$  τα δεδομένα κατασκεύασε τον πίνακα και υπολόγισε το  $gini$  ευρετήριο χρόνος  $O(N)$
  - $O(N^2)$  Υπολογιστικά μη αποδοτικό! Επανάληψη υπολογισμού.

Tid	Refund	Marital Status	Taxable Income	Chest
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



## Δέντρο Απόφασης: GINI

- Για ποιο αποδοτικό υπολογισμό, για κάθε γνώρισμα
  - Ταξινόμηση σε γνώρισμα -  $O(N \log N)$
  - Σειριακή διάσχιση των τιμών, ενημερώνοντας κάθε φορά τον πίνακα με τους μετρητές και υπολογίζοντας το ευρετήριο Gini
  - Επιλογή του διαχωρισμού με το μικρότερο ευρετήριο Gini

Παράδειγμα - Διαχωρισμός στο γνώρισμα Income

Taxable Income	No		Yes		No		Yes		No		Yes		No		Yes	
	Count	Percentage	Count	Percentage	Count	Percentage	Count	Percentage	Count	Percentage	Count	Percentage	Count	Percentage	Count	Percentage
60	55	0.300	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000
70	65	0.361	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000
75	72	0.400	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000
80	80	0.444	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000
85	87	0.483	1	0.056	1	0.056	1	0.056	1	0.056	1	0.056	1	0.056	1	0.056
90	92	0.511	2	0.111	2	0.111	2	0.111	2	0.111	2	0.111	2	0.111	2	0.111
95	97	0.539	3	0.167	3	0.167	3	0.167	3	0.167	3	0.167	3	0.167	3	0.167
100	110	0.611	4	0.222	4	0.222	4	0.222	4	0.222	4	0.222	4	0.222	4	0.222
120	122	0.678	5	0.278	5	0.278	5	0.278	5	0.278	5	0.278	5	0.278	5	0.278
125	125	0.704	6	0.333	6	0.333	6	0.333	6	0.333	6	0.333	6	0.333	6	0.333
220	230	0.783	7	0.389	7	0.389	7	0.389	7	0.389	7	0.389	7	0.389	7	0.389
Gini		0.420	0.400	0.375	0.343	0.417	0.400	0.309	0.343	0.375	0.400	0.420				

ID	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Για <55, δεν υπάρχει εγγραφή οπότε 0  
Για <65, κοιτάμε το μικρότερο το 60, NO 0->1, 7->6 YES δεν αλλάζει

Για <72, κοιτάμε το μικρότερο το 70, NO 1->2 6->5, YES δεν αλλάζει

κακ

Καλύτερα: Αγνοούμε τα σημεία στα οποία δεν υπάρχει αλλαγή κλάσης (αυτά δε μπορεί να είναι σημεία διαχωρισμού)

Άρα, στο παράδειγμα, αγνοούνται τα σημεία 55, 65, 72, 87, 92, 122, 172, 230

Από 11 πιθανά σημεία διαχωρισμού μας μένουν μόνο 2

Taxable Income	No		Yes		No		Yes		No		Yes		No		Yes	
	Count	Percentage	Count	Percentage	Count	Percentage	Count	Percentage	Count	Percentage	Count	Percentage	Count	Percentage	Count	Percentage
60	55	0.300	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000
70	65	0.361	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000
75	72	0.400	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000
80	80	0.444	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000	0	0.000
85	87	0.483	1	0.056	1	0.056	1	0.056	1	0.056	1	0.056	1	0.056	1	0.056
90	92	0.511	2	0.111	2	0.111	2	0.111	2	0.111	2	0.111	2	0.111	2	0.111
95	97	0.539	3	0.167	3	0.167	3	0.167	3	0.167	3	0.167	3	0.167	3	0.167
100	110	0.611	4	0.222	4	0.222	4	0.222	4	0.222	4	0.222	4	0.222	4	0.222
120	122	0.678	5	0.278	5	0.278	5	0.278	5	0.278	5	0.278	5	0.278	5	0.278
125	125	0.704	6	0.333	6	0.333	6	0.333	6	0.333	6	0.333	6	0.333	6	0.333
220	230	0.783	7	0.389	7	0.389	7	0.389	7	0.389	7	0.389	7	0.389	7	0.389
Gini		0.420	0.400	0.375	0.343	0.417	0.400	0.309	0.343	0.375	0.400	0.420				

## Δέντρο Απόφασης: Κατασκευή Δέντρου

### Μέτρα μη Καθαρότητας

- Ευρετήριο Gini - Gini Index
- Εντροπία - Entropy
- Λάθος ταξινομήσεις - Misclassification error

