# Ensembles



# Coneixement, Raonament i Incertesa.
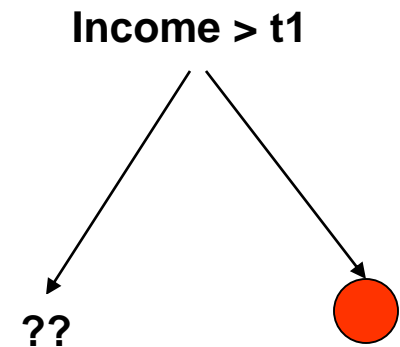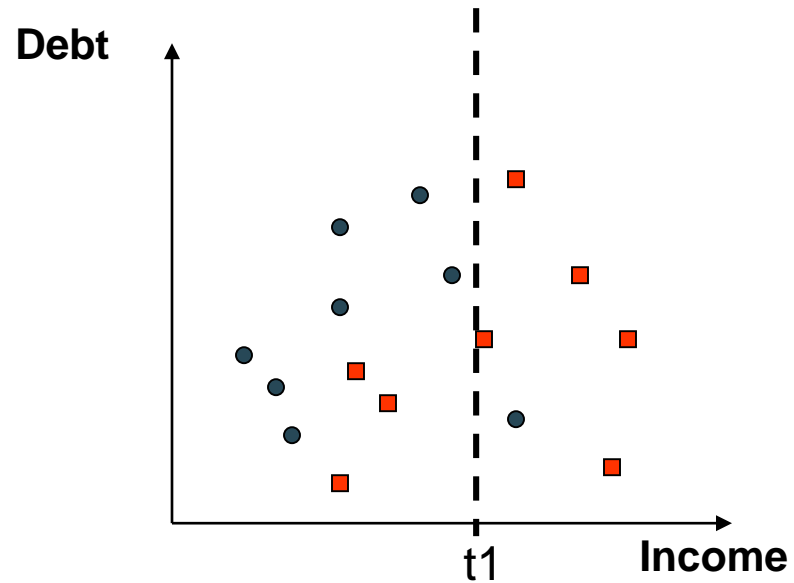
# Decision Tree Example



**Debt**

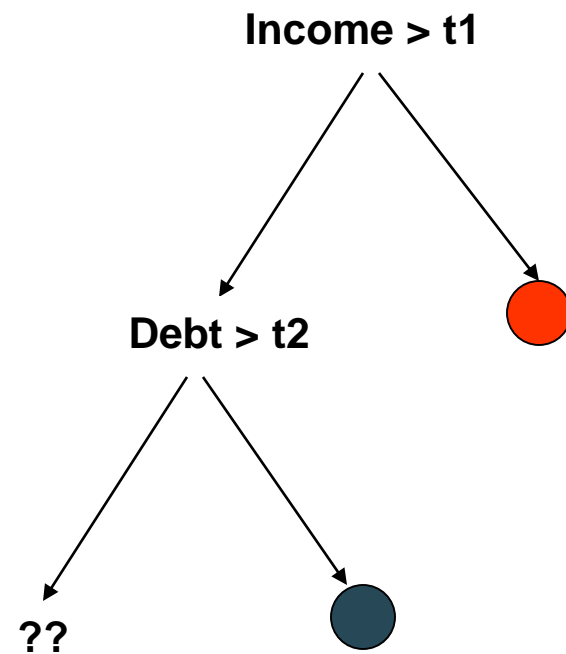**Income**

# Decision Tree Example

**Debt**

t1

**Income**

**Income > t1**

**??**

# Decision Tree Example

# Decision Tree Example



**Debt**

t2

t3    t1    **Income**

**Income > t1**

**Debt > t2**

**Income > t3**

# Decision Tree Example

**Debt**

t2

t3    t1    **Income**

Note: tree boundaries are linear and axis-parallel

**Income > t1**

**Debt > t2**

**Income > t3**

# Minimum Distance Classifier

# Local Decision Boundaries

Boundary? Points that are equidistant
between points of class 1 and 2
Note: locally the boundary is linear

1 —— 2

Feature 2

1

2

?          2

1

Feature 1

# Finding the Decision Boundaries

Feature 2

1

2

1

1

2

?

