

# Some enhancements in Decision Trees

**Djamel A. Zighed**

ERIC Laboratory  
University of Lyon 2 (France)  
zighed@univ-lyon2.fr

# Outline

1. Decision tree / some recalls
2. Insensitivity of the criterion to the sample size
3. Entropy measure sensitive to the sample size
4. Lattice for mining small sample
5. Example on brandy control
6. Conclusion

# 1. Some recalls / Decision tree

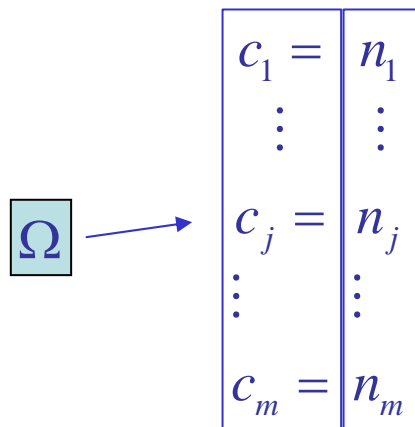
## Classification Task

Let's consider a learning sample  $\Omega$ ;

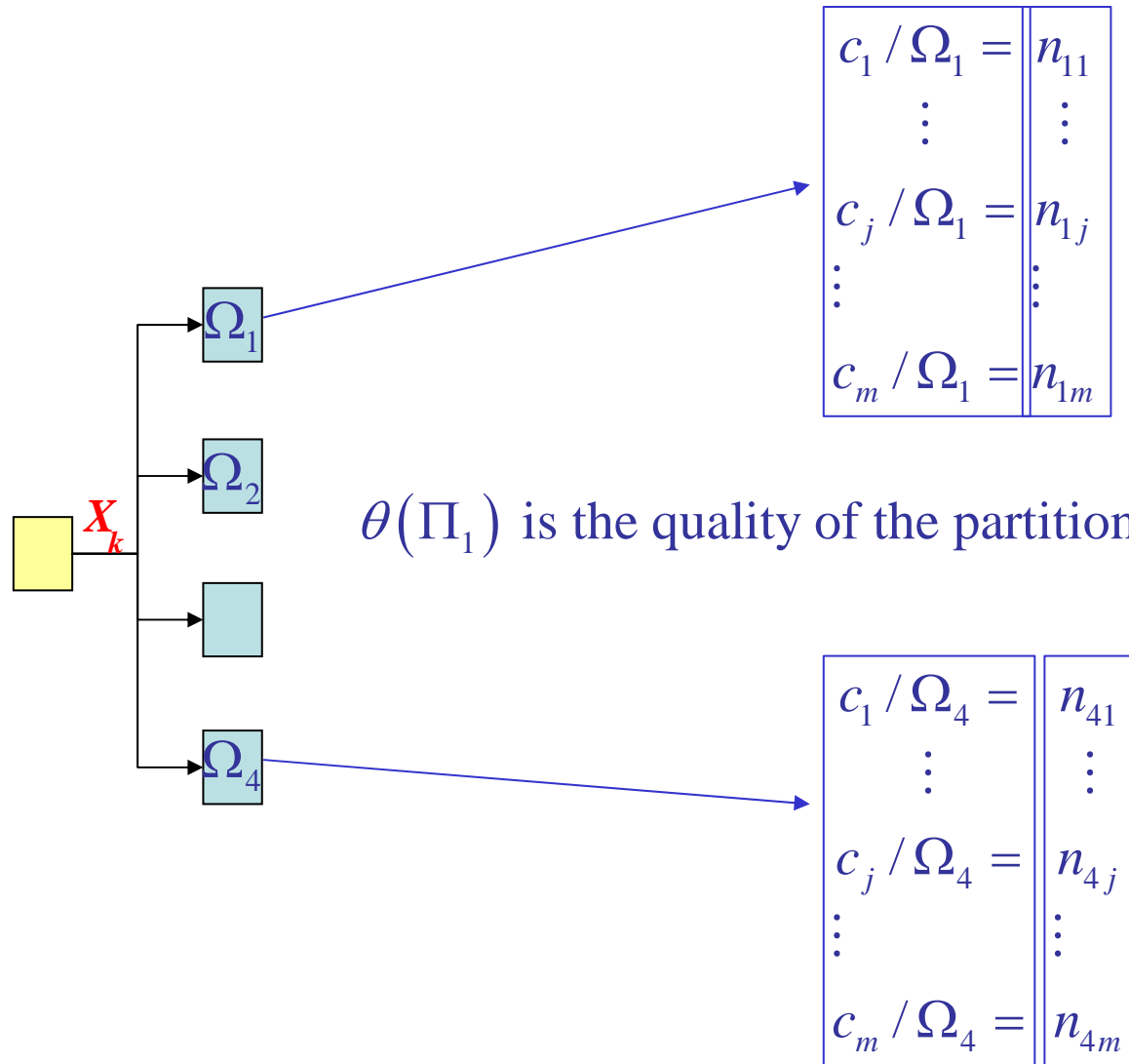
The size of  $\Omega$  :  $|\Omega| = n$

$\forall \omega \in \Omega; X(\omega) = (X_1(\omega), \dots, X_p(\omega))$  : Predictive attributes

$\forall \omega \in \Omega; C(\omega)$  : Predicted attribute;  $C(\omega) \in \{c_1, c_2, \dots, c_m\}$

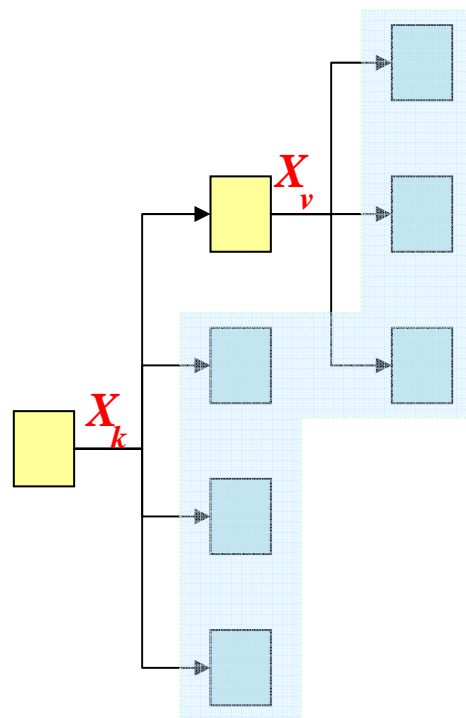


$X_k$  : Brings about a partition  $\Pi_1$  on  $\Omega$



$\theta(\Pi_1)$  is the quality of the partition  $\Pi_1$  brought about by  $X_k$

$X_k; X_v$  : Bring about a partition  $\Pi_2$  on  $\Omega$

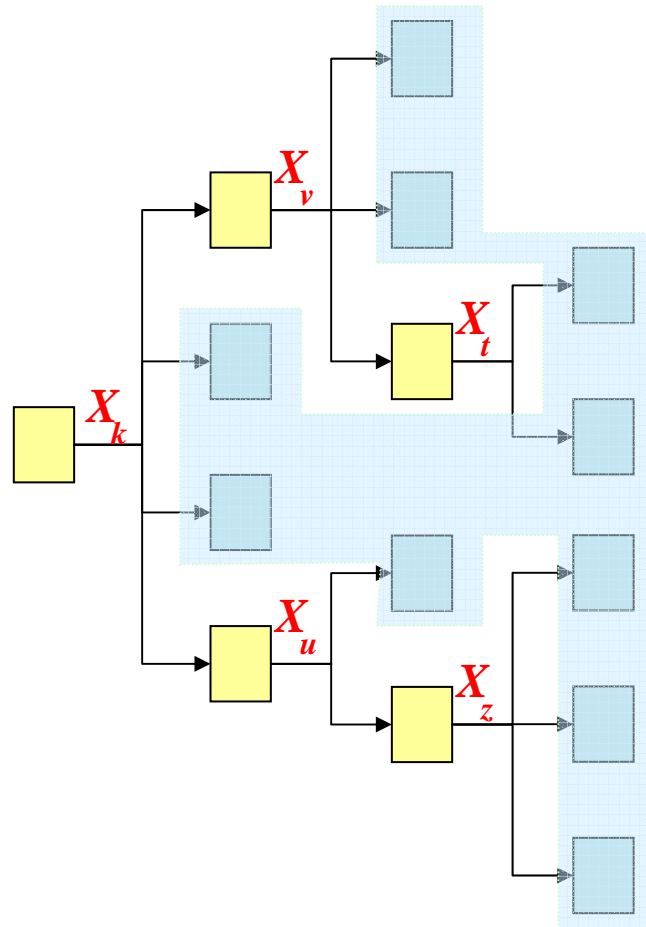


$\theta(\Pi_2)$

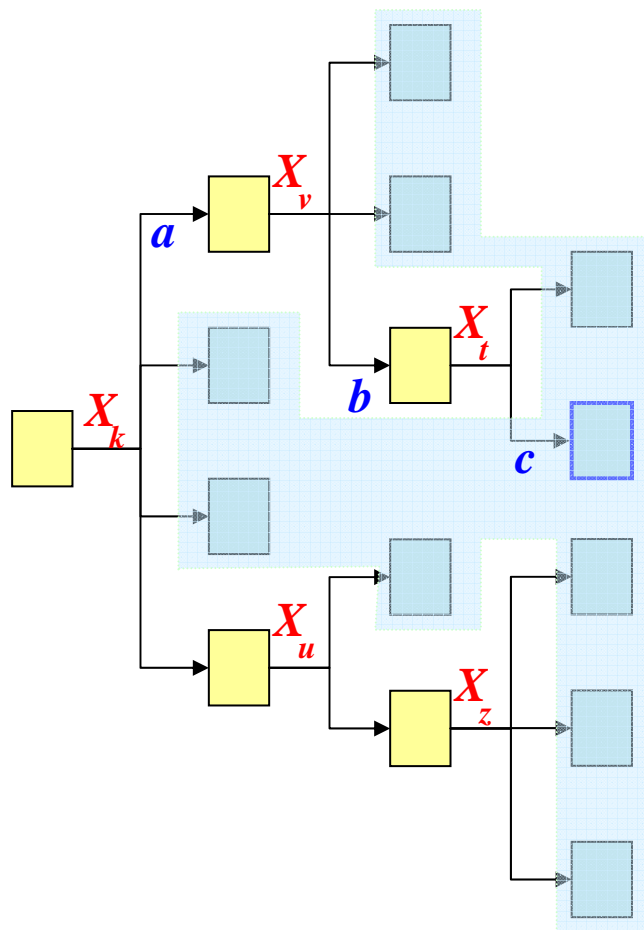
$X_k; X_v, X_t, X_u, X_z$  : Bring about a partition  $\Pi_s$  on  $\Omega$

The size :  $|\Pi_s| = \pi = 10$

$\theta(\Pi_s)$



## Criterion for Growing the tree



$$\theta(\Pi_s) = \sum_{s \in S} w_s h(s)$$

were

$w_s$  is the weight of  $s$

and

$h(s)$  is the purity of  $s$

$$h(s) = - \sum_{i=1}^m p(c_i / s) \log p(c_i / s)$$

$$h(s) = \sum_{i=1}^m p(c_i / s) (1 - p(c_i / s))$$

$w_s$  = The probability of the path  $s$  :

$$w_s = P(X_k = x; X_v = y, \dots) = \frac{|\Pi_s|}{|\Pi|}$$

$$p(c_i / s) = \frac{n_{is}}{n_{.s}}$$

## Entropy measure :

$$h : \Pi \rightarrow \mathbb{R}^+$$

$$\Pi = \bigcup_n S_n$$

$$S_n = \left\{ \gamma \in \mathbb{R}^n / \gamma_i \geq 0 (i = 1, \dots, n) \text{ and } \sum_{i=1}^n \gamma_i = 1 \right\}$$