## Δέντρο Απόφασης: Εντροπία

Εντροπία για τον κόμβο t :

$$Entropy(t) = -\sum_{j=1}^c p(j|t) \log p(j|t)$$

$p(j|t)$  σχετική συχνότητα της κλάσης j στον κόμβο t  
c αριθμός κλάσεων

Μετράει την ομοιογένεια ενός κόμβου

**Μέγιστη τιμή**  $\log(c)$  όταν όλες οι εγγραφές είναι ομοιόμορφα κατανομημένες στις κλάσεις (που σημαίνει τη λιγότερη ενδιαφέρουσα πληροφορία)

**Ελάχιστη τιμή** (0.0) όταν όλες οι εγγραφές ανήκουν σε μία κλάση (που σημαίνει την πιο ενδιαφέρουσα πληροφορία)

## Δέντρο Απόφασης: Εντροπία

Παράδειγματα

$$Entropy(t) = -\sum_j p(j|t) \log_2 p(j|t)$$

C1	0	P(C1) = 0/6 = 0	P(C2) = 6/6 = 1
C2	6	Entropy = -0 log 0 - 1 log 1 = -0 - 0 = 0	

C1	1	P(C1) = 1/6	P(C2) = 5/6
C2	5	Entropy = -(1/6) log <sub>2</sub> (1/6) - (5/6) log <sub>2</sub> (5/6) = 0.65	

C1	2	P(C1) = 2/6	P(C2) = 4/6
C2	4	Entropy = -(2/6) log <sub>2</sub> (2/6) - (4/6) log <sub>2</sub> (4/6) = 0.92	

## Δέντρο Απόφασης: Εντροπία

Και σε αυτήν την περίπτωση, όταν ένας κόμβος p διασπάται σε k σύνολα (παιδιά), η ποιότητα του διαχωρισμού υπολογίζεται ως:

$$GAIN_{split} = Entropy(p) - \left( \sum_{i=1}^k \frac{n_i}{n} Entropy(i) \right)$$

όπου,  $n_i$  = αριθμός εγγραφών του παιδιού i,  
 $n$  = αριθμός εγγραφών του κόμβου p.

Χρησιμοποιείται στα ID3 and C4.5

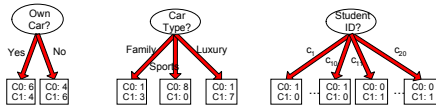
Όταν χρησιμοποιούμε την εντροπία για τη μέτρηση της μη καθαρότητας τότε η διαφορά καλείται **κέρδος πληροφορίας (information gain)**



### Δέντρο Απόφασης

$$\Delta = I(\text{parent}) - \sum_{i=1}^k \frac{N(u_i)}{N} I(u_i)$$

Τείνει να ευνοεί διαχωρισμούς που καταλήγουν σε μεγάλο αριθμό από διασπάσεις που η κάθε μία είναι μικρή αλλά καθαρή



Μπορεί να καταλήξουμε σε πολύ μικρούς κόμβους (με πολύ λίγες εγγραφές) για αξιόπιστες προβλέψεις

Στο παράδειγμα, το student-id είναι κλειδί, όχι χρήσιμο για προβλέψεις

### Δέντρο Απόφασης: Λόγος Κέρδους

- Μία λύση είναι να έχουμε μόνο δυαδικές διασπάσεις
- Εναλλακτικά, μπορούμε να λάβουμε υπό όψιν μας τον αριθμό των κόμβων

$$\text{GainRATIO}_{\text{split}} = \frac{\text{GAIN}_{\text{split}}}{\text{SplitINFO}}$$

Όπου:

$$\text{SplitINFO} = - \sum_{i=1}^k \frac{n_i}{n} \log \frac{n_i}{n}$$

**SplitINFO:** εντροπία της διάσπασης

Μεγάλος αριθμός μικρών διασπάσεων (υψηλή εντροπία) τιμωρείται

Χρησιμοποιείται στο C4.5

### Δέντρο Απόφασης: Λόγος Κέρδους

$$\text{GainRATIO}_{\text{split}} = \frac{\text{GAIN}_{\text{split}}}{\text{SplitINFO}}$$

$$\text{SplitINFO} = - \sum_{i=1}^k \frac{n_i}{n} \log \frac{n_i}{n}$$

Παράδειγμα

Έστω N εγγραφές αν τις χωρίσουμε

Σε 3 κόμβους SplitINFO =  $-\log(1/3) = \log 3$

Σε 2 κόμβους SplitINFO =  $-\log(1/2) = \log 2 = 1$

Άρα οι 2 ευνοούνται

### Δέντρο Απόφασης: Κατασκευή Δέντρου

#### Μέτρα μη Καθαρότητας

- Ευρετήριο Gini - Gini Index
- Εντροπία - Entropy
- Λάθος ταξινομήσεις - Misclassification error