2

1

Feature 1

# Finding the Decision Boundaries

# Finding the Decision Boundaries

# Overall Boundary = Piecewise Linear

Decision Region
for Class 1

Decision Region
for Class 2

Feature 2

1

2

1

2

2

1

2

Feature 1

- Bagging and Boosting

    ➔ Aggregating Classifiers



FINAL RULE

Breiman (1996) found gains in accuracy by aggregating predictors built from reweighed versions of the learning set

# Bagging and Boosting:
## Aggregating Classifiers

*3 questions:*

? How to **reweigh** ?

? How to **aggregate** ?

? Which type of **gain**

in accuracy ?

# Bagging

- *Bagging* = Bootstrap Aggregating

- Reweighing of the learning sets is done by drawing at random with replacement from the learning sets

- Predictors are aggregated by plurality voting

# The Bagging Algorithm

- B bootstrap samples

- From which we derive:

  - **B Classifiers** $\in \{-1, 1\}: c^1, c^2, c^3, ..., c^B$

  - **B Estimated probabilities** $\in [0, 1]: p^1, p^2, p^3, ..., p^B$

The aggregate classifier becomes:

$$c_{bag}(x) = sign\left(\frac{1}{B}\sum_{b=1}^{B} c^b(x)\right)$$ or $$p_{bag}(x) = \frac{1}{B}\sum_{b=1}^{B} p^b(x)$$

# Bagging Example (Opitz, 1999)

| Original | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| Training set 1 | 2 | 7 | 8 | 3 | 7 | 6 | 3 | 1 |
| Training set 2 | 7 | 8 | 5 | 6 | 4 | 2 | 7 | 1 |
| Training set 3 | 3 | 6 | 2 | 7 | 5 | 6 | 2 | 2 |
| Training set 4 | 4 | 5 | 1 | 4 | 6 | 4 | 3 | 8 |

# Aggregation

## Sign

**Initial set**
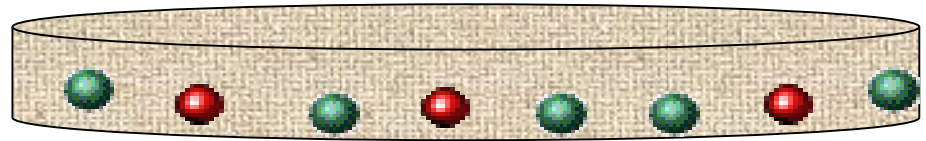
Classifier 1

+

Classifier 2

+

Classifier 3

+

...

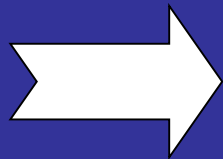+

Classifier T

## Final rule

# Boosting

- Freund and Schapire (1997), Breiman (1998)

- Data *adaptively resampled*

Previously misclassified observations ➜ weights ⬆

Previously wellclassified observations ➜ weights ⬇

➡ **Predictor aggregation done by *weighted voting***

# AdaBoost

$$y_i \in \{-1, +1\}$$

- Initialize weights: $w_i^1 = 1/N$

- Fit a classifier with these weights
- Give predicted probabilities to observations according to this classifier

$$p_b(x) = \hat{P}_w(y = 1 | x) \in [0,1]$$

- Compute "pseudo probabilities": $f_b(x) = \frac{1}{2} \log\left( \frac{p_b(x)}{1 - p_b(x)} \right) \in \Re$

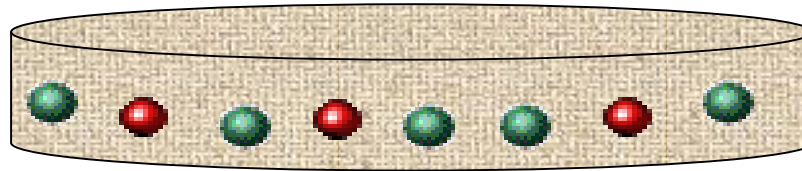- Get new weights: $w_i^{b+1} = w_i^b \exp\left[ - y_i f_b(x_i) \right]$

& "Normalize" it (i.e., rescale so that it sums to 1)