⇒ Symetry

$$h(\gamma_1, \gamma_2, \dots, \gamma_n) = h(\gamma_{\sigma_1}, \gamma_{\sigma_2}, \dots, \gamma_{\sigma_n})$$

⇒ Minimality

$$h(\gamma_1, \gamma_2, \dots, \gamma_n) = 0 \Leftrightarrow \exists i / \gamma_i = 1 \quad \forall j \neq i ; \gamma_j = 0$$

⇒ Maximality

$$h(\gamma_1, \gamma_2, \dots, \gamma_n) = \text{Max} \Leftrightarrow \forall i (i = 1, \dots, n) ; \gamma_i = \frac{1}{n}$$

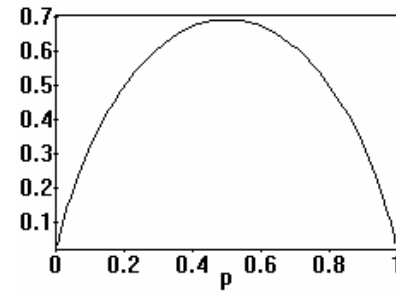
⇒ Continuity

⇒ Concavity



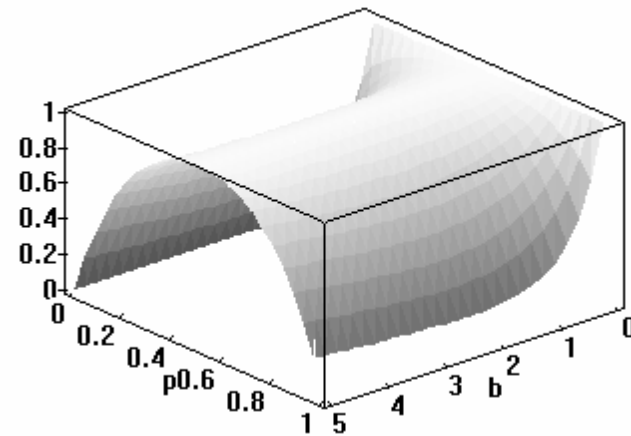
$$h(s) = h(p_1, p_2)$$

$$h(s) = -\sum_{i=1}^2 p_i \log p_i$$



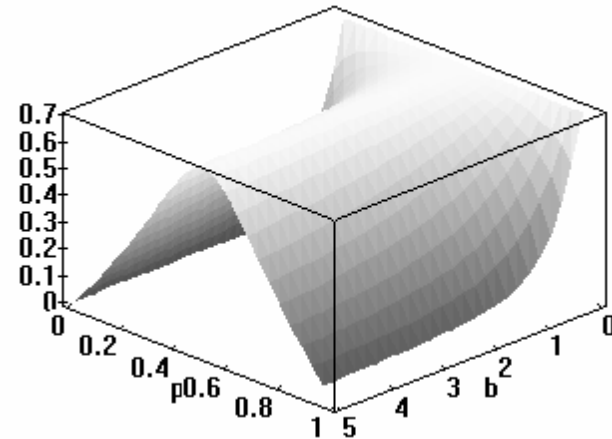
$$h(s) = \frac{1}{2^{\beta-1} - 1} \left[ \left( \sum_{i=1}^2 p_i^{\beta} - 1 \right) \right]$$

$$\beta > 0; \beta \neq 1$$

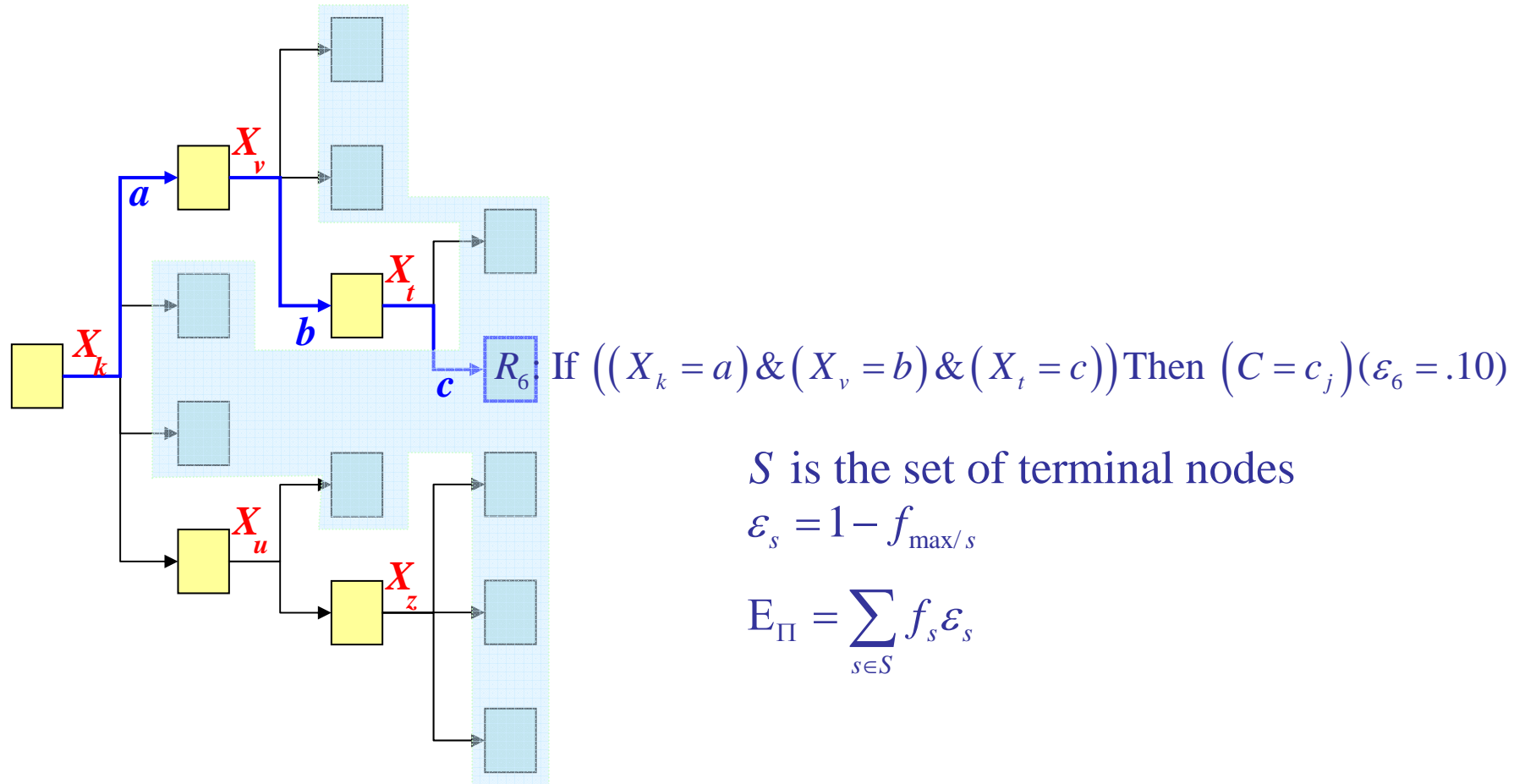


$$h(s) = \frac{1}{1-\beta} \log \left( \sum_{i=1}^2 p_i^{\beta} \right)$$

$$\beta > 0; \beta \neq 1$$



Each partition  $\Pi$  on  $\Omega$  is described by a model  $\varphi$  : a set of rules  
 $R_i$ : If (condition) Then (Conclusion) ( $\varepsilon$ )



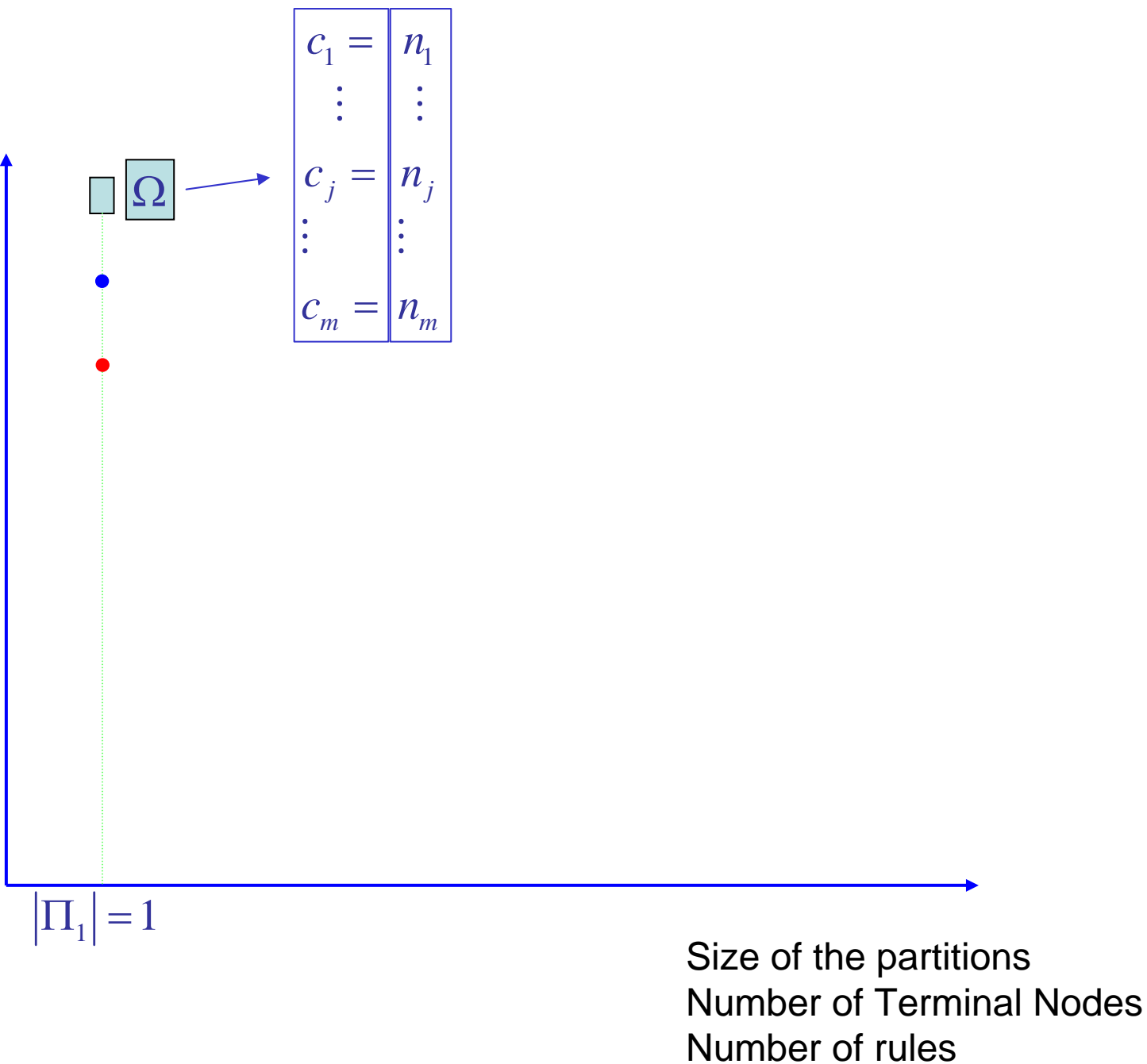
$S$  is the set of terminal nodes

$$\varepsilon_s = 1 - f_{\max/s}$$

$$E_{\Pi} = \sum_{s \in S} f_s \varepsilon_s$$

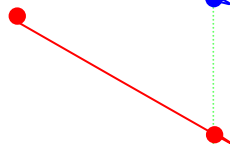
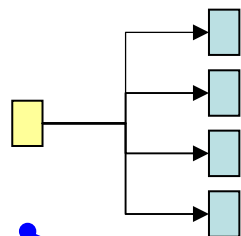
$$\theta(\Pi_s) = \sum_{s \in \mathcal{S}} w_s h(s)$$