### Δέντρο Απόφασης: Λάθος Ταξινόμησης

Λάθος ταξινόμησης (classification error) για τον κόμβο t :

$$\text{Error}(t) = 1 - \max_{\text{class } i} P(i | t)$$

Μετράει το λάθος ενός κόμβου

Παράδειγμα

C1	0
C2	6

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$\text{Error} = 1 - \max(0, 1) = 1 - 1 = 0$$

C1	1
C2	5

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$\text{Error} = 1 - \max(1/6, 5/6) = 1 - 5/6 = 1/6$$

C1	2
C2	4

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

$$\text{Error} = 1 - \max(2/6, 4/6) = 1 - 4/6 = 1/3$$

### Δέντρο Απόφασης: Λάθος Ταξινόμησης

Λάθος ταξινόμησης (classification error) για τον κόμβο t :

$$\text{Error}(t) = 1 - \max_{\text{class } i} P(i | t)$$

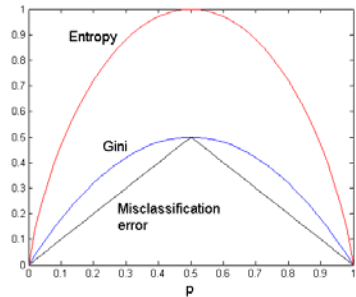
Μετράει το λάθος ενός κόμβου

**Μέγιστη τιμή** 1-1/c όταν όλες οι εγγραφές είναι ομοιόμορφα κατανεμημένες στις κλάσεις (που σημαίνει τη λιγότερο ενδιαφέρουσα πληροφορία)

**Ελάχιστη τιμή** (0.0) όταν όλες οι εγγραφές ανήκουν σε μία κλάση (που σημαίνει την πιο ενδιαφέρουσα πληροφορία)

## Δέντρο Απόφασης: Σύγκριση

Για ένα πρόβλημα δύο κλάσεων



$p$  ποσοστό εγγραφών που ανήκει σε μία από τις δύο κλάσεις ( $p$  κλάση +,  $1-p$  κλάση -)

Όλες την μεγαλύτερη τιμή για 0.5 (ομοιόμορφη κατανομή)

Όλες μικρότερη τιμή όταν όλες οι εγγραφές σε μία μόνο κλάση (0 και στο 1)

## Δέντρο Απόφασης: Σύγκριση

▪ Όπως είδαμε και στα παραδείγματα οι τρεις μετρήσεις είναι συνεπής μεταξύ τους, πχ  $N1$  μικρότερη τιμή από το  $N2$  και με τις τρεις μετρήσεις

▪ Ωστόσο το γνώρισμα που θα επιλεγεί για τη συνθήκη ελέγχου εξαρτάται από το ποια μέτρηση χρησιμοποιείται

## Δέντρο Απόφασης: Κατασκευή Δέντρου

- Greedy στρατηγική.
  - Διάσπαση εγγραφών με βάση έναν έλεγχο γνωρίσματος που βελτιστοποιεί ένα συγκεκριμένο κριτήριο
- Θέματα
  - Καθορισμός του τρόπου διαχωρισμού των εγγραφών
    - Καθορισμός του ελέγχου γνωρίσματος
    - Καθορισμός του βέλτιστου διαχωρισμού
  - **Πότε θα σταματήσει ο διαχωρισμός (συνθήκη τερματισμού)**

## Δέντρο Απόφασης: Κριτήρια Τερματισμού

- Σταματάμε την επέκταση ενός κόμβου όταν όλες οι εγγραφές του ανήκουν στην ίδια κλάση
- Σταματάμε την επέκταση ενός κόμβου όταν όλα τα γνωρίσματα έχουν τις ίδιες τιμές
- Γρήγορος τερματισμός

## Δέντρο Απόφασης

### Πλεονεκτήματα Δέντρων Απόφασης

- Μη παραμετρική προσέγγιση: Δε στηρίζεται σε υπόθεση εκ των προτέρων γνώσης σχετικά με τον τύπο της κατανομής πιθανότητας που ικανοποιεί η κλάση ή τα άλλα γνωρίσματα
- Η κατασκευή του βέλτιστου δέντρου απόφασης είναι ένα NP-complete πρόβλημα. Ευριστικοί: Αποδοτική κατασκευή ακόμα και στην περίπτωση πολύ μεγάλου συνόλου δεδομένων
- Αφού το δέντρο κατασκευαστεί, η ταξινόμηση νέων εγγραφών πολύ γρήγορη  $O(h)$  όπου  $h$  το μέγιστο ύψος του δέντρου
- Εύκολα στην κατανόηση (ιδιαίτερα τα μικρά δέντρα)
- Η ακρίβεια τους συγκρίσιμη με άλλες τεχνικές για μικρά σύνολα δεδομένων

