

CART : Classification And Regression Trees

The context :

We have data $\mathcal{L} = \{(x_1, y_1), \dots, (x_n, y_n)\}$
where $x_i = (x_{i1}, \dots, x_{ip})$: the vector that
contains the observations

of the P explanatory variables

for the individual i .

y_i : the observation of the response
variable for the individual i .

In regression: $\forall i \in \{1, \dots, n\},$

$$y_i \in \mathbb{R}$$

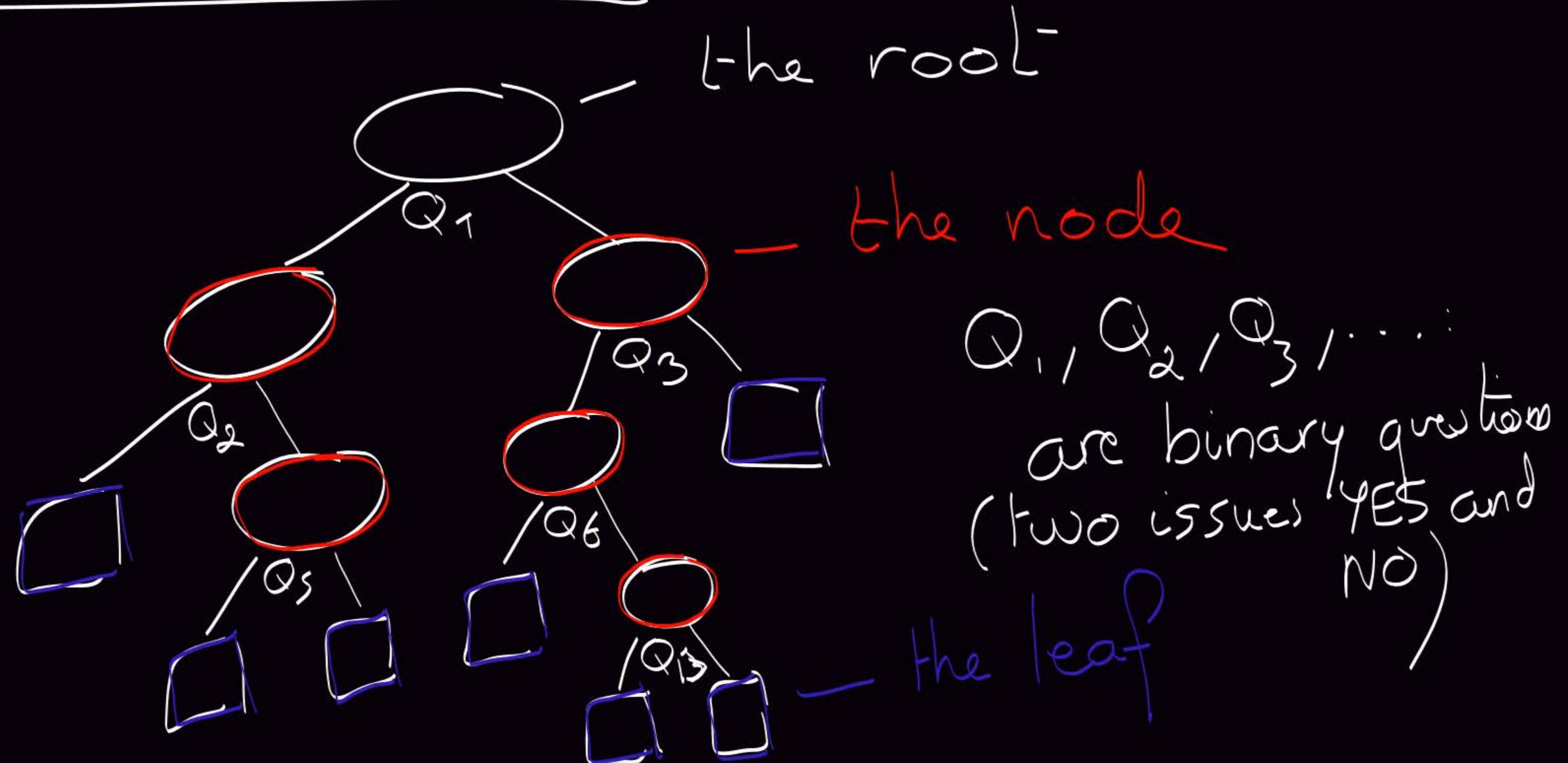
In classification: $\forall i \in \{1, \dots, n\},$

$$y_i \in \mathcal{C} = \{1, \dots, I\}$$

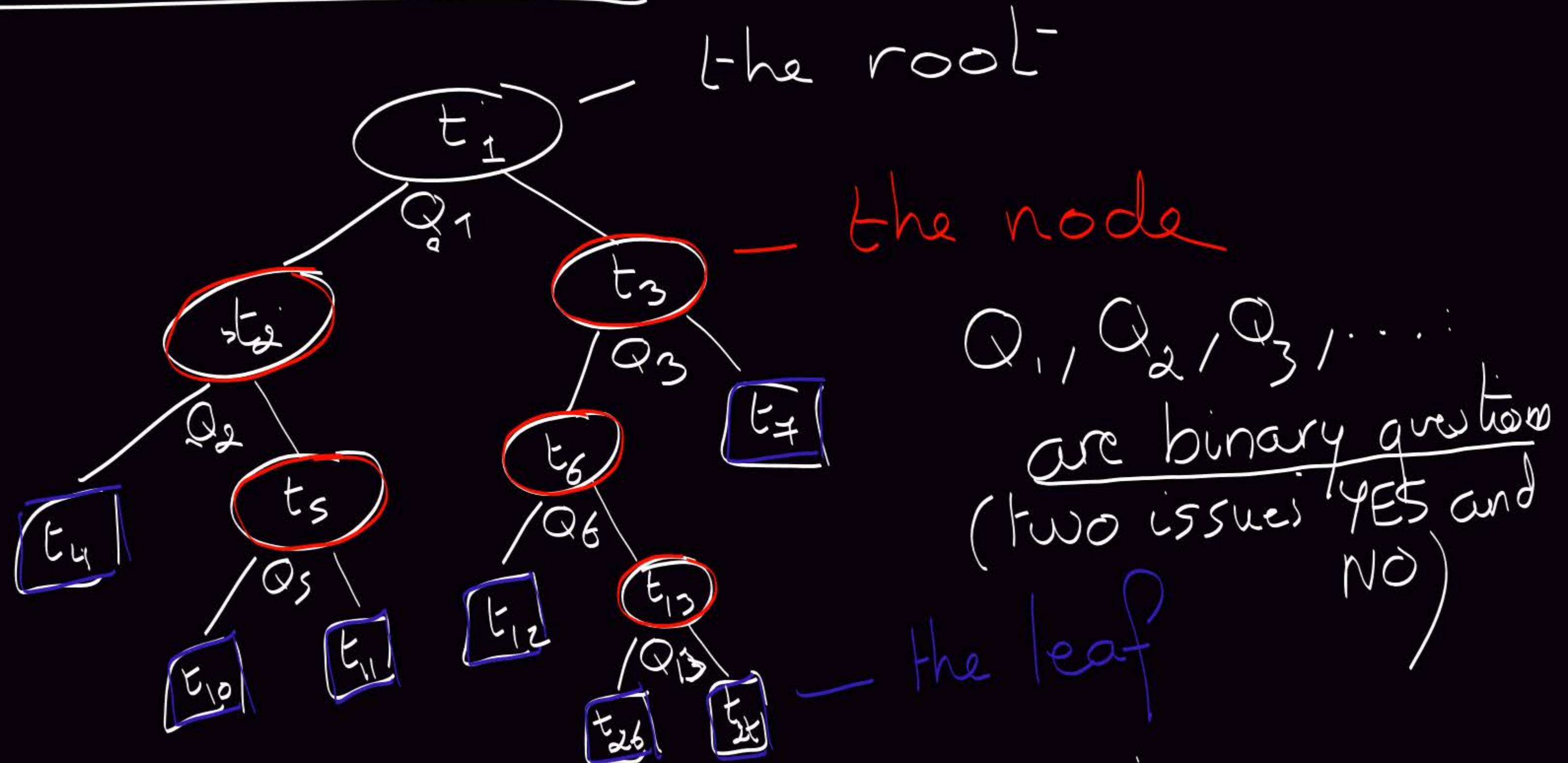
$\Rightarrow y$ is a qualitative variable.