$$E_{\Pi} = \sum_{s \in \mathcal{S}} f_s \varepsilon_s$$



$$\theta(\Pi_s) = \sum_{s \in \mathcal{S}} w_s h(s)$$

$$E_{\Pi} = \sum_{s \in \mathcal{S}} f_s \varepsilon_s$$

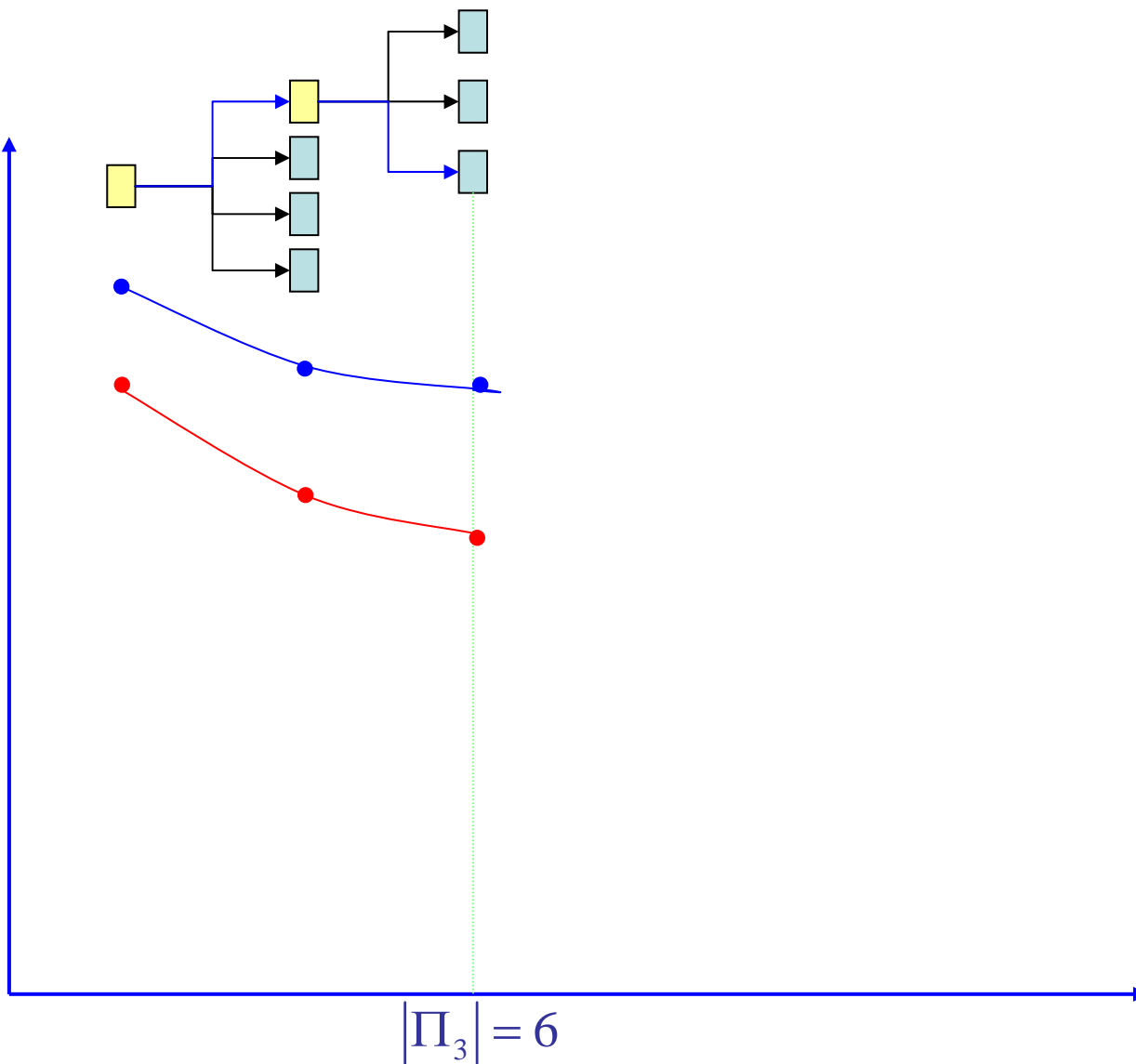


$$|\Pi_2| = 4$$

Size of the partitions  
 Number of Terminal Nodes  
 Number of rules

$$\theta(\Pi_s) = \sum_{s \in \mathcal{S}} w_s h(s)$$

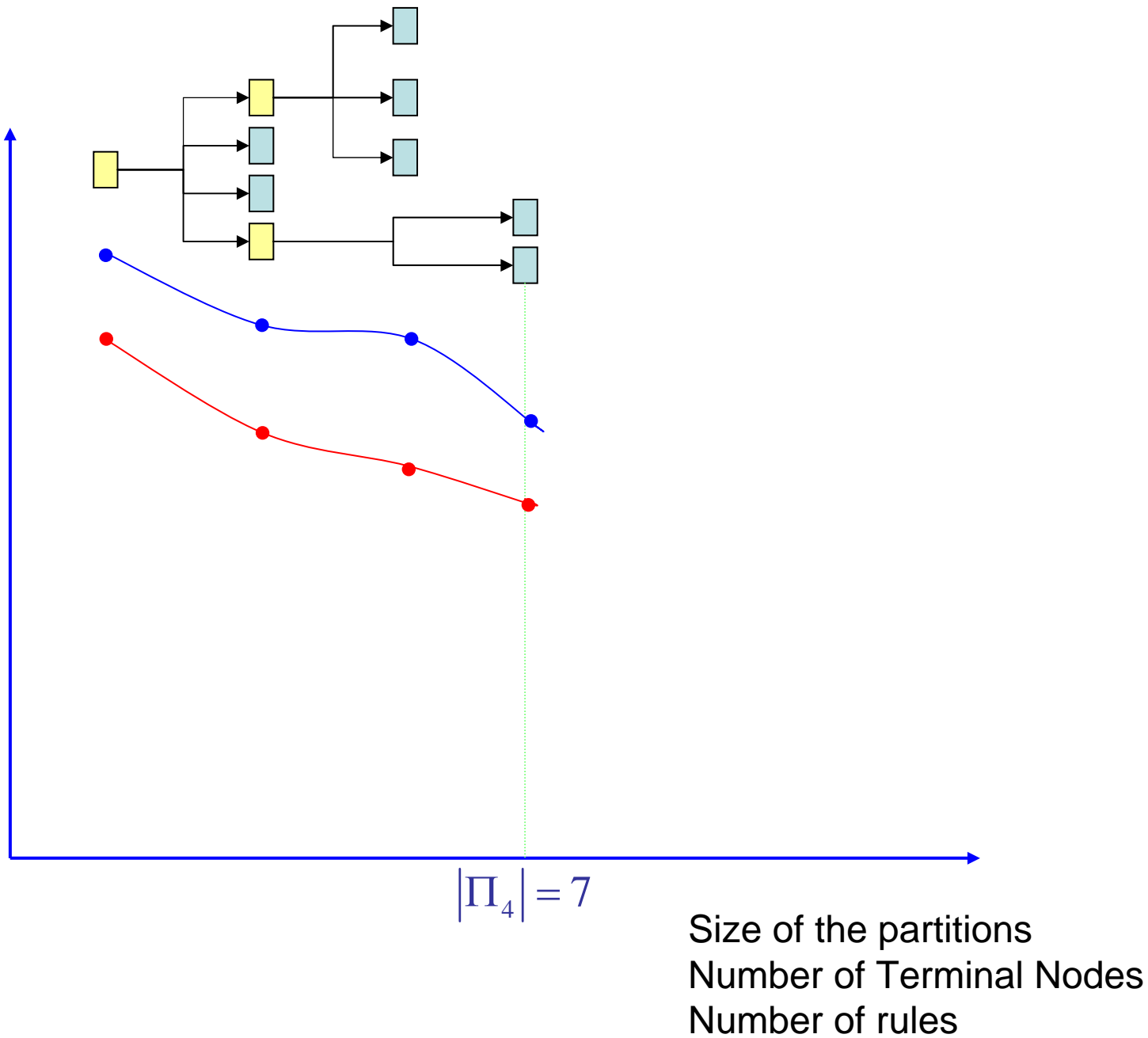
$$E_{\Pi} = \sum_{s \in \mathcal{S}} f_s \varepsilon_s$$



Size of the partitions  
 Number of Terminal Nodes  
 Number of rules

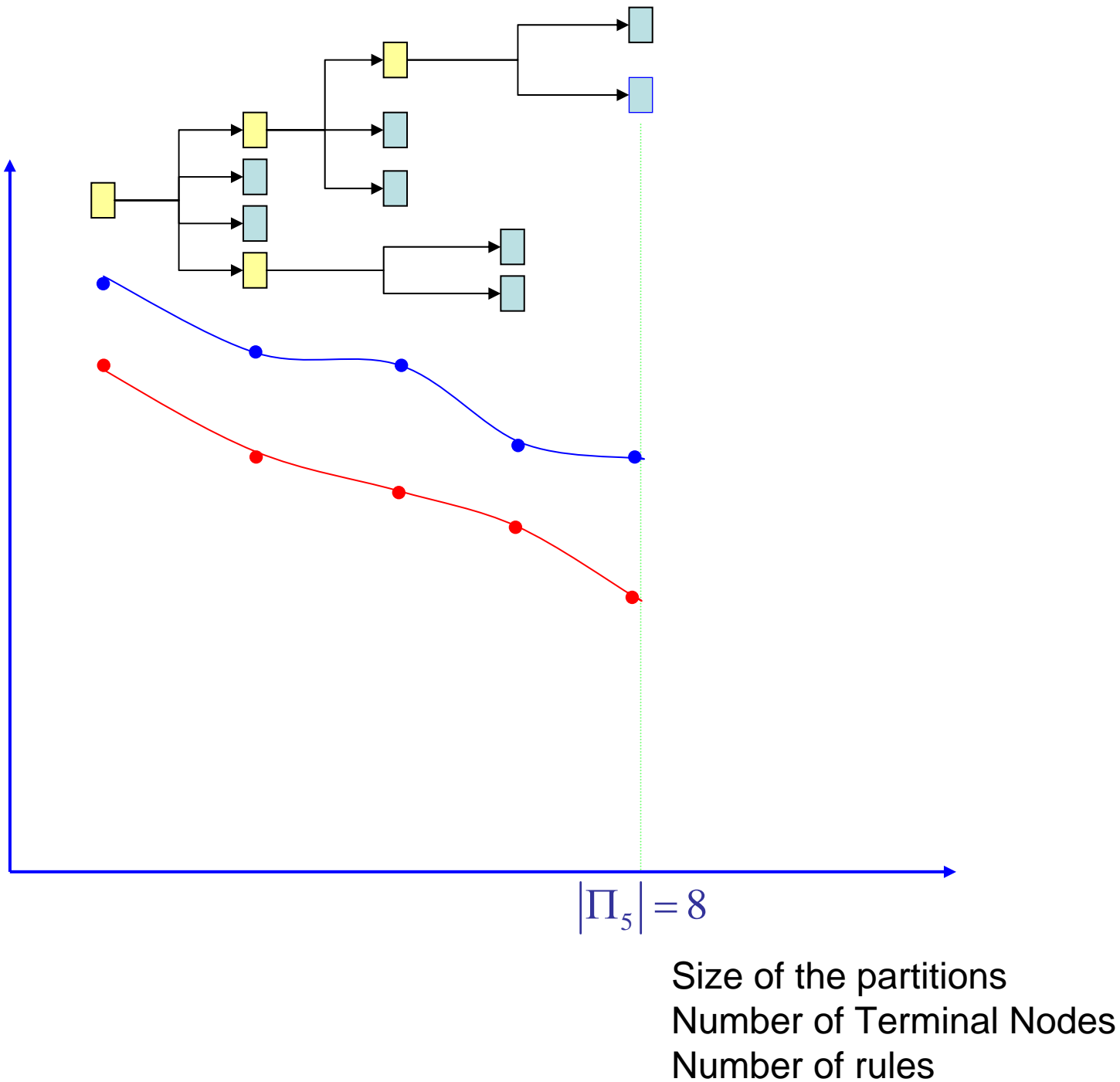
$$\theta(\Pi_s) = \sum_{s \in \mathcal{S}} w_s h(s)$$