## Δέντρο Απόφασης

### Πλεονεκτήματα

- Καλή συμπεριφορά στο θόρυβο
- Η ύπαρξη πλεοναζόντων γνωρισμάτων (γνωρίσματα των οποίων η τιμή εξαρτάται από κάποιο άλλο) δεν είναι καταστροφική για την κατασκευή. Χρησιμοποιείται ένα από τα δύο. Αν πάρα πολλά, μπορεί να οδηγήσουν σε δέντρα πιο μεγάλα από ότι χρειάζεται

## Δέντρο Απόφασης

### Στρατηγική αναζήτησης

- Ο αλγόριθμος που είδαμε χρησιμοποιεί μια greedy, top-down, αναδρομική διάσπαση για να φτάσει σε μια αποδεκτή λύση
- Άλλες στρατηγικές?
  - Bottom-up (από τα φύλλα, αρχικά κάθε εγγραφή και φύλλο)
  - Bi-directional

## Δέντρο Απόφασης

### Εκφραστικότητα

- Δυνατότητα αναπαράστασης για συναρτήσεις διακριτών τιμών, αλλά δε δουλεύουν σε κάποια είδη δυαδικών προβλημάτων - πχ, parity  $O(1)$  αν υπάρχει μονός (ζυγός) αριθμός από δυαδικά γνωρίσματα  $2^d$  κόμβοι για  $d$  γνωρίσματα
- Όχι καλή συμπεριφορά για συνεχείς μεταβλητές  
Ιδιαίτερα όταν η συνθήκη ελέγχου αφορά ένα γνώρισμα τη φορά

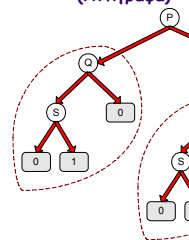
## Δέντρο Απόφασης

### Data Fragmentation - Διάσπαση Δεδομένων

- Ο αριθμός των εγγραφών μειώνεται όσο κατεβαίνουμε στο δέντρο
- Ο αριθμός των εγγραφών στα φύλλα μπορεί να είναι πολύ μικρός για να πάρουμε οποιαδήποτε στατιστικά σημαντική απόφαση
- Μπορούμε να αποτρέψουμε την περαιτέρω διάσπαση όταν ο αριθμός των εγγραφών πέσει κάτω από ένα όριο

## Δέντρο Απόφασης

### Tree Replication (Αντίγραφα)



Το ίδιο υπο-δέντρο να εμφανίζεται πολλές φορές σε ένα δέντρο απόφασης  
Αυτό κάνει το δέντρο πιο περίπλοκο και πιθανών δυσκολότερο στην κατανόηση

Σε περιπτώσεις διάσπασης ενός γνωρίσματος σε κάθε εσωτερικό κόμβο - ο ίδιος έλεγχος σε διαφορετικά σημεία

## Δέντρο Απόφασης

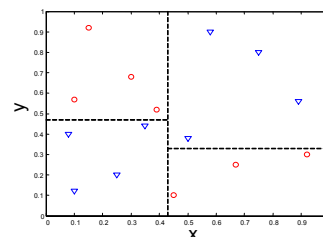
### Decision Boundary

Μέχρι στιγμής είδαμε ελέγχους που αφορούν μόνο ένα γνώρισμα τη φορά, μπορούμε να δούμε τη διαδικασία ως τη διαδικασία *διαμερισμού του χώρου* των γνωρισμάτων σε ξένες περιοχές μέχρι κάθε περιοχή να περιέχει εγγραφές που να ανήκουν στην ίδια κλάση

Η οριακή γραμμή (Border line) μεταξύ δυο γειτονικών περιοχών που ανήκουν σε διαφορετικές κλάσεις ονομάζεται και **decision boundary** (όριο απόφασης)

## Δέντρο Απόφασης

Όταν η συνθήκη ελέγχου περιλαμβάνει μόνο ένα γνώρισμα τη φορά τότε το Decision boundary είναι παράλληλη στους άξονες (τα decision boundaries είναι ορθογώνια παραλληλόγραμμα)



### Δέντρο Απόφασης

Οβlique (πλάγιο) Δέντρο Απόφασης

$x + y < 1$

Class = +      Class = ●

- Οι συνθήκες ελέγχου μπορούν να περιλαμβάνουν περισσότερα από ένα γνωρίσματα
- Μεγαλύτερη εκφραστικότητα
- Η εύρεση βέλτιστων συνθηκών ελέγχου είναι υπολογιστικά ακριβή

Εξόρυξη Δεδομένων: Ακ. Έτος 2007-2008      ΤΑΞΙΝΟΜΗΣΗ      67

### Δέντρο Απόφασης

Constructive induction

Κατασκευή σύνθετων γνωρισμάτων ως αριθμητικών ή λογικών συνδυασμών άλλων γνωρισμάτων

Εξόρυξη Δεδομένων: Ακ. Έτος 2007-2008      ΤΑΞΙΝΟΜΗΣΗ      68

### Δέντρο Απόφασης: C4.5

- Simple depth-first construction.
- Uses Information Gain
- Sorts Continuous Attributes at each node.
- Needs entire data to fit in memory.
- Unsuitable for Large Datasets.
- Needs out-of-core sorting.