In both cases, the goal is the same:
Find a function that makes a link
between the response variable and
the explanatory ones.

what is the Tree ?



what is the Tree ?



- ?
 - 1) How to define the binary questions to split the node?
 - 2) How to decide that a node is a leaf?
 - 3) How this tree can help me to have an idea of the link between the response and the explanatory variables?

↳ How to define the binary questions?

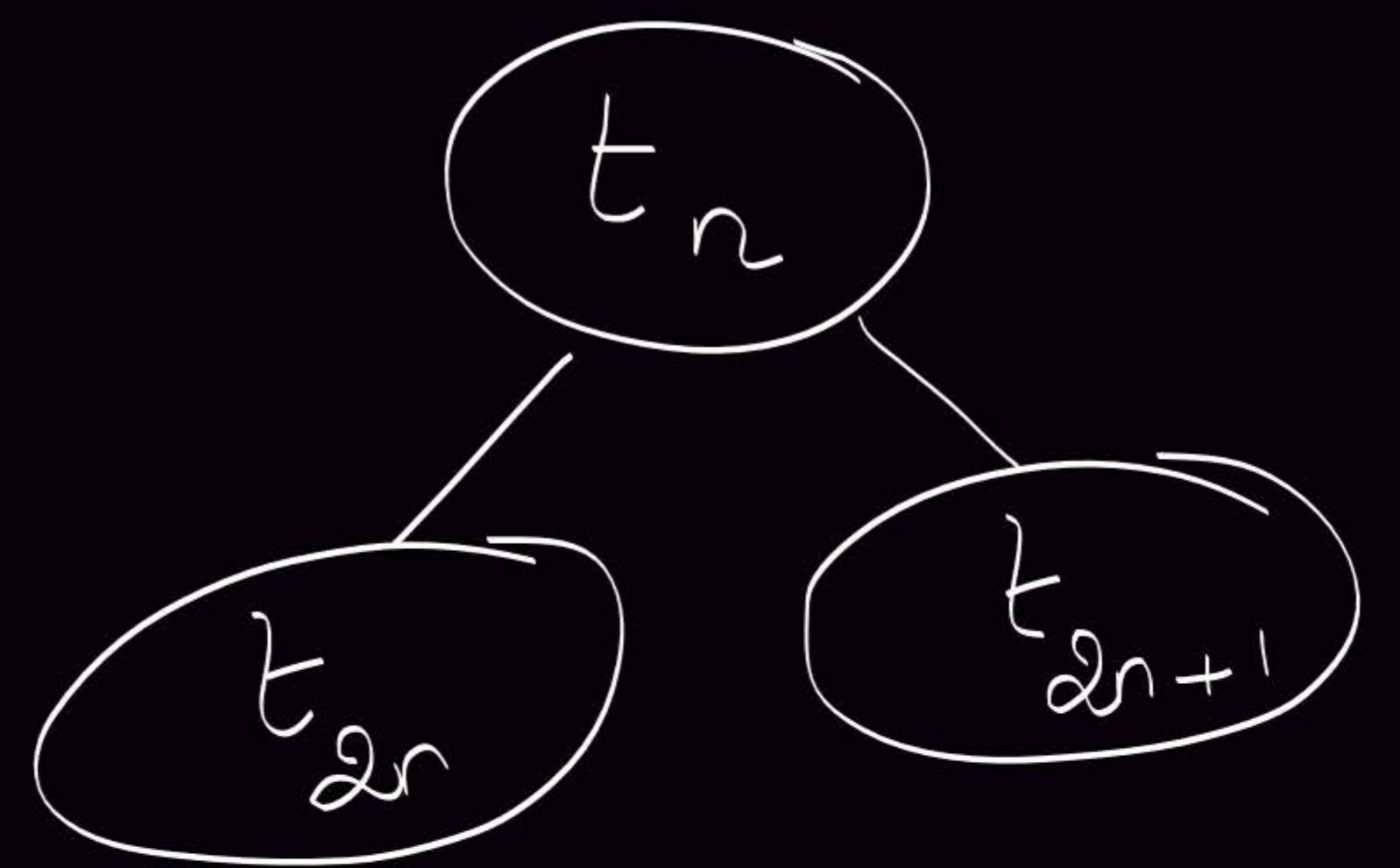
The binary questions are

$(x^i < a_i)$ or $(x^i > a_i)$ ← for x^i
which numeric

or $(x^i \in G_i)$ ← for x^i
which is qualitative

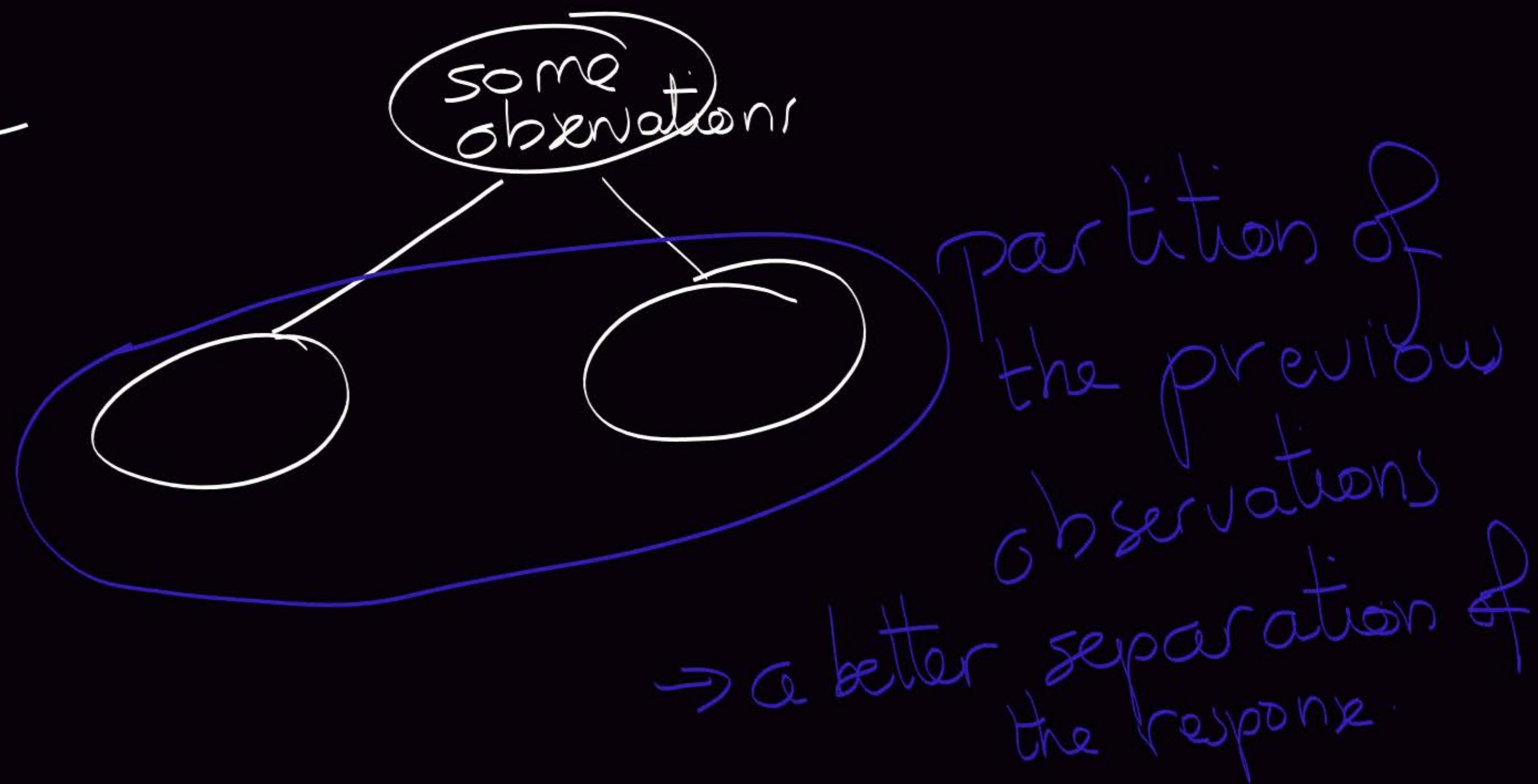
2 questions now that we
know the form of the binary
question:

- 1) How to define the variable
involved in the question?
- 2) How to define a_i or g_i ?



The answer to those two questions
will depend on the framework.

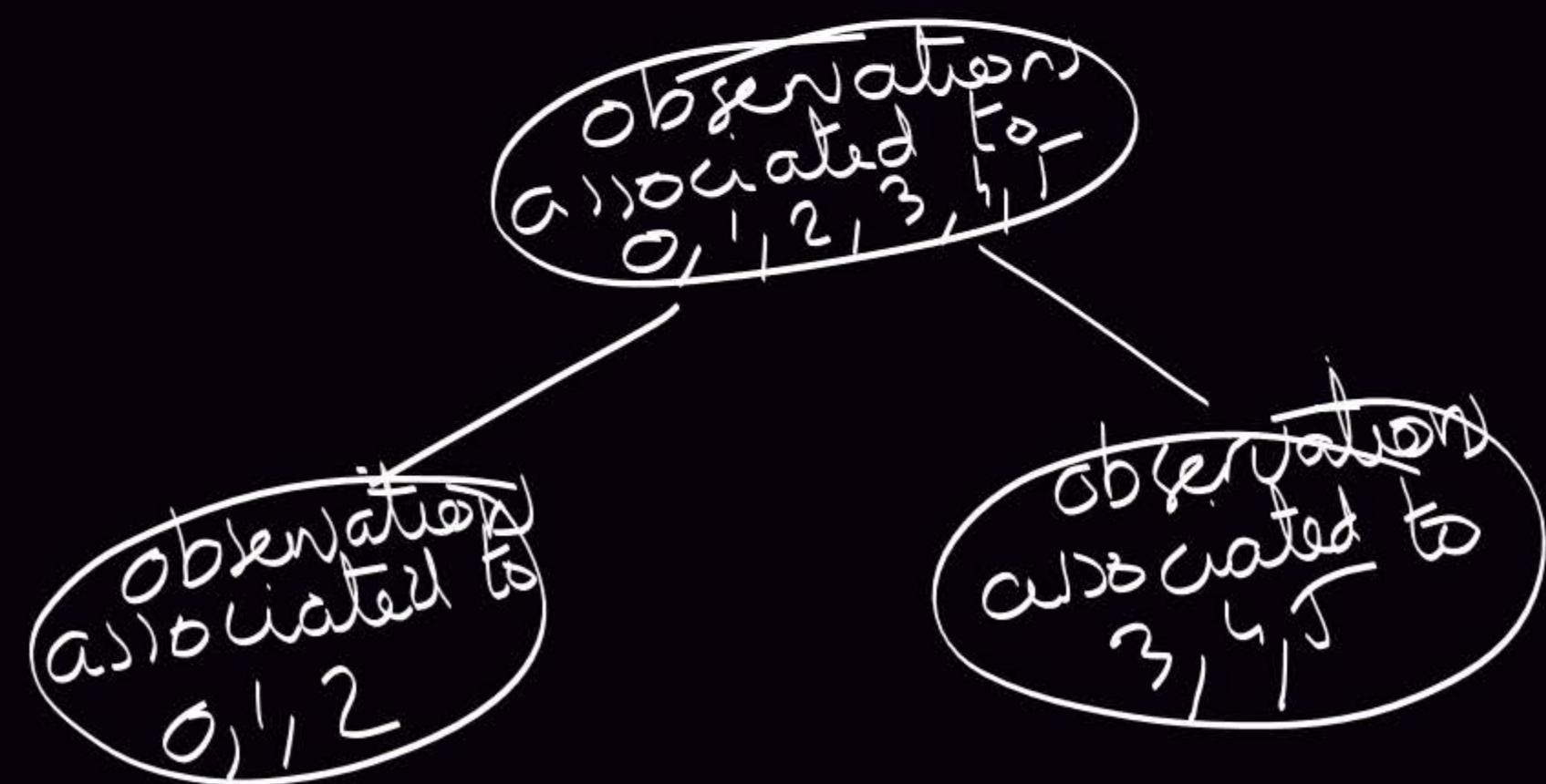
Why?



Example :

In classification
judgment for a film.