$$E_{\Pi} = \sum_{s \in \mathcal{S}} f_s \varepsilon_s$$



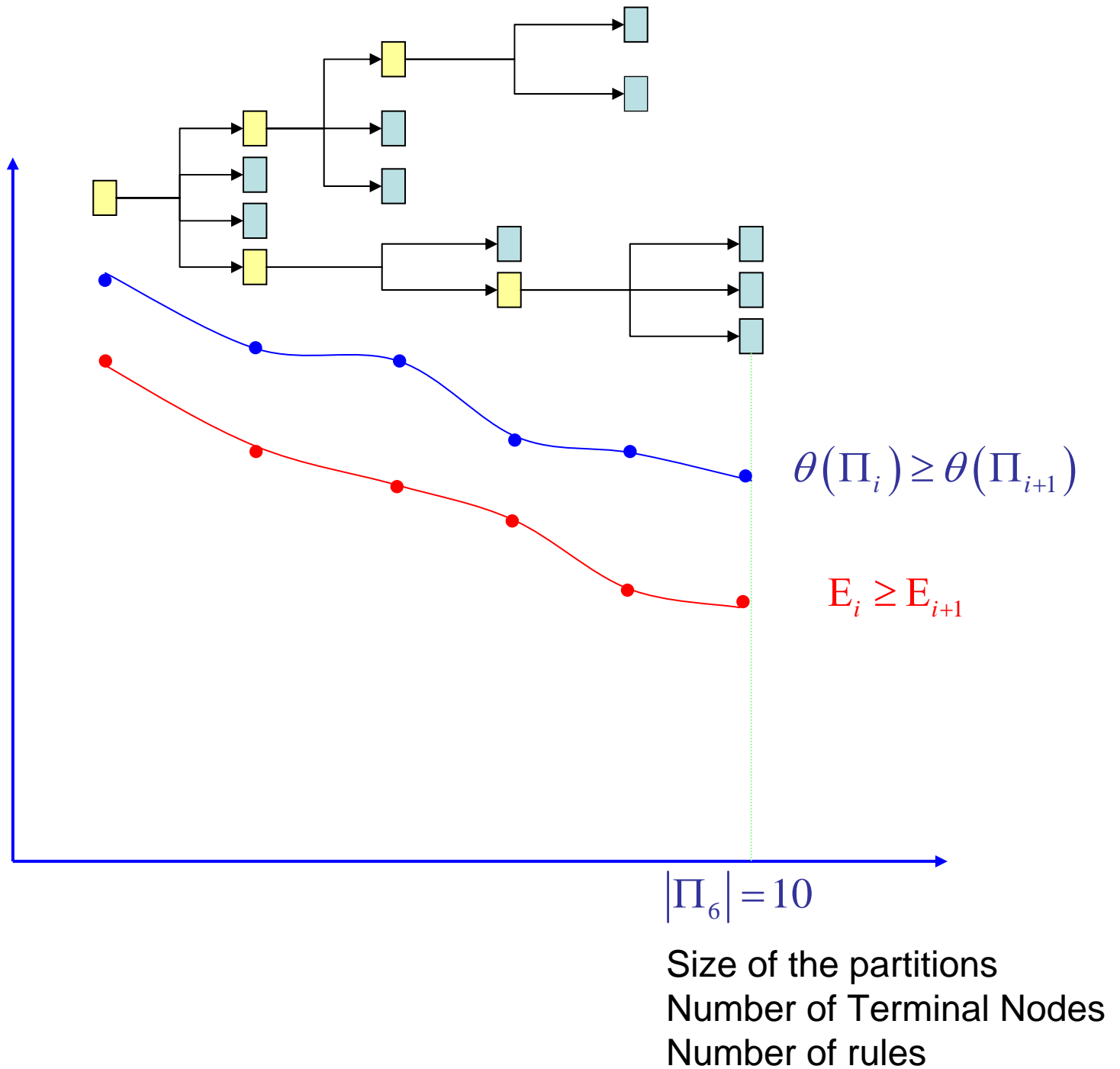
$$\theta(\Pi_s) = \sum_{s \in \mathcal{S}} w_s h(s)$$

$$E_{\Pi} = \sum_{s \in \mathcal{S}} f_s \varepsilon_s$$



$$\theta(\Pi) = \sum_{s \in \mathcal{S}} w_s h(s)$$

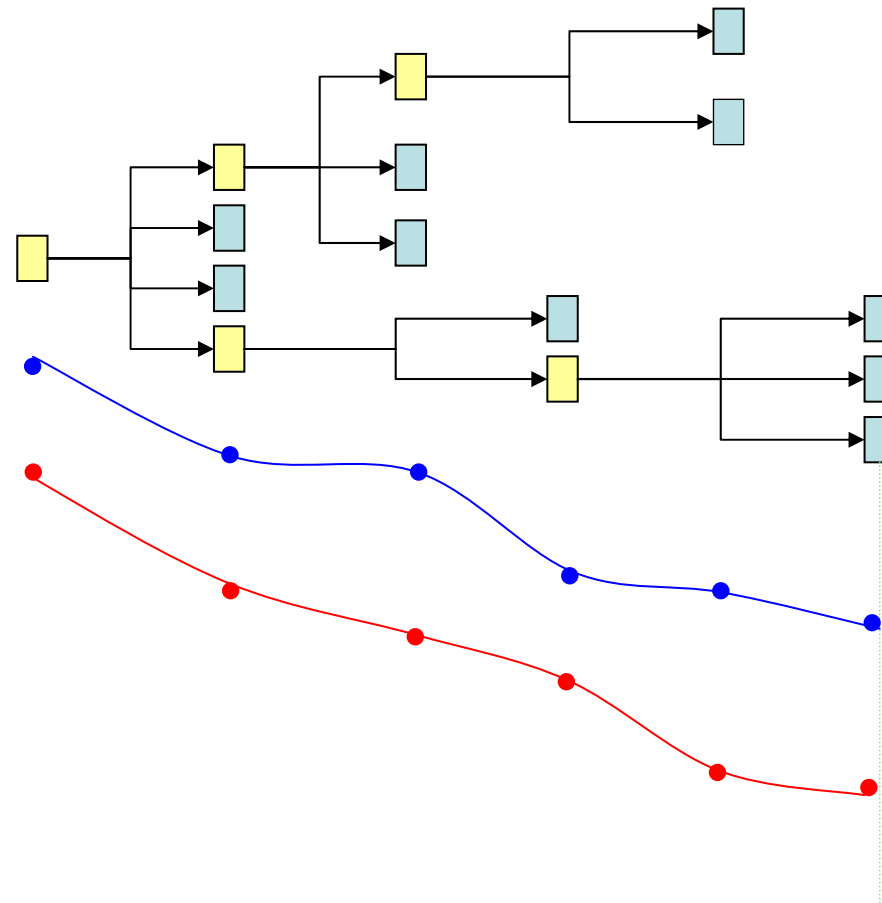
$$E_{\Pi} = \sum_{s \in \mathcal{S}} f_s \varepsilon_s$$





$$\theta(\Pi) = \sum_{s \in \mathcal{S}} w_s h(s)$$

$$E_{\Pi} = \sum_{s \in \mathcal{S}} f_s \varepsilon_s$$



$$\theta(\Pi_i) \geq \theta(\Pi_{i+1})$$

$$E_i \geq E_{i+1}$$

The largest tree will have the lowest value of  $\theta$  and error  $E$ ,  
 Shall we consider the predictive model associated to the  
 largest tree as the best one ?

$$\theta(\Pi) = \sum_{s \in S} w_s h(s)$$

$$E_{\Pi} = \sum_{s \in S} f_s \varepsilon_s$$

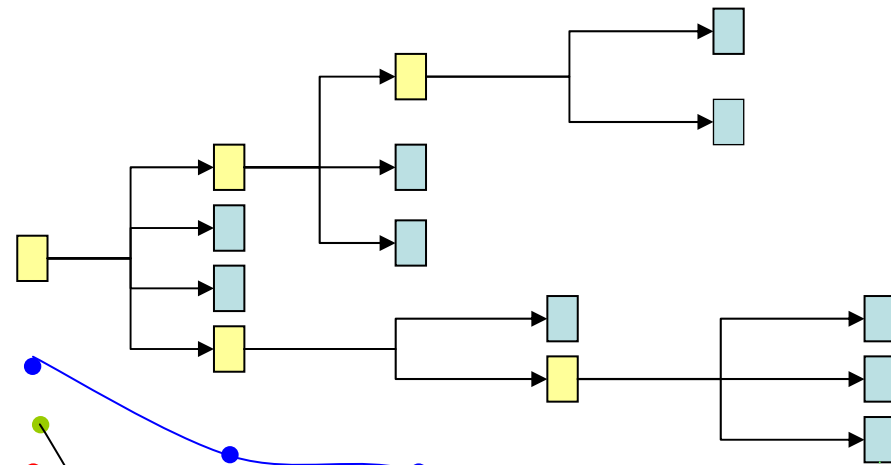
$\Omega_t$  sample test

$$\delta = \frac{\sum_{\omega \in \Omega_t} (C(\omega) \neq \varphi(\omega))}{|\Omega_t|}$$

$$(C(\omega) \neq \varphi(\omega)) = \begin{cases} = 1 \\ = 0 \end{cases}$$

$\Pi_*$

Size of the partitions  
Number of Terminal Nodes  
Number of rules



$$\theta(\Pi_i) \geq \theta(\Pi_{i+1})$$

$$E_i \geq E_{i+1}$$

$$\theta(\Pi) = \sum_{s \in S} w_s h(s)$$

$$E_{\Pi} = \sum_{s \in S} f_s \varepsilon_s$$

$\Omega_t$  sample test

$$\delta = \frac{\sum_{\omega \in \Omega_t} (C(\omega) \neq \varphi(\omega))}{|\Omega_t|}$$