You can download the software from:  
<http://www.cse.unsw.edu.au/~quinlan/c4.5r8.tar.gz>

Εξόρυξη Δεδομένων: Ακ. Έτος 2007-2008      ΤΑΞΙΝΟΜΗΣΗ      69

## Θέματα στην Ταξινόμηση

Εξόρυξη Δεδομένων: Ακ. Έτος 2007-2008      ΤΑΞΙΝΟΜΗΣΗ      70

### Θέματα Ταξινόμησης

- Underfitting and Overfitting
- Εκτίμηση Λάθους
- Τιμές που λείπουν

Εξόρυξη Δεδομένων: Ακ. Έτος 2007-2008      ΤΑΞΙΝΟΜΗΣΗ      71

### Overfitting

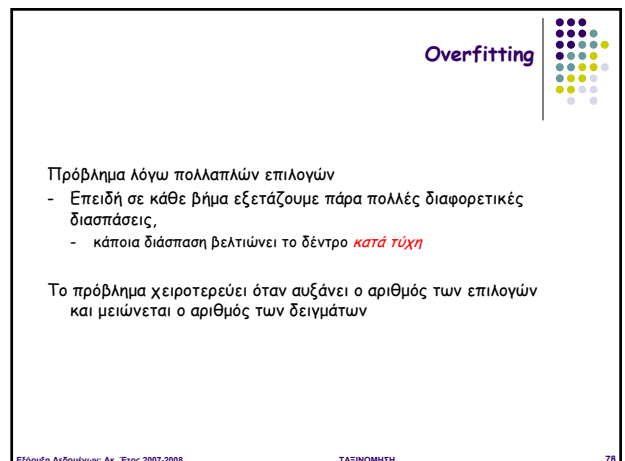
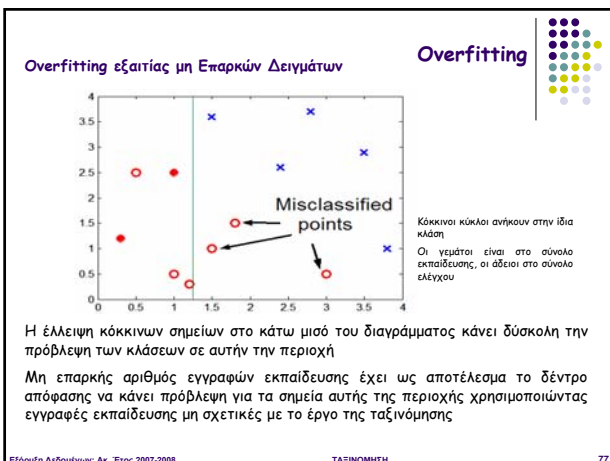
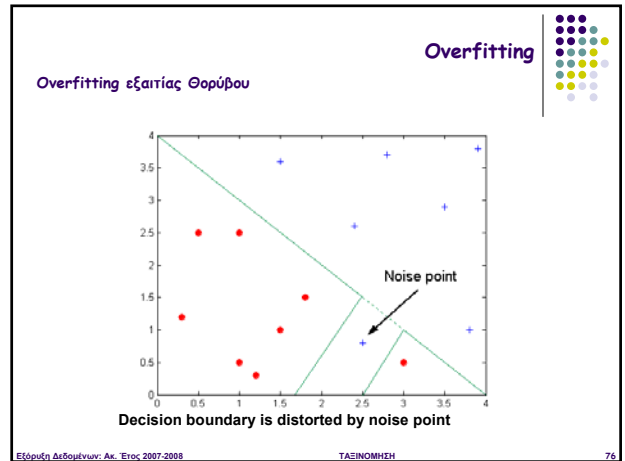
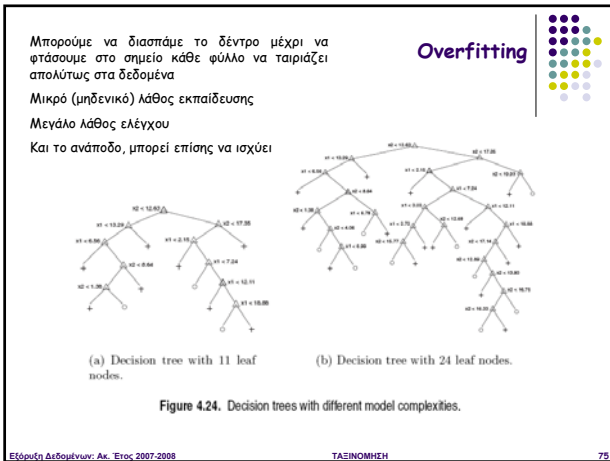
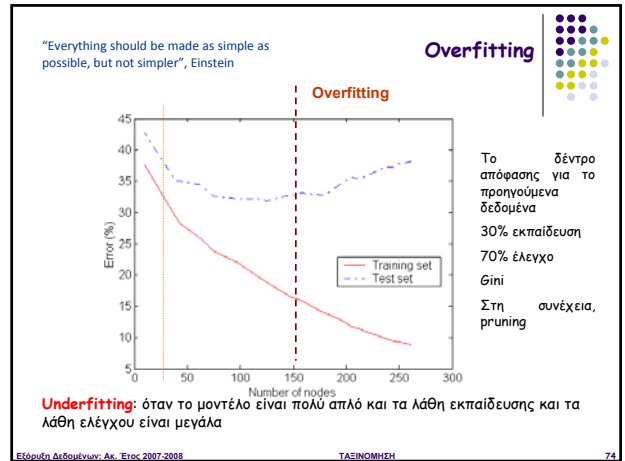
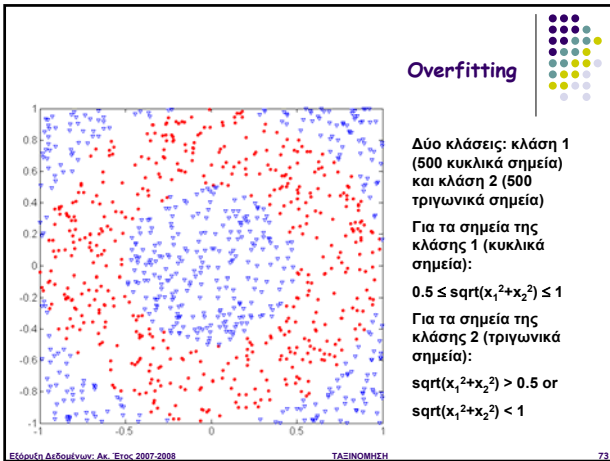
Λάθη