↳ 0 1 2 3 4 5

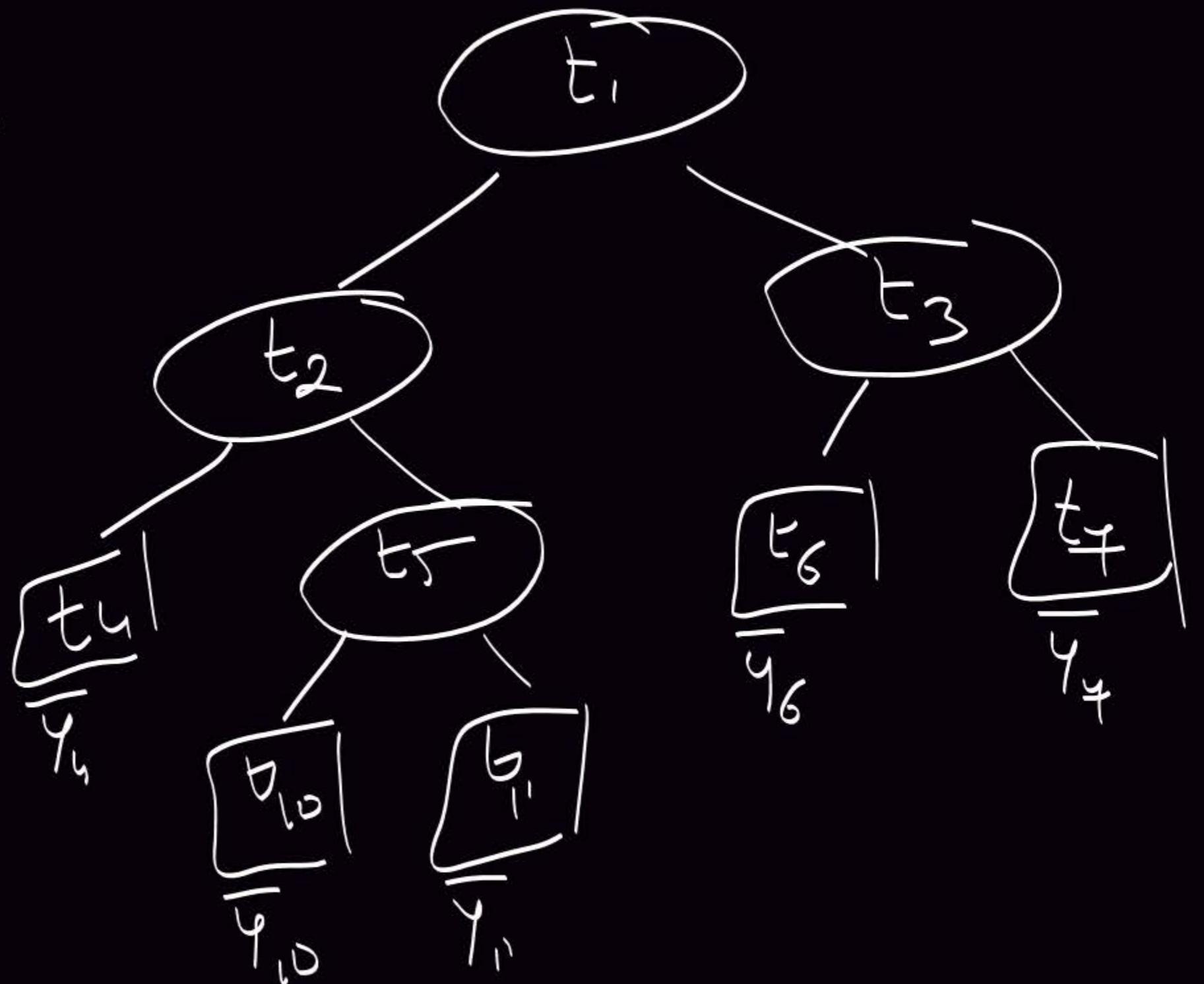


to be able to give an answer
to those two points, we need

to define :

- the notion of predicted value associated to a node
- a criteria to measure the purity of a node