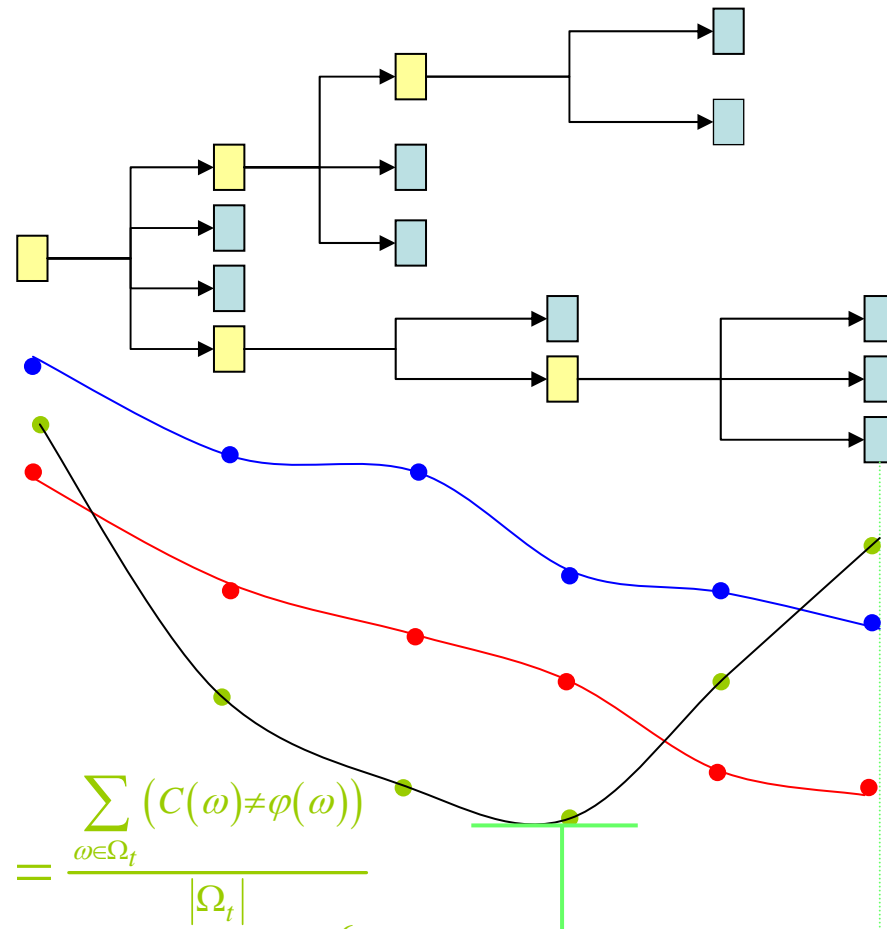
$$\theta(\Pi_i) \geq \theta(\Pi_{i+1})$$

$$E_i \geq E_{i+1}$$

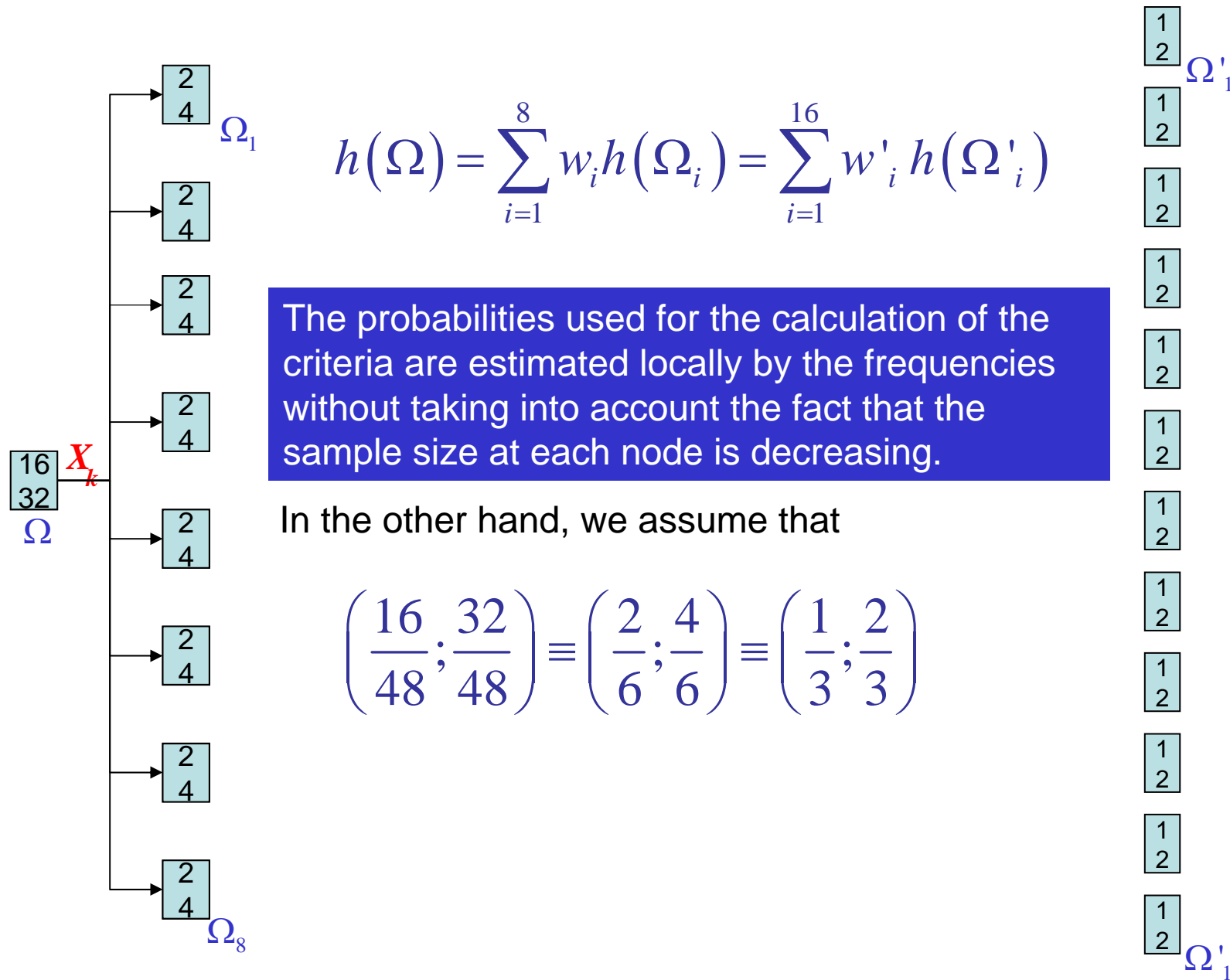
$$\theta(\Pi) = \sum_{s \in S} w_s h(s)$$

$$E_{\Pi} = \sum_{s \in S} f_s \varepsilon_s$$

aren't suitable for building the right tree



## 2. Insensitivity of the criterion to the sample size



The insensitivity of the criteria to the sample size may lead to different assumptions and solutions :

- We don't need any criteria nor algorithm for growing the tree, develop the largest one, we assume that there is a **large data set in each node** to make the estimates of the probabilities reliable ;
- Suppose that the **predictive attributes are independent** then, we can get a reliable estimates of the probabilities at each node by applying Bayes Formula;
- Start from the largest tree and apply one of the pruning techniques to find the right tree ;
- Fix a **minimum size** in each node before splitting;
- Introduce a **penalty parameter for the complexity** of the tree (complexity is the size of current partition), and allow merging modalities of predictive attributes, thus we better control the size of the current partition, we better use the sample;
- Use statistical criteria to stop the growing process at the right deep;

In all of these cases, the criteria for growing the tree remains insensitive to the sample size.

We aim to formulate a new criteria which depends on both the frequencies vector and the size of the sample on which the frequencies are estimated, thus we may avoid the over fitting without post treatment nor external parameterization.

### 3. Entropy measure sensitive to the sample size

$$\hat{h} : \Pi \rightarrow \mathbb{R}^+ \quad \Pi = \bigcup_n s_n$$

$$s_n = \left\{ \gamma \in \mathbb{R}^n / \gamma_i \geq 0 (i = 1, \dots, n) \text{ and } \sum_{i=1}^n \gamma_i = 1 \right\}$$

⇒ **Symetry**  $\hat{h}(\gamma_1, \gamma_2, \dots, \gamma_n) = \hat{h}(\gamma_{\sigma_1}, \gamma_{\sigma_2}, \dots, \gamma_{\sigma_n})$

⇒ **Minimality**  $\hat{h}(\gamma_1, \gamma_2, \dots, \gamma_n) = \text{Min} \Leftrightarrow \exists i / \gamma_i = 1 \forall j \neq i; \gamma_j = 0$

⇒ **Maximality**  $\hat{h}(\gamma_1, \gamma_2, \dots, \gamma_n) = \text{Max} \Leftrightarrow \forall i (i = 1, \dots, n) ; \gamma_i = \frac{1}{n}$

➡ **Sensitivité to the sample size**  $\hat{h}(\gamma_1, \gamma_2, \dots, \gamma_n, n_1) > \hat{h}(\gamma_1, \gamma_2, \dots, \gamma_n, n_2)$   
if  $n_1 < n_2$

➡ **Convergence**  $\lim_{n \mapsto \infty} \hat{h}(\gamma_1, \gamma_2, \dots, \gamma_n, n) \rightarrow h(p_1, p_2, \dots, p_n)$

⇒ **Concavity**

$$\hat{h}(\Omega) = \hat{h}(\gamma_1, \dots, \gamma_i, \dots, \gamma_m, n) = h(\hat{\gamma}_1, \dots, \hat{\gamma}_i, \dots, \hat{\gamma}_m)$$

were

$h$  is any entropy measure

$$\hat{\gamma}_i = \frac{n_i + \lambda}{n + m.\lambda}$$

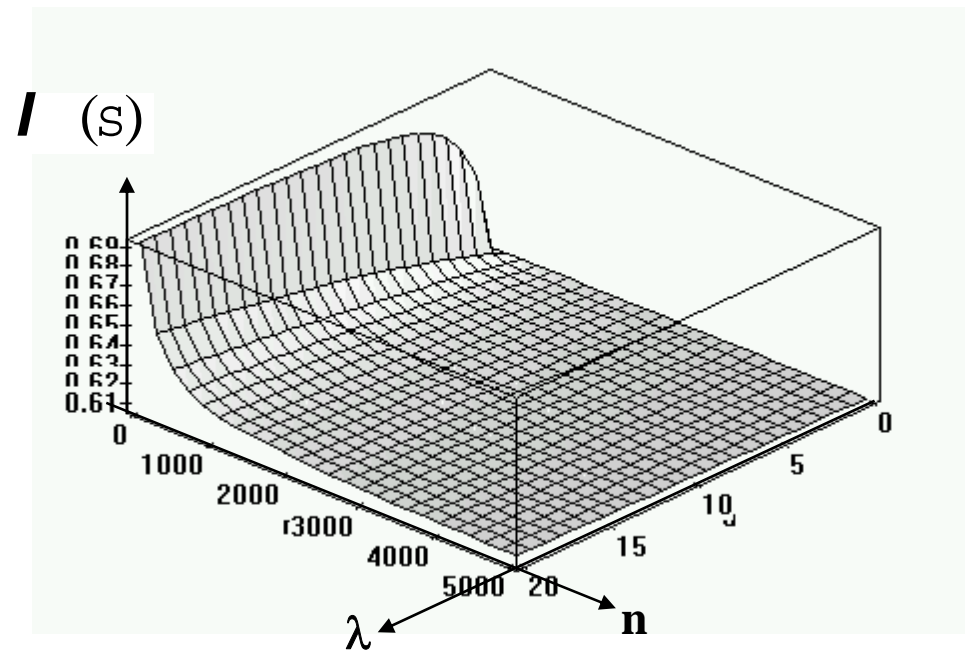
were  $\lambda > 0$

$$\hat{h}(\Omega) = - \sum_{i=1}^m \frac{n_i + \lambda}{n + m.\lambda} \log \frac{n_i + \lambda}{n + m.\lambda}$$

$$\hat{h}(\Omega) = \sum_{i=1}^m \frac{n_i + \lambda}{n + m.\lambda} \left( 1 - \frac{n_i + \lambda}{n + m.\lambda} \right)$$

$$\hat{h}(\Omega) = \sum_{j=1}^k \frac{n_{\cdot j}}{n} \left[ - \sum_{i=1}^m \frac{n_{ij} + \lambda}{n_{\cdot j} + m\lambda} \text{Log} \left( \frac{n_{ij} + \lambda}{n_{\cdot j} + m\lambda} \right) \right]$$