- **Εκπαίδευσης** (training, resubstitution, apparent): λάθη ταξινόμησης στα δεδομένα του συνόλου εκπαίδευσης (ποσοστό δεδομένων εκπαίδευσης που ταξινομούνται σε λάθος κλάση)
- **Γενίκευσης** (generalization): τα αναμενόμενα λάθη ταξινόμησης του μοντέλου σε δεδομένα που δεν έχει δει

**Overfitting**

Μπορεί ένα μοντέλο που ταιριάζει πολύ καλά με τα δεδομένα εκπαίδευσης να έχει μεγαλύτερο λάθος γενίκευσης από ένα μοντέλο που ταιριάζει λιγότερο καλά στα δεδομένα εκπαίδευσης

Εξόρυξη Δεδομένων: Ακ. Έτος 2007-2008      ΤΑΞΙΝΟΜΗΣΗ      72



## Overfitting

- Το overfitting έχει ως αποτέλεσμα δέντρα απόφασης που είναι πιο περίπλοκα από ό,τι χρειάζεται
- Τα λάθη εκπαίδευσης δεν αποτελούν πια μια καλή εκτίμηση για τη συμπεριφορά του δέντρου σε εγγραφές που δεν έχει δει ξανά
- Νέοι μέθοδοι για την εκτίμηση του λάθους

## Πολυπλοκότητα Μοντέλου

### Occam's Razor

- Δοθέντων δυο μοντέλων με παρόμοια λάθη γενίκευσης, πρέπει να προτιμάται το απλούστερο από το πιο περίπλοκο
- Ένα πολύπλοκο μοντέλο είναι πιο πιθανό να έχει ταιριαστεί (Fitted) τυχαία λόγω λαθών στα δεδομένα
- Για αυτό η πολυπλοκότητα του μοντέλου θα πρέπει να αποτελεί έναν από τους παράγοντες της αξιολόγησής του

## Εκτίμηση του Λάθους Γενίκευσης

- Re-substitution errors:** Λάθος στην εκπαίδευση ( $\sum e(t)$ )
- Generalization errors:** Λάθος στον έλεγχο ( $\sum e'(t)$ )