- Combine the "pseudo probabilities": $c_{Boost} = sign\left[ \sum_{b=1}^{B} f_b(x) \right]$

# Weighting

Initial set



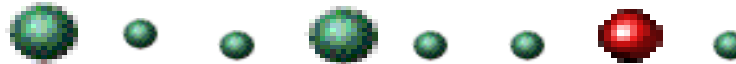**Classifier 1**

$$f_1(x)$$

**Checking & Modification**

**Classifier 2**

$$f_2(x)$$

**Checking & Modification**

$+$

$+$

$...$

# Aggregation

**Initial set**

**Sign**

$$f_1(x)$$

$$+$$

$$f_2(x)$$

$$+$$

$$f_3(x)$$

$$+$$

$$...$$

$$+$$

$$f_B(x)$$

$$=$$

**Classifier 1**

**Classifier 2**

**Classifier 3**

**Classifier ………**

**Classifier B**

**Final rule**

# Boosting

- Definition of Boosting:

  Boosting refers to a general method of producing a very accurate prediction rule by combining rough and moderately inaccurate rules-of-thumb.

- Intuition:

  1) No learner is always the best;

  2) Construct a set of base-learners which when combined achieves higher accuracy

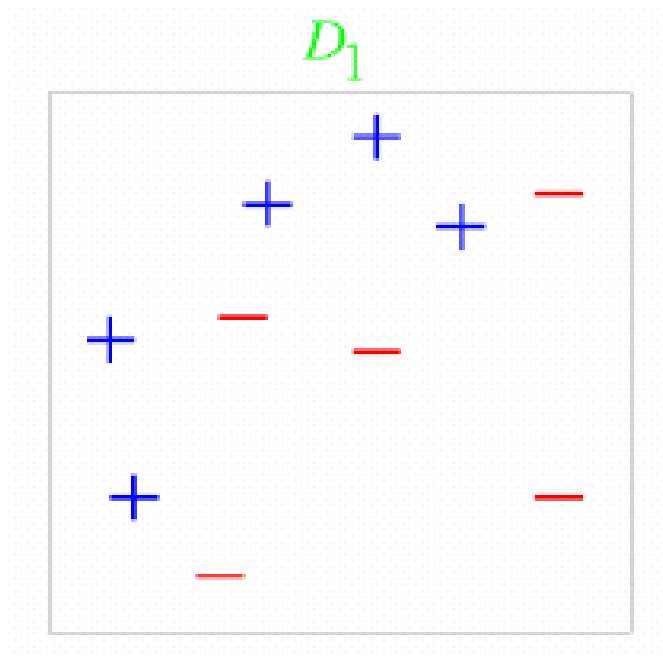# Boosting

3) Different learners may:

  --- Be trained by different algorithms
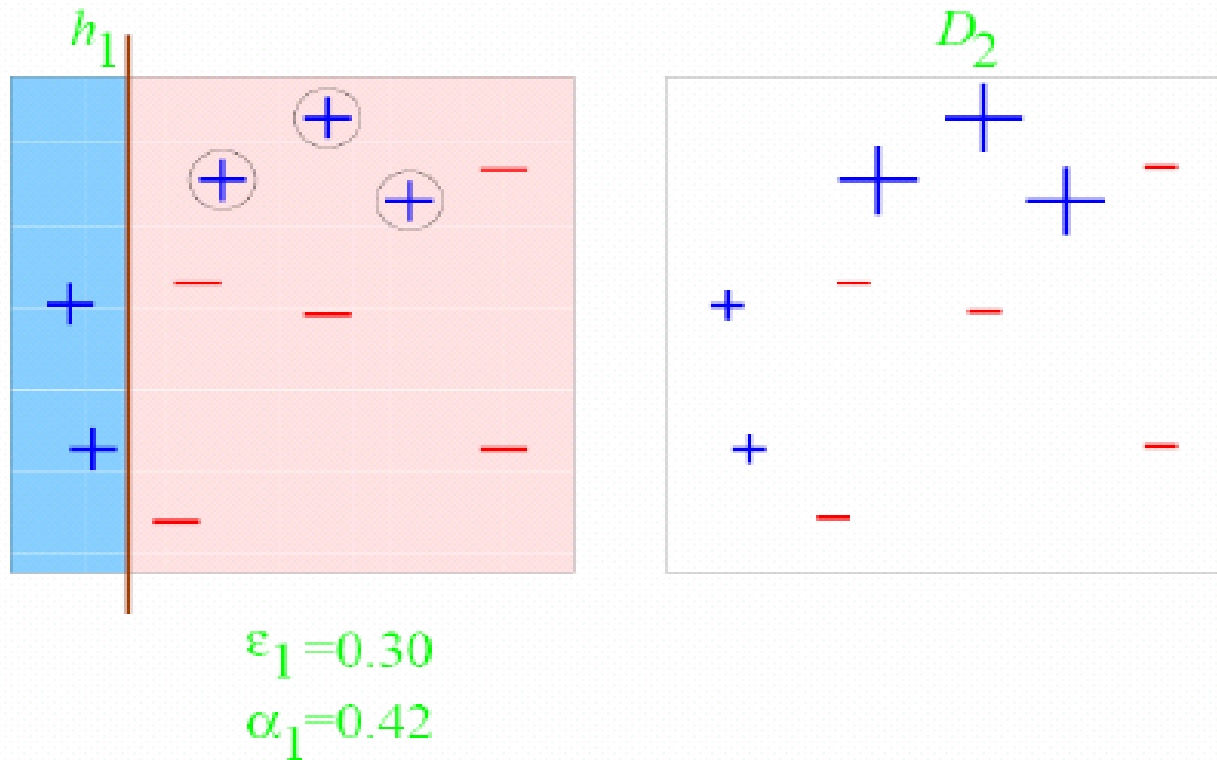
  --- Use different modalities(features)

  --- Focus on different subproblems

  --- ……

4) A week learner is "rough and moderately inaccurate" predictor but one that can predict better than chance.

# A toy example



$D_1$

# A toy example(cont'd)



$\varepsilon_1 = 0.30$

$\alpha_1 = 0.42$

# A toy example(cont'd)



$h_2$

$D_3$

$\varepsilon_2 = 0.21$

$\alpha_2 = 0.65$

# A toy example(cont'd)



$\varepsilon_3 = 0.14$

$\alpha_3 = 0.92$

# A toy example(cont'd)



$$H_{\text{final}} = \text{sign}\left( 0.42 \quad + 0.65 \quad + 0.92 \right)$$

$$=$$

# Ensembles Methods

Funcionament:

- Aprendre multiples definicions alternatives d'un concepte usant **diferents dades d'aprenentatge** o **diferents algorismes d'aprenentatge**.
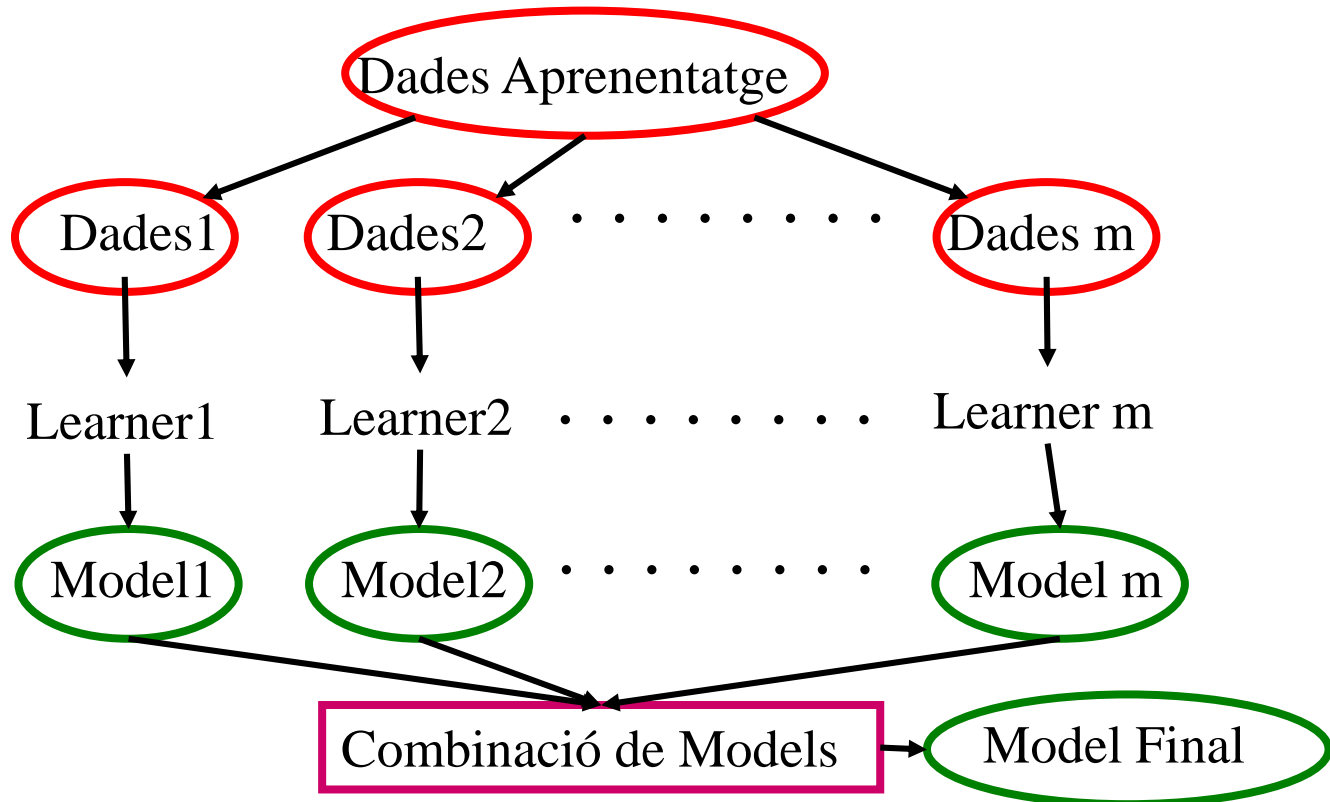
- **Combinar** les decisions de multiples definicions, p.ex. Usant el vot pesat.
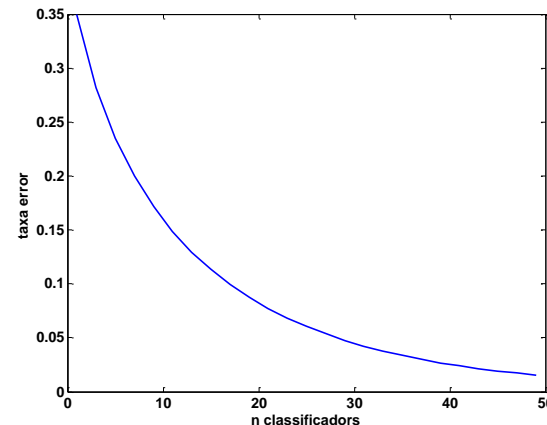
# Perque funcionen?

Suposem que tenim 25 classificadors base

- Cada classificador té un taxa d'error, $\varepsilon = 0.35$
- Suposem que els classificadors són independents
- La probabilitat que el 'ensemble classifier' faci una predicció erronia (si s'equivoca en 13 de les 25 prediccions) :

$$\sum_{i=13}^{25} \binom{25}{i} \varepsilon^i (1-\varepsilon)^{25-i} = 0.06$$

# Valor dels 'Ensembles'

- Quan combinem múltiples decisions **independents** i **diverses** cada un de les cuals és millor que l'atzar, els errors deguts a atzar es cancel·len els uns als altres, i les decisions correctes es reforcen.

# Ensembles Homogenis

Utilitzar un **únic, algorisme d'aprenentatge arbitrari** però **<span style="color:red">manipular les dades d'aprenentatge</span>** per a fer-lo aprendre multiples models.

- Data1 $\neq$ Data2 $\neq$ ... $\neq$ Data m
- Learner1 = Learner2 = ... = Learner m
- Model 1 $\neq$ Model 2 $\neq$ ... $\neq$ Model m

Mètodes per canviar les dades d'aprenentatge:

- Bagging: Re-mostreigar les dades d'aprenentatge
- Boosting: Re-pesar les dades d'aprenentatge
- Decorate: Afegir dades d'aprenentatge adicionals artificials