$S_1$	$S_2$
0.1	0.4
0.3	0.2





# How to fix $\lambda$

$\tau$  soft constraint

$$\begin{array}{cc}
 \begin{array}{c} \mathbf{v} \\ \boxed{\begin{array}{c} \tau + 1 \\ 0 \\ \vdots \\ 0 \end{array}} \\ n_{\cdot \mathbf{v}} = \tau + 1 \end{array} & 
 \begin{array}{c} \mathbf{u} \\ \boxed{\begin{array}{c} \tau \\ 0 \\ \vdots \\ 0 \end{array}} \\ n_{\cdot \mathbf{u}} = \tau \end{array}
 \end{array}$$

$$\lambda^* : \hat{h}(u_{\lambda^*}) - \hat{h}(v_{\lambda^*}) = \max_{\lambda} (\hat{h}(u_{\lambda}) - \hat{h}(v_{\lambda}))$$

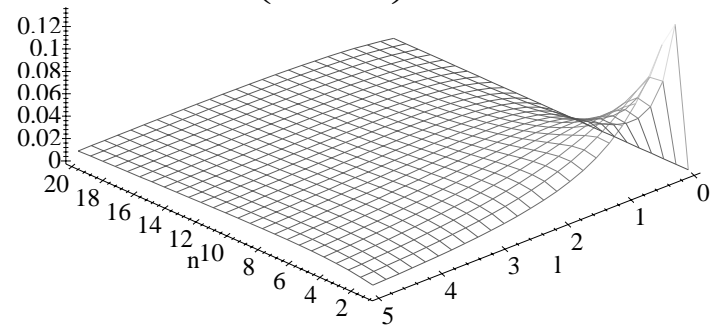
$$\begin{aligned}
 f(\lambda) &= (\hat{h}(u_{\lambda}) - \hat{h}(v_{\lambda})) \\
 &= (m-1) \frac{2\tau^2 + 2\tau + m\lambda(1+2\tau)}{(\tau + m\lambda)^2 (\tau + 1 + m\lambda)^2}
 \end{aligned}$$

$$\text{Find } \lambda^* : f(\lambda^*) \geq f(\lambda); \lambda \neq \lambda^*$$

$$\frac{\partial f(\lambda)}{\partial \lambda} = 0$$

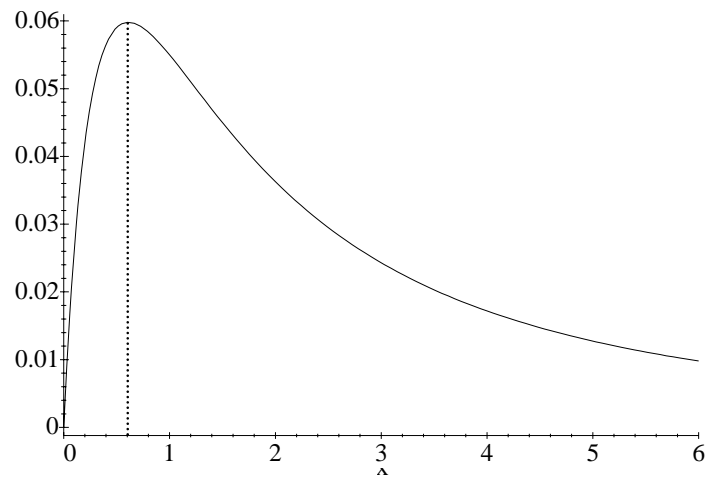
$$m = 3$$

$$\varphi(\lambda, \tau)$$



$$m = 2$$

$$\tau = 2$$



$$\lambda^* = 0.61098$$

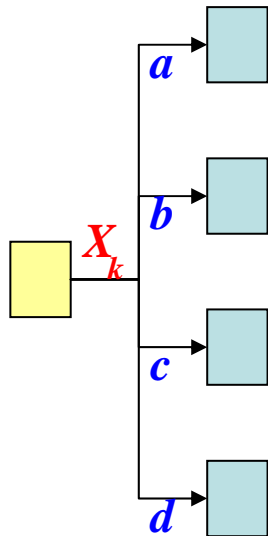
$$f(\lambda) = \lambda \left[ \frac{5\lambda + 6}{2 \left( 1 + \lambda^2 (3 + 2\lambda)^2 \right)} \right]$$

# 4. Lattice to learn from finite data set

Merging nodes reduces the size of the partition and allows deeper exploration

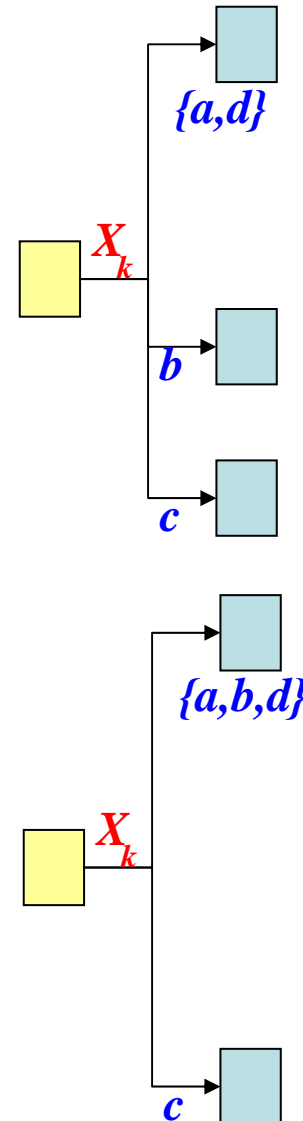
## CHAID (Kass, 1980)

Looks for the best **partition** on the set of observed modalities of  $X_k$  in the current node



Looks for the the best **bi-partition** on the set of observed modalities of  $X_k$  in the current node

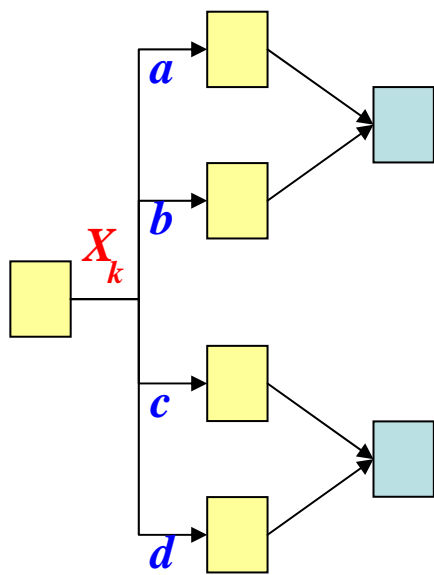
## CART (Breiman et al. 84)