Ως λάθος μετράμε το ποσοστό των εγγραφών που ο ταξινομητής τοποθετεί σε λάθος κλάση

Μέθοδοι εκτίμησης του λάθους γενίκευσης:

### 1. Optimistic approach - Αισιόδοξη προσέγγιση:

$$e'(t) = e(t)$$

## Εκτίμηση του Λάθους Γενίκευσης

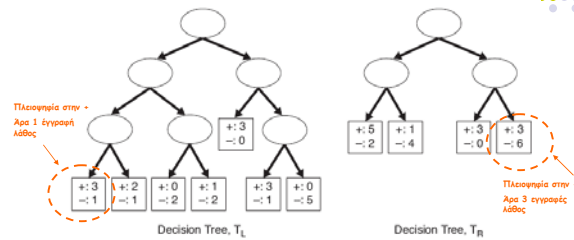


Figure 4.27. Example of two decision trees generated from the same training data.

Με βάση το λάθος εκπαίδευσης

Αριστερό  $4/24 = 0.167$

Δεξί:  $6/24 = 0.25$

## Εκτίμηση του Λάθους Γενίκευσης

### 2. Pessimistic approach - Απαισιόδοξη προσέγγιση:

$k$ : αριθμός φύλλων,  
για κάθε φύλλο  $t_i$  προσθέτουμε ένα κόστος  $V(t_i)$

$$e'(T) = \frac{\sum_{i=1}^k [e(t_i) + V(t_i)]}{\sum_{i=1}^k n(t_i)}$$

Αν για κάθε φύλλο  $t$ :  $e'(t) = e(t) + 0.5$   
Συνολικό λάθος:  $e'(T) = e(T) + k \times 0.5$  ( $k$ : αριθμός φύλλων)

Για ένα δέντρο με 30 φύλλα και 10 λάθη στο σύνολο εκπαίδευσης (από σύνολο 1000 εγγραφών):  
Training error =  $10/1000 = 1\%$   
Generalization error =  $(10 + 30 \times 0.5)/1000 = 2.5\%$

Το 0.5 σημαίνει ότι διαχωρισμός ενός κόμβου δικαιολογείται αν βελτιώνει τουλάχιστον μία εγγραφή

## Εκτίμηση του Λάθους Γενίκευσης

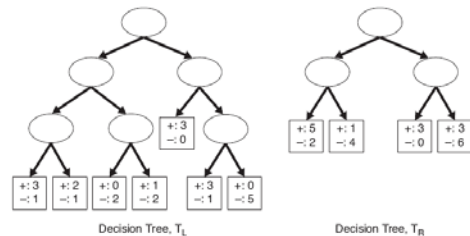


Figure 4.27. Example of two decision trees generated from the same training data.

Με βάση το λάθος εκπαίδευσης

Αριστερό  $(4 + 7 \times 0.5)/24 = 0.3125$

Δεξί:  $(6 + 4 \times 0.5)/24 = 0.3333$

Αντί για 0.5, κάτι μεγαλύτερο;

## Εκτίμηση του Λάθους Γενίκευσης

### 3. Reduced error pruning (REP):

- χρήση ενός συνόλου επαλήθευσης για την εκτίμηση του λάθους γενίκευσης

Χώρισε τα δεδομένα εκπαίδευσης:

2/3 εκπαίδευση

1/3 (σύνολο επαλήθευσης - validation set) για υπολογισμό λάθους

Χρήση για εύρεση του κατάλληλου μοντέλου

## Αντιμέτωπιση Overfitting

### Pre-Pruning (Early Stopping Rule)

Σταμάτα τον αλγόριθμο πριν σχηματιστεί ένα πλήρες δέντρο

Συνήθεις συνθήκες τερματισμού για έναν κόμβο:

- Σταμάτα όταν όλες οι εγγραφές ανήκουν στην ίδια κλάση
- Σταμάτα όταν όλες οι τιμές των γνωρισμάτων είναι οι ίδιες

Πιο περιοριστικές συνθήκες:

- Σταμάτα όταν ο αριθμός των εγγραφών είναι μικρότερος από κάποιο προκαθορισμένο κατώφλι

- Σταμάτα όταν η επέκταση ενός κόμβου δεν βελτιώνει την καθαρότητα (π.χ., *gain* ή *information gain*) ή το λάθος γενίκευσης περισσότερο από κάποιο κατώφλι. (-) δύσκολος ο καθορισμός του κατωφλιού, (-) αν και το κέρδος μικρό, κατοπινοί διαχωρισμοί μπορεί να καταλήξουν σε καλύτερα δέντρα