Regression framework

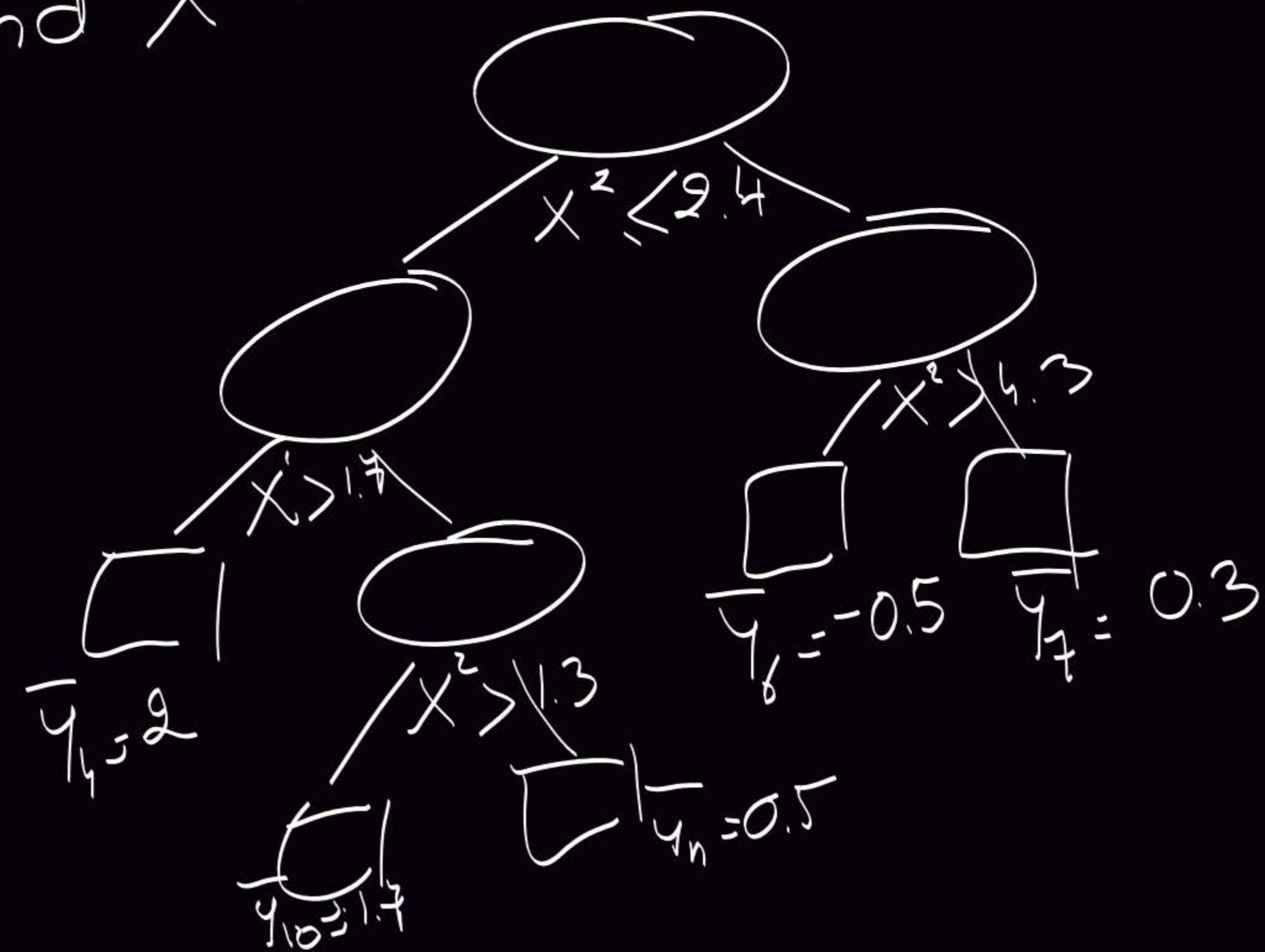


$$Y = f(X) + \epsilon$$

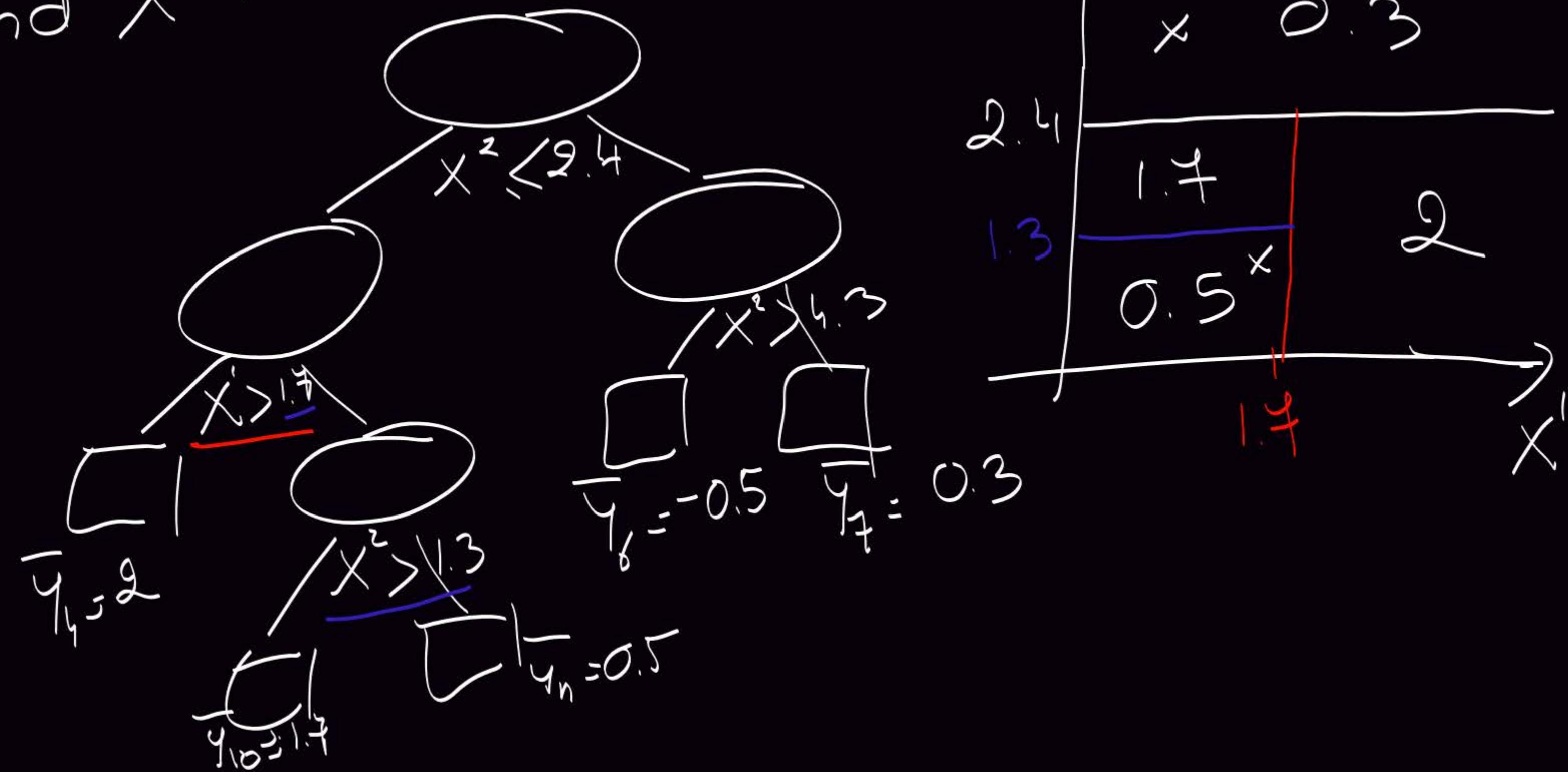
with $X = (X^1, X^P)$

$$\hat{Y} = \begin{cases} \bar{Y}_4 & \text{if } Q_1 \text{ is YES} \\ \bar{Y}_{10} & \text{if } Q_1 \text{ is YES and } Q_2 \text{ is YES} \\ \bar{Y}_{10} & \text{if } Q_1 \text{ is YES and } Q_2 \text{ is NO} \\ \bar{Y}_7 & \text{if } Q_2 \text{ is YES and } Q_5 \text{ is YES} \\ \bar{Y}_{11} & \dots \\ \bar{Y}_6 & \dots \\ \bar{Y}_7 & \dots \end{cases}$$

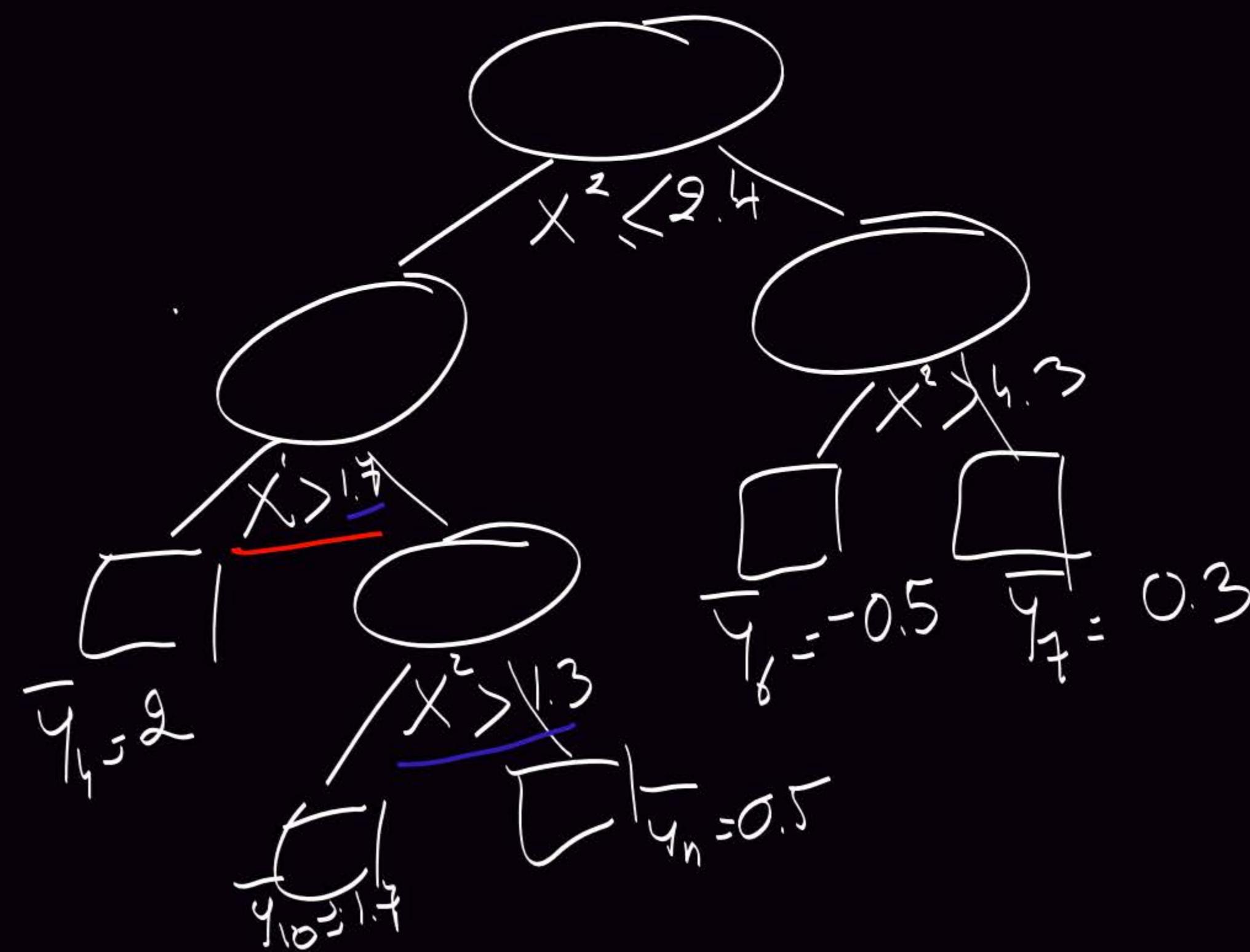
Example:
a data set with 2 quantitative variables X^1
and X^2 .



Example:
 a data set with 2 quantitative variables X^1
 and X^2 .



$$\hat{g}(x_1, x_2) = \begin{cases} 2 & \text{if } x_2 \leq 2.4 \\ & \text{and } x_1 \geq 1.7 \\ 1.7 & \text{if } 1.3 \leq x_2 \leq 2.4 \\ & \text{and } x_1 \leq 1.7 \\ 0.5 & \text{if } x_1 \leq 1.7 \\ & \text{and } x_2 \leq 1.3 \\ -0.5 & \text{if } x_2 > 1.3 \\ 0.3 & \text{if } 2.1 \leq x_2 \leq 4.3 \end{cases}$$



Assume that the form of the tree
is determined.