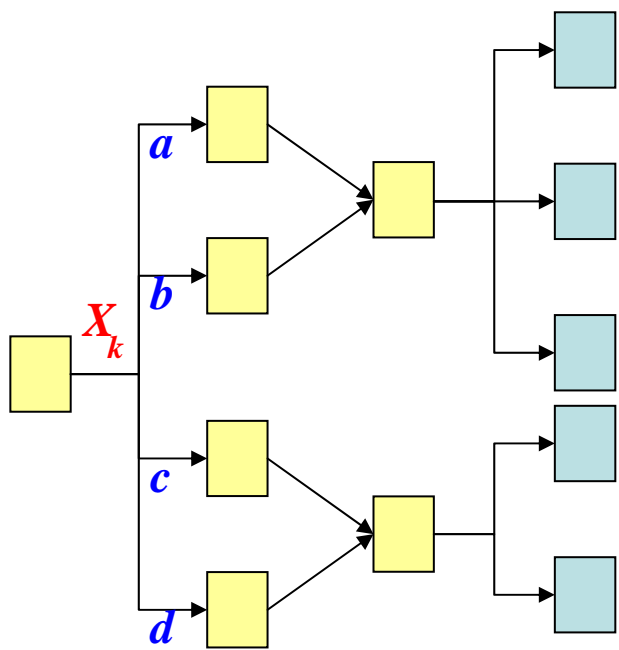


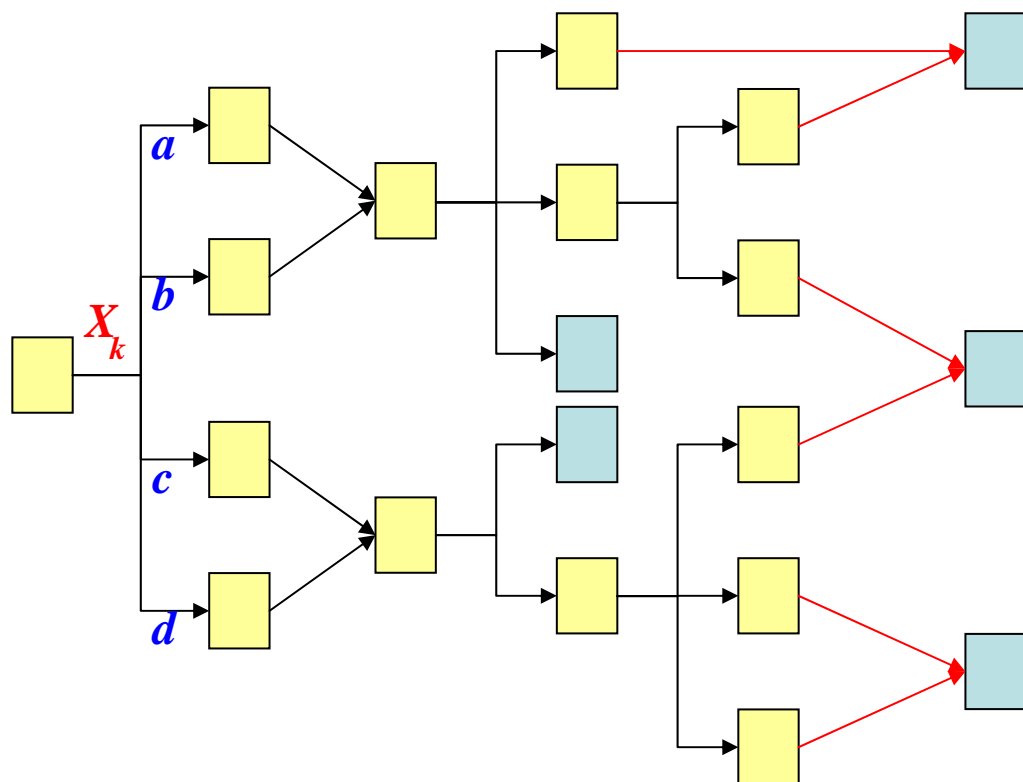
N-ary trees

The merge is allowed only within brother nodes

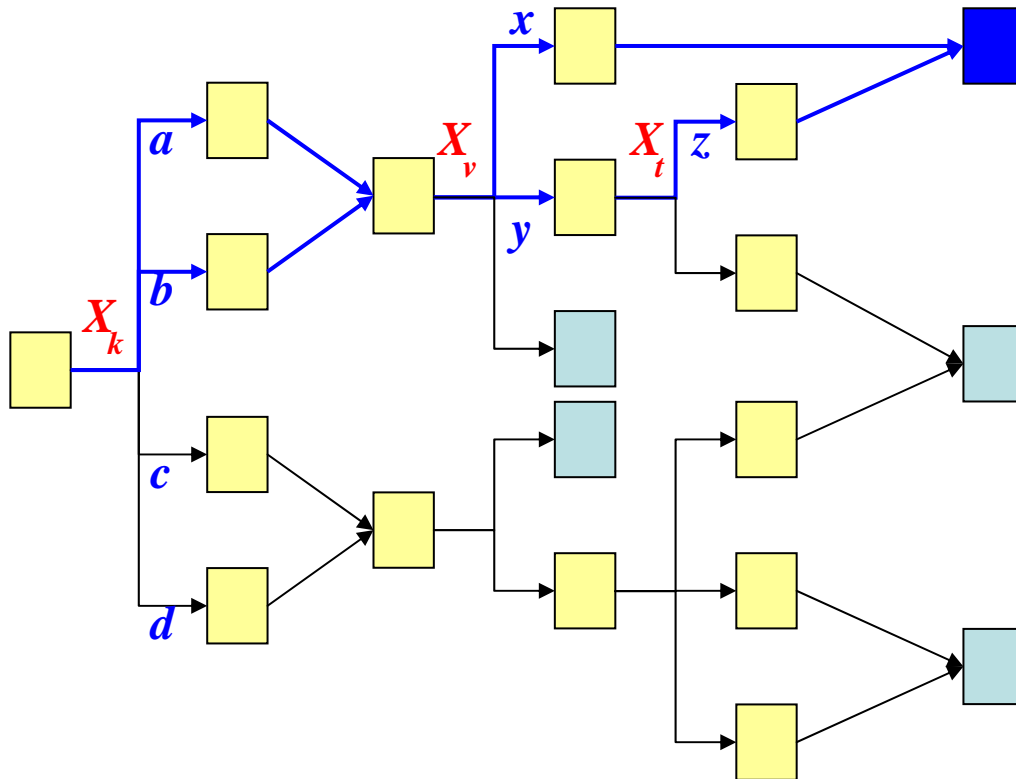
Binary trees

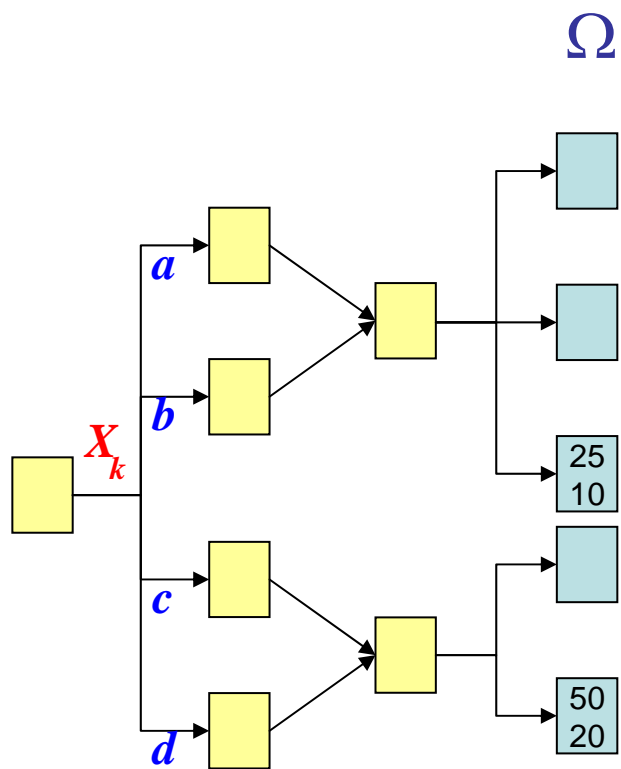






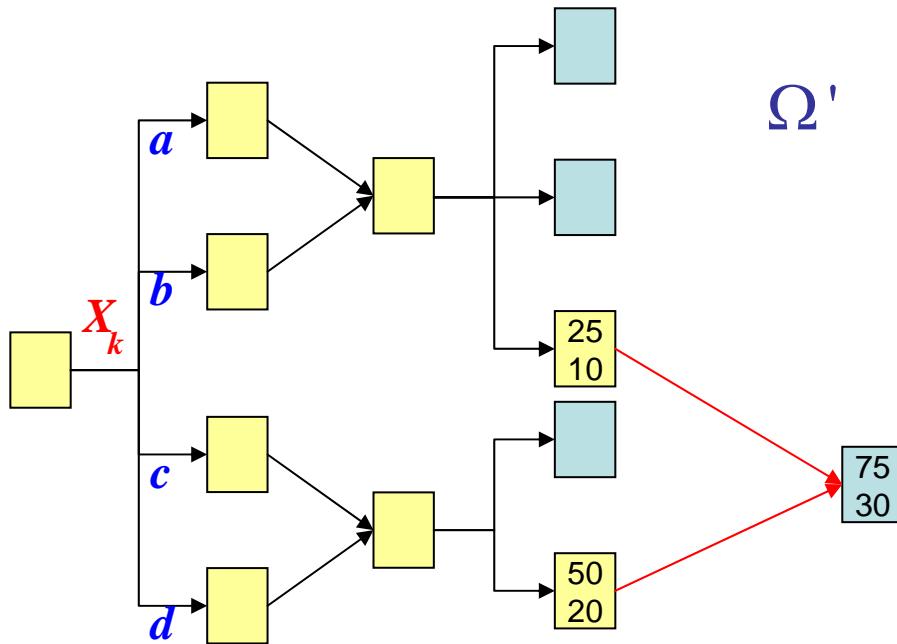
If  $\left[ \left( X_k \in \{a, b\} \right) \wedge \left( \left( X_v = x \right) \vee \left( \left( X_v = y \right) \wedge \left( X_t = z \right) \right) \right) \right]$  Then  $C = c_i$





$$\hat{h}(\Omega) = \sum_{i=1}^5 w_i \hat{h}(\Omega_i)$$





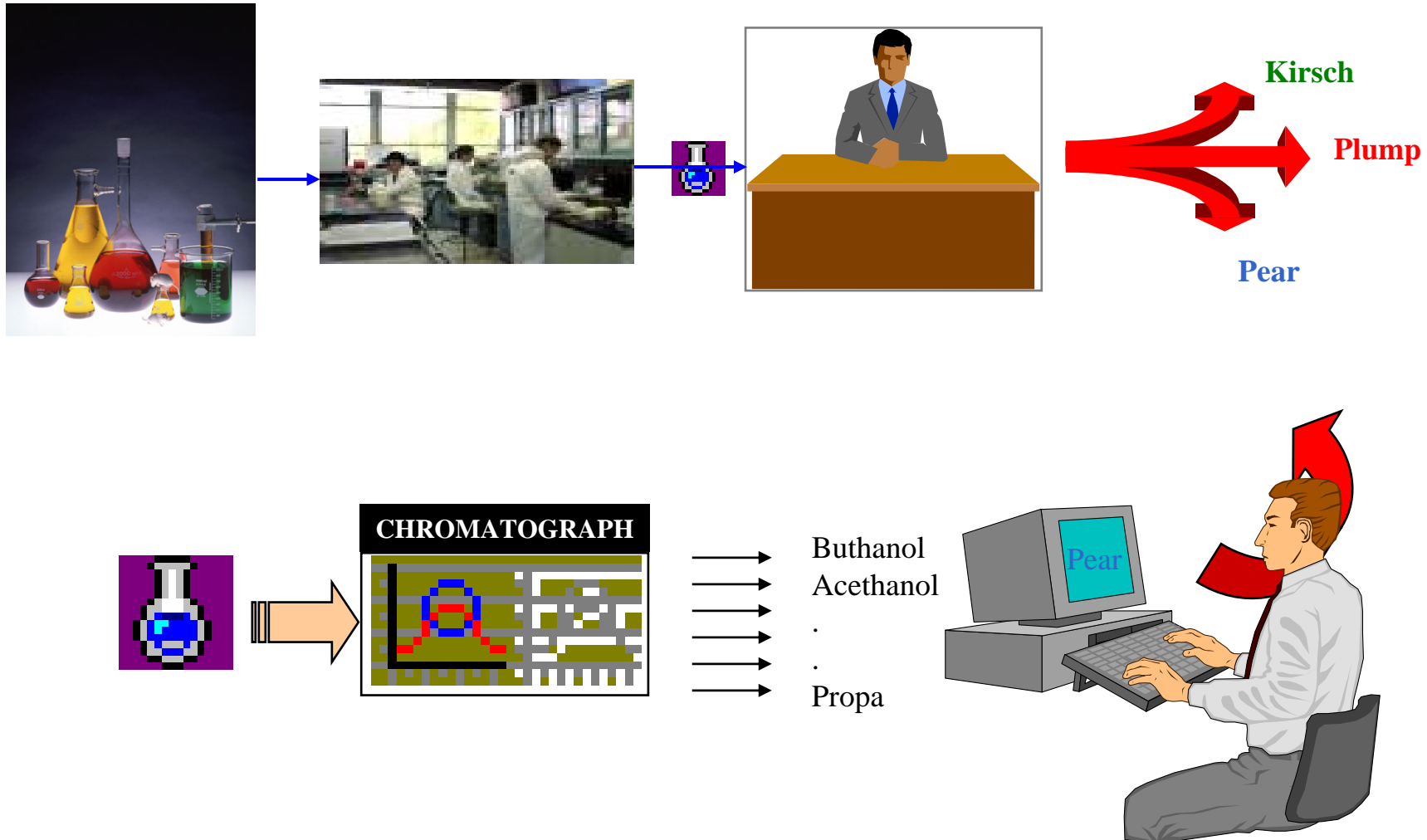
$$\hat{h}(\Omega') = \sum_{i=1}^4 w_i \hat{h}(\Omega_i)$$