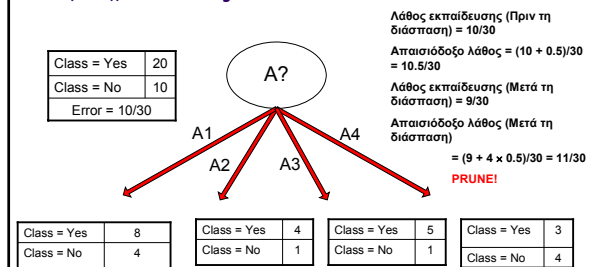
## Overfitting

### Post-pruning

- Ανάπτυξε το δέντρο πλήρως
- Trim - ψαλίδιασε τους κόμβους bottom-up
- Αν το λάθος γενίκευσης μειώνεται με το ψαλίδισμα, αντικατέστησε το υποδέντρο με
  - ένα φύλλο - οι ετικέτες κλάσεων του φύλλου καθορίζεται από την πλειοψηφία των κλάσεων των εγγραφών του υποδέντρου (subtree replacement)
  - ένα από τα κλαδιά του (Branch), αυτό που χρησιμοποιείται συχνότερα (subtree raising)

## Overfitting

### Παράδειγμα Post-Pruning



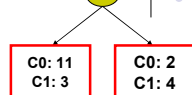
### Παράδειγμα post-pruning

## Overfitting

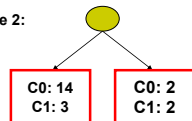
- Αισιόδοξη προσέγγιση?  
Όχι διάσπαση
- Απαισιόδοξη προσέγγιση?  
όχι case 1, ναι case 2
- REP?

Εξαρτάται από το σύνολο επαλήθευσης

Case 1:



Case 2:



## Τιμές που λείπουν

- Οι τιμές που λείπουν επηρεάζουν την κατασκευή του δέντρου με τρεις τρόπους:
  - Πώς υπολογίζονται τα μέτρα καθαρότητας
  - Πώς κατανέμονται στα φύλλα οι εγγραφές με τιμές που λείπουν
  - Πώς ταξινομείται μια εγγραφή εκπαίδευσης στην οποία λείπει μια τιμή

### Τιμές που λείπουν

**Υπολογισμοί μέτρων καθαρότητας**

Tid	Refund	Marital Status	Taxable Income	Class
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	?	Single	90K	Yes

Πριν τη διάσπαση:  
 $Entropy(\text{Parent}) = -0.3 \log(0.3) - (0.7) \log(0.7) = 0.8813$

	Class = Yes	Class = No
Refund=Yes	0	3
Refund=No	2	4
Refund=?	1	0

Διάσπαση στο Refund:  
 $Entropy(\text{Refund}=\text{Yes}) = 0$   
 $Entropy(\text{Refund}=\text{No}) = -(2/6) \log(2/6) - (4/6) \log(4/6) = 0.9183$   
 $Entropy(\text{Children}) = 0.3(0) + 0.6(0.9183) = 0.551$   
 $Gain = 0.9 \times (0.8813 - 0.551) = 0.3303$

Missing value

Εξέταση Δεξοτέρων Ακ. Έτος 2007-2008 ΤΑΞΙΝΟΜΗΣΗ 91

### Τιμές που λείπουν

Tid	Refund	Marital Status	Taxable Income	Class
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	?	Single	90K	Yes

Σε ποιο φύλλο;  
**Πιθανότητα Refund=Yes is 3/9 (3 από τις 9 εγγραφές έχουν refund=Yes)**  
**Πιθανότητα Refund=No is 6/9**  
**Ανάθεση εγγραφής στο αριστερό παιδί με βάρος 3/9 και στο δεξί παιδί με βάρος 6/9**

Yes No

**Refund**

Class=Yes	0	Class=Yes	2
Class=No	3	Class=No	4

Yes No

**Refund**

Class=Yes	0 + 3/9	Class=Yes	2 + 6/9
Class=No	3	Class=No	4

Εξέταση Δεξοτέρων Ακ. Έτος 2007-2008 ΤΑΞΙΝΟΜΗΣΗ 92

### Τιμές που λείπουν

**Νέα εγγραφή**

Tid	Refund	Marital Status	Taxable Income	Class
11	No	?	85K	?

	Married	Single	Divorced	Total
Class=No	3	1	0	4
Class=Yes	6/9	1	1	2.67
Total	3.67	2	1	6.67

Refund

Yes: NO

No: MarSt

Single, Divorced: TaxInc

< 80K: NO

> 80K: YES

Married: NO

Πιθανότητα οικογενειακή κατάσταση (MarSt) = Married is 3.67/6.67

Πιθανότητα οικογενειακή κατάσταση (MarSt) = {Single, Divorced} is 3/6.67

Εξέταση Δεξοτέρων Ακ. Έτος 2007-2008 ΤΑΞΙΝΟΜΗΣΗ 93