How to define the response (predicted)

associated to each leaf?

$$\text{let } R(\hat{f}) = \frac{1}{n} \sum_{i=1}^n (\hat{f}(x_i) - f(x_i))^2$$

we want to minimize $R(\hat{f})$

Since $R(\hat{f})$ is unknown, we
replace $R(f)$ by $R_n(\hat{f})$

defined by :

$$R_n(\hat{f}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2.$$

Prop:

Let \bar{T} denotes a tree.

Let \tilde{T} denotes the leaves of \bar{T} .

We have $\hat{f}(x) = \sum_{t \in \tilde{T}} \bar{y}_t \times 1|_{x \in t}$

To minimize $R_n(\hat{f})$, we need to have

$$\bar{y}_t = \frac{1}{n_t} \times \sum_{\substack{i \in \{1, \dots, n\} \\ x_i \in t}} y_i$$

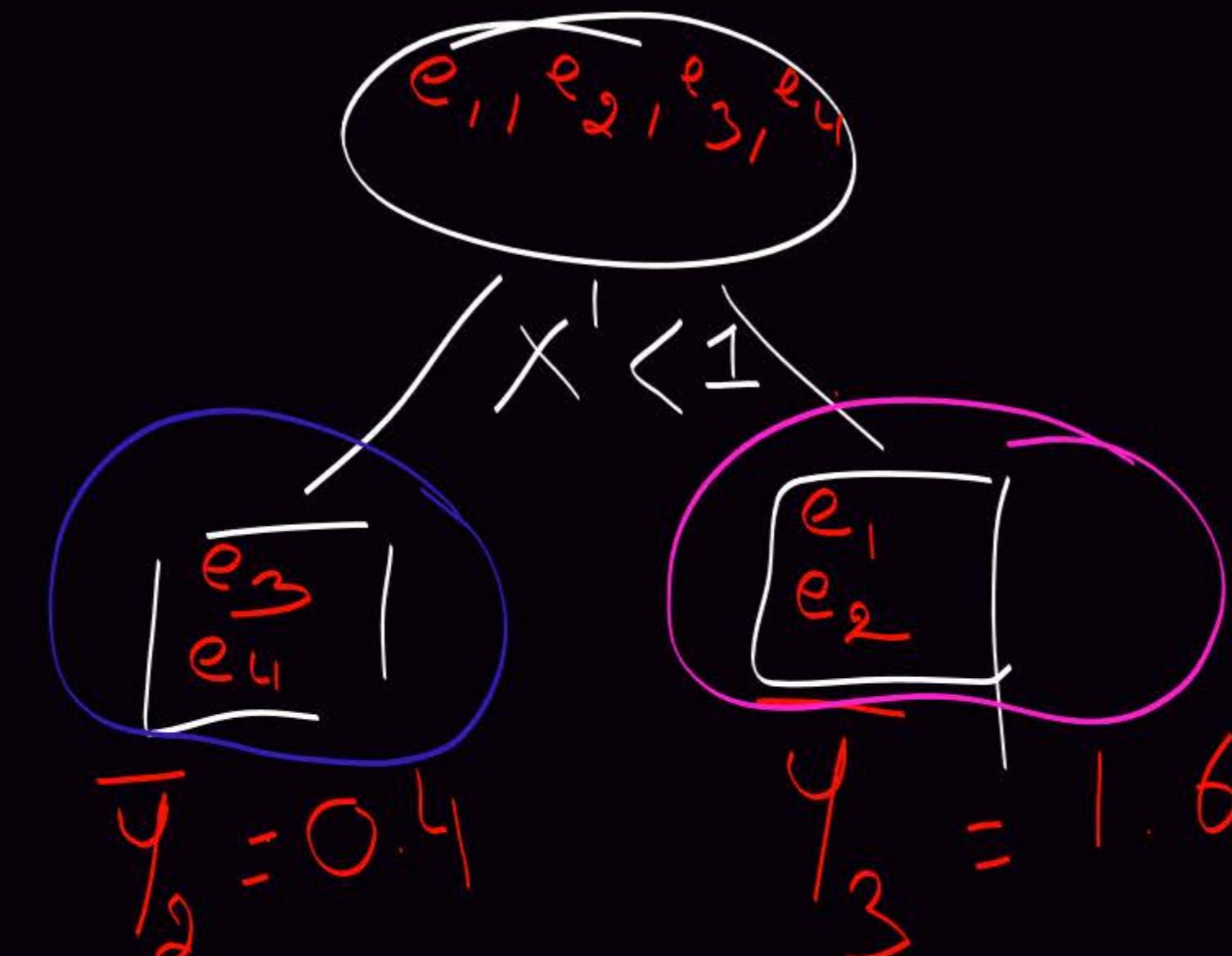
with n_t : number of observations from \mathcal{S} that are in t

$$\frac{R_B}{\pi} : \\ x \in A = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{otherwise} \end{cases}$$

Example :

dataset

	x^1	x^2	y
e_1	1.4	0.2	1.2
e_2	1.5	-0.2	2
e_3	0.8	1.2	0.5
e_4	0.4	2.4	0.3



Proof:

$$\begin{aligned} R_n(\hat{f}) &= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \\ &= \frac{1}{n} \sum_{t \in T} \left(\sum_{i \in \{1, \dots, n\}} (y_i - \hat{f}(x_i))^2 \right) \end{aligned}$$

To minimize $R_n(\hat{f})$, we just have to minimize each term of the sum (because each term is positive and does not involve observation used in the other terms).

$\forall t \in T$, we have to minimize

$$\sum_{\substack{i \in \{1, \dots, n\} \\ x_i \in t}} (y_i - \hat{f}(x_i))^2$$

But

$$\sum_{\substack{i \in \{1, \dots, n\} \\ x_i \in t}} (y_i - \hat{f}(x_i))^2 = \sum_{\substack{i \in \{1, \dots, n\} \\ x_i \in t}} (y_i - \bar{y}_t)^2$$

$$= \sum_{x_i \in t} (y_i^2 - 2y_i \bar{y}_t + \bar{y}_t^2)$$

$$= n_t \bar{y}_t^2 - 2\bar{y}_t \sum_{x_i \in t} y_i + \sum_{x_i \in t} y_i^2$$

$$\begin{aligned}
 &= n_t \left(\bar{y}_t^2 - 2\bar{y}_t \left(\frac{1}{n_t} \sum_{\substack{i \\ X_i \in t}} y_i \right) + \sum_{\substack{i \\ X_i \in t}} y_i^2 \right) \\
 &= n_t \left(\left(\bar{y}_t - \frac{1}{n_t} \sum_{\substack{i \\ X_i \in t}} y_i \right)^2 + \frac{1}{n_t} \sum_{\substack{i \\ X_i \in t}} y_i^2 \right)
 \end{aligned}$$

This should be equal to 0 to solve the problem of minimization

$$g(\bar{y}_t) = f(\bar{y}_t) + \beta \left(\begin{array}{l} \text{does not} \\ \text{depend} \\ \text{on } \bar{y}_t \end{array} \right)$$

To construct the tree:

The criteria is proportional to
the variance.

let t denotes a node of \bar{T}

$$R(t) = \frac{1}{n} \sum_{x_i \in t} (\bar{y}_t - y_i)^2$$

let t_L and t_R two children of
the node t , children generated by
a binary question δ .