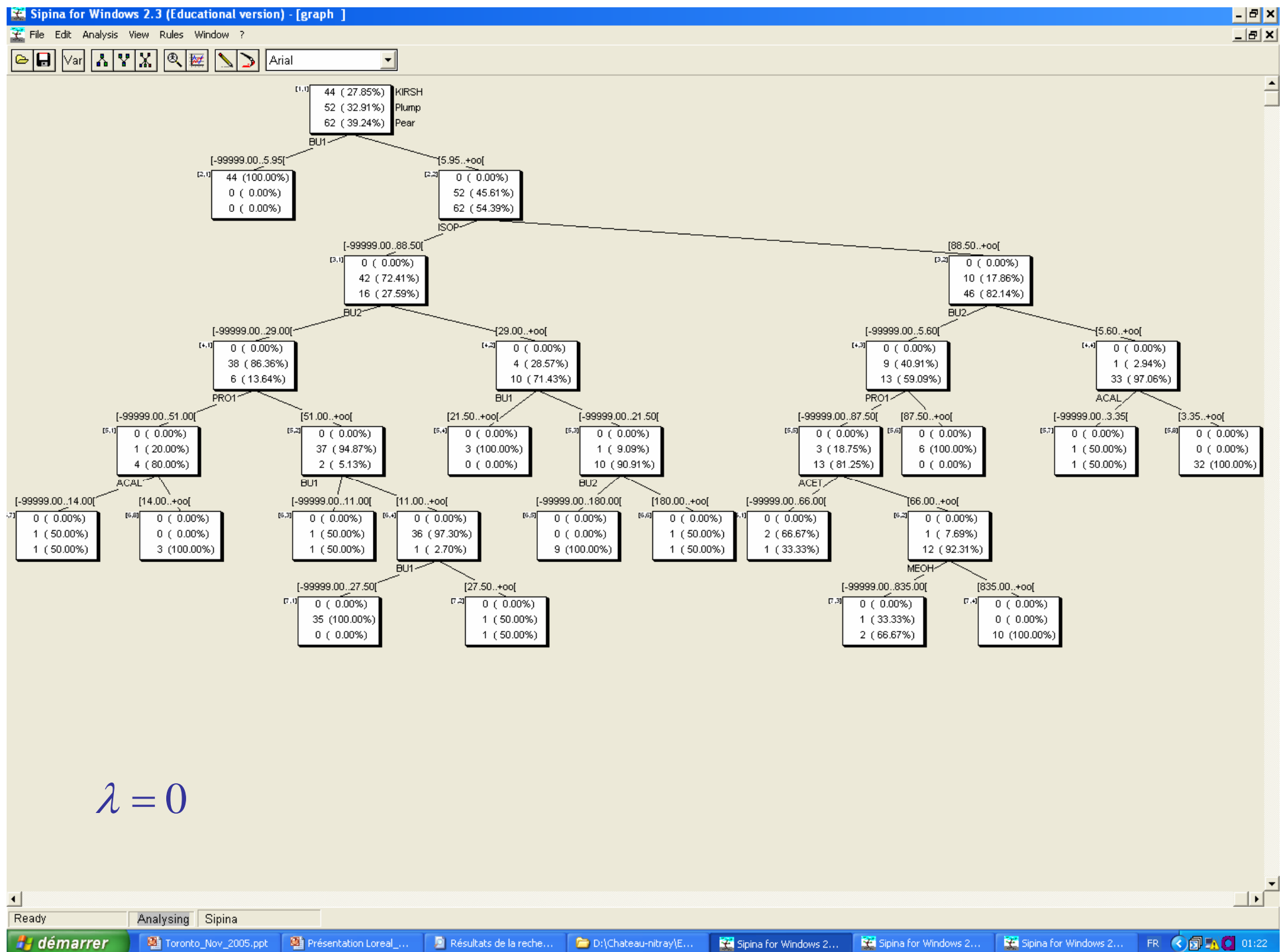
Merging such nodes must decrease the entropy

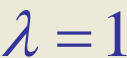
$$\hat{h}(\Omega') < \hat{h}(\Omega)$$

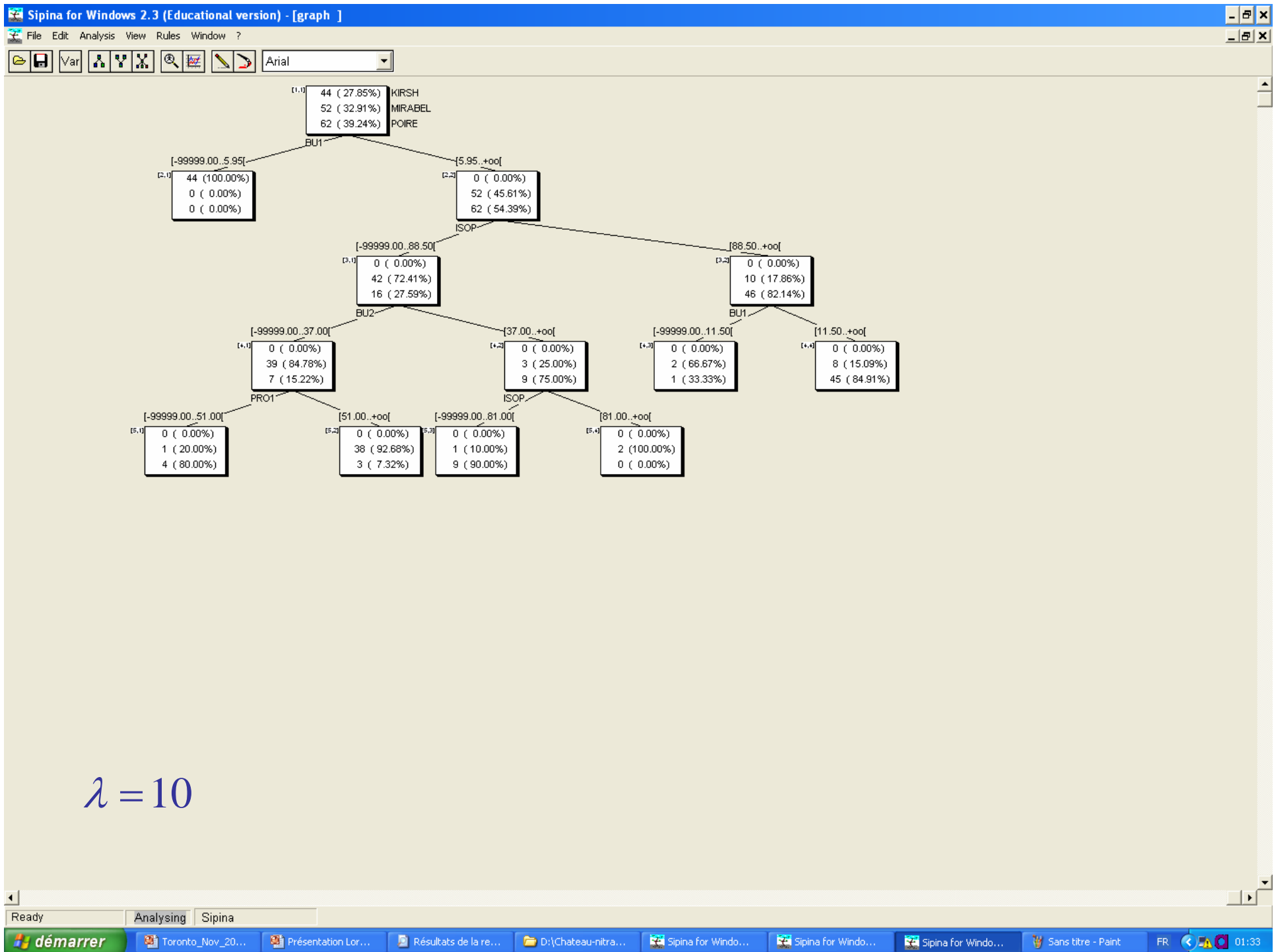
# 5. Example

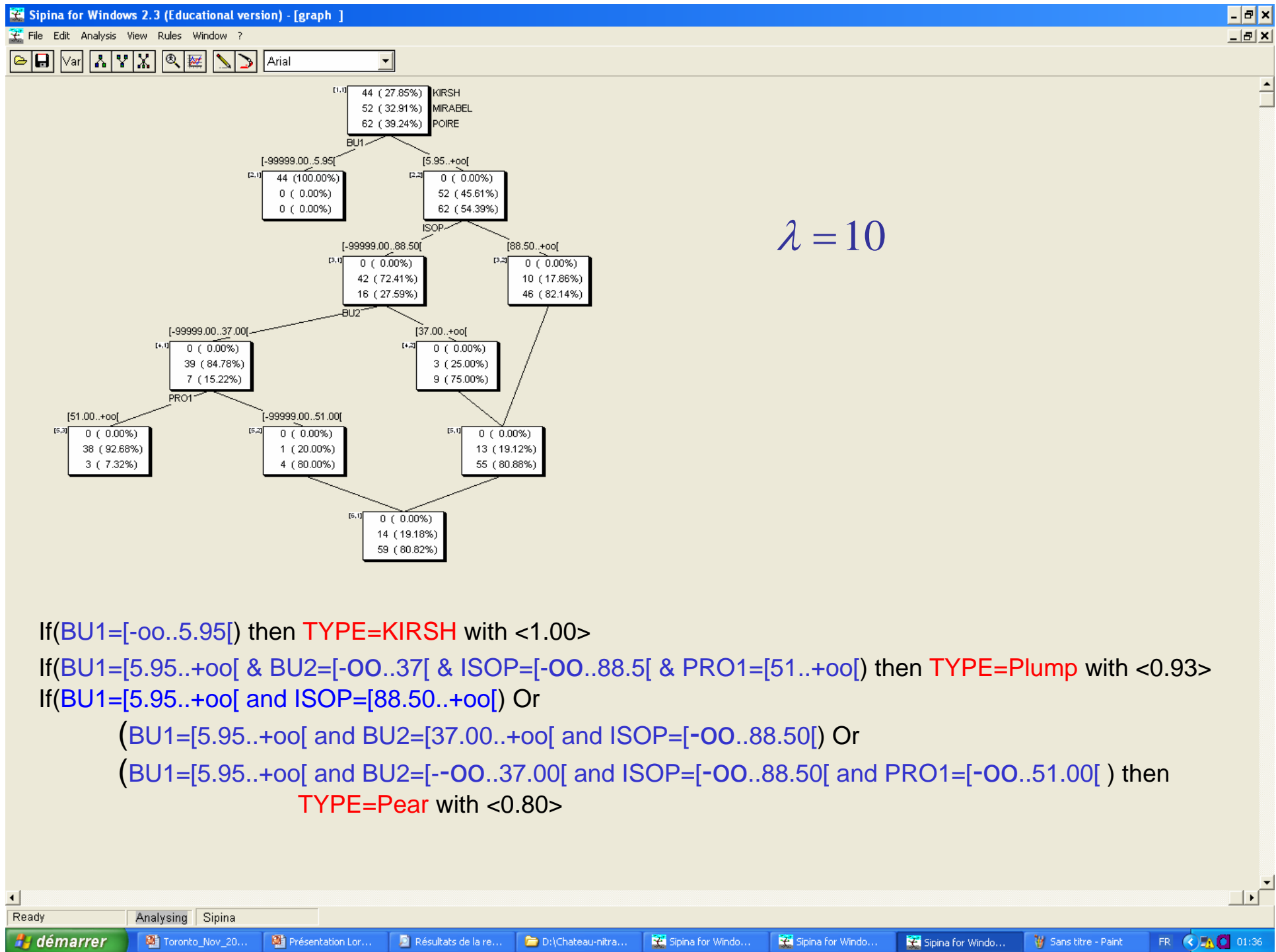
## Quality Control of brandy











## 6. Conclusion