The set of all possible binary questions
is S

$$S = \bigcup_{i \in \{I, P\}} \left\{ \begin{array}{l} (x^i \leq a_i, a_i \in \{s_i^{(1)}, \dots, s_i^{(n-1)}\}) \\ \text{or } (x^i > a_i, a_i \in \{s_{i+1}^{(1)}, \dots, s_{i+1}^{(n-1)}\}) \\ \text{or } (x^i \in \mathcal{E}^i, \mathcal{E}^i \subseteq \{I, \overline{I}\}) \end{array} \right\}$$

where $\{I, \overline{I}\}$: the set of modalities for x^i
 when x^i is a qualitative variable

x_1^c, x_n^c are the observations
of the variable X^c

let $x_{(1)}^c \leq x_{(2)}^c \leq \dots \leq x_{(n)}^c$ the same
observations
written in
ascending order.

$$S_i^{(k)} = \frac{1}{2} \left(x_{(k)}^c + x_{(k+1)}^c \right)$$

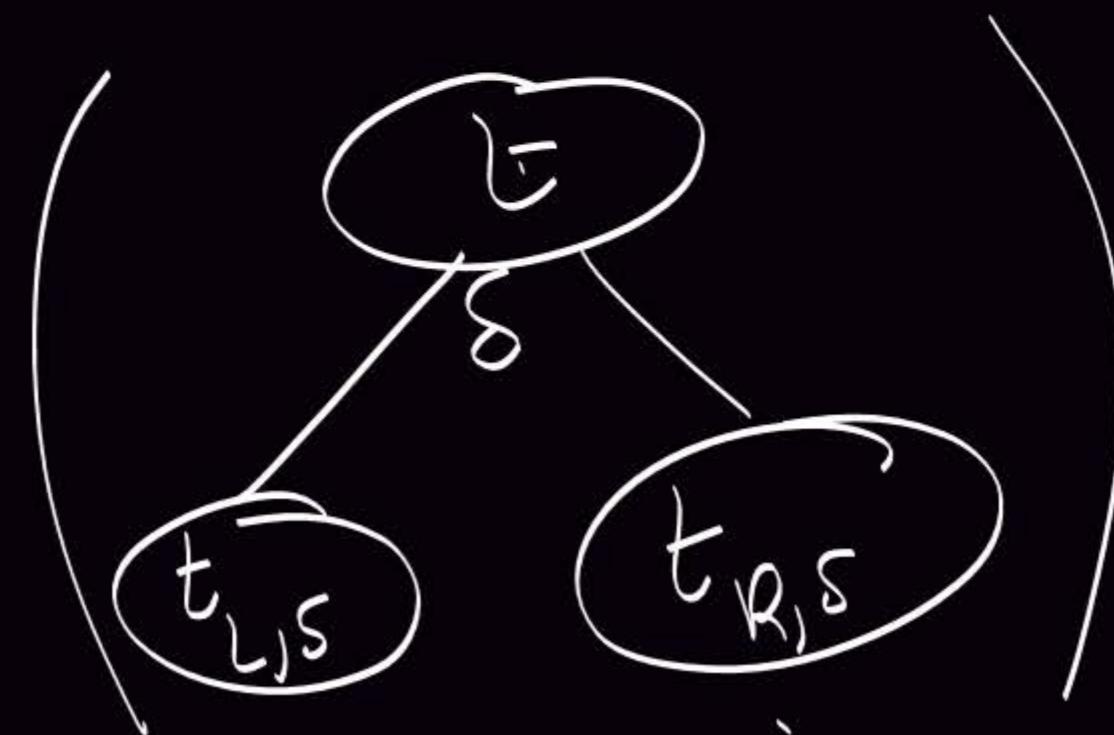
Rk: let X' be a quantitative variable

At the maximum, we have $I(n-1)$ possible divisions involving X'

Let X' be a qualitative variable.
We have $2^{I_i} - 2$ possible divisions
(I_i being the number of modalities of X')

Let t a node that we want to split.

For each possible division s of t
we compute $\Delta R(t, \delta) = R(t) - R(t_{l,s}) + R(t_{r,s})$



The final decision $\delta^*(t)$, it means
the decision that appears in
the tree to split the node t ,

defined by:

$$\delta^*(t) = \underset{\delta \in S}{\operatorname{argmax}} \Delta R(t, \delta)$$

$$\Delta R(t, \delta)$$

$\Leftrightarrow s^*(t)$ is a deviation of s

that maximizes

$$\Delta R(t, s)$$

How to say that a node is a Leaf?

- Stopping rule: a node is a Leaf if:
- the node contains just one observation
 - all the observations of the node have the same value for y .

At the end, we have a maximal tree!

Q_b : the error, computed on the training sampling, for a maximal tree is O^4 .

Example :

dataset

	x^1	x^2	y
e ₁	1.9	0.2	1.2
e ₂	1.5	-0.2	2
e ₃	0.8	1.2	0.5
e ₄	0.4	2.4	0.3

$$S = \left\{ \begin{array}{l} (x^1 \leq 0.6), (x^1 < 1.1) \\ (x^1 \leq 1.45), (x^1 > 0.6) \\ (x^1 > 1.1), (x^1 > 1.45) \\ (x^2 \leq 0), (x^2 \leq 0.7) \\ (x^2 \leq 1.8), (x^2 > 0) \\ (x^2 > 0.7), (x^2 > 1.8) \end{array} \right\}$$

$$\Delta R(t_1, S_1)$$

$$\Delta R(t_1, S_2)$$

$$\Delta R(t_1, S_3)$$

$$\Delta R(t_1, S_{12})$$

$$(x' \leq 1.45)$$

The best division

is

maximal value

$$S = \{ S_1 : (x' \leq 0.6), S_2 : (x' < 1.1) \}$$

$$S_3 : (x' \leq 1.45), S_4 : (x' > 0.6)$$

$$S_5 : (x' > 1.1), S_6 : (x' > 1.45)$$

$$S_7 : (x^2 \leq 0), S_8 : (x^2 \leq 0.7)$$

$$S_9 : (x^2 \leq 1.8), S_{10} : (x^2 > 0)$$

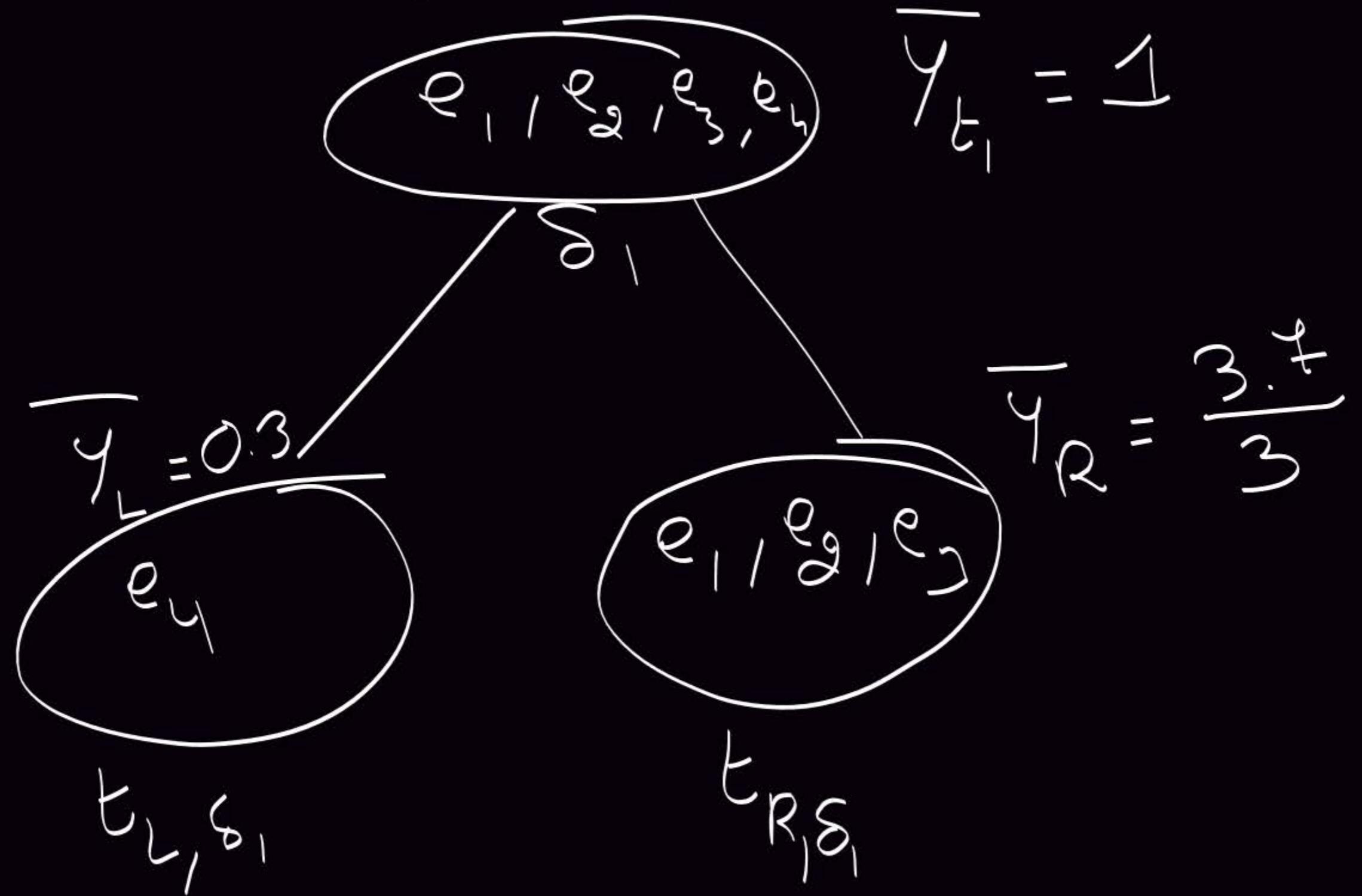
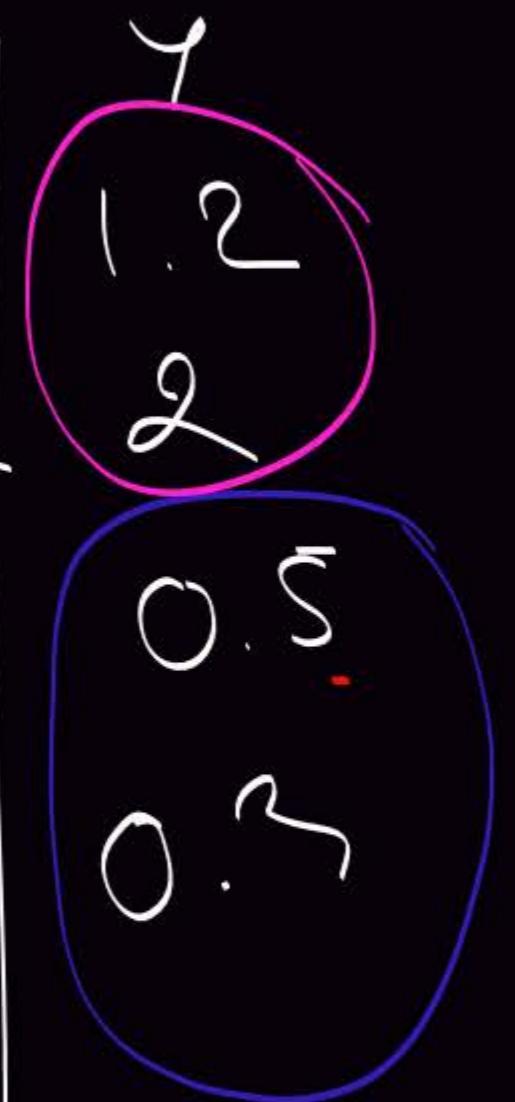
$$S_{11} : (x^2 > 0.7), S_{12} : (x^2 > 1.8)$$

Example :

$$\delta_1 = (X^1 \leq 0.6)$$

dataset

	X^1	X^2
e_1	1.4	0.2
e_2	1.5	-0.2
e_3	0.8	1.2
e_4	0.4	2.4



Example :

	x^1	x^2	y
e ₁	1.4	0.2	1.2
e ₂	1.5	-0.2	2
e ₃	0.8	1.2	0.5
e ₄	0.4	2.4	0.3

$$S_1 = (X^1 \leq 0.6) \quad \bar{y}_{t_1} = 1$$

$$(e_1) \frac{\bar{y}_L}{\bar{y}_R} = 0.3$$

$$\frac{\bar{y}_L}{\bar{y}_R} = \frac{3.4}{3}$$

$$B(t_1) = \frac{1}{4} \left((1 - 1.2)^2 + (1 - 2)^2 + (1 - 0.5)^2 + (1 - 0.3)^2 \right)$$

$$B(t_L) = \frac{1}{4} \left((0.3 - 0.3)^2 \right) = 0$$

$$B(t_R) = \frac{1}{4} \left(\left(1.2 - \frac{3.4}{3} \right)^2 + \left(2 - \frac{3.4}{3} \right)^2 + \left(0.5 - \frac{3.4}{3} \right)^2 \right)$$

what happens in the classification framework?

Let $\{1, \dots, I\}$ denotes the set of modalities
of γ .

$$\forall i \in \{1, \dots, I\}, n_i = \#\{j \in \{1, \dots, n\} \text{ s.t. } \gamma_j = i\}$$
$$\bar{\pi}_i = \frac{n_i}{n} \quad \begin{matrix} \text{prior probability} \\ \text{of class } i \end{matrix}$$

let t be a node of a tree \overline{T}

$$n_t = \#\left\{ i \in \{1, \dots, n\} \text{ s.t. } x_i \in t \right\}$$

$$i \in \{1, \overline{1}\}, n_{i,t} = \#\left\{ j \in \{1, \dots, n\} \text{ s.t. } x_j \in t \text{ and } y_j = i \right\}$$

we estimate:

. the probability that an observation
is in t and of class i by:

$$P(t, i) = \pi_i \times \frac{n_{i,t}}{n_i}$$

. the probability that an observation is in t by:

$$P(t) = \sum_{i=1}^I P(t, i)$$

the posterior probability of class i
in node t by:

$$P(i|t) = \frac{P(i,t)}{P(t)}$$

definition 1:

We say that a function h is
an heterogeneity function on

$$\{(P_1, P_I) \in \mathbb{R}^I \text{ s.t. } \forall i \in \{1, \dots, I\}, P_i \geq 0 \text{ and } \sum P_i = 1\}$$
 if:

- h is symmetric w.r.t (P_1, P_I)
- h is maximal at the point $(\frac{1}{I}, \dots, \frac{1}{I})$
- h is minimal at the points $(1, 0, \dots, 0)$, $(0, 1, 0, \dots, 0)$, ..., $(0, \dots, 0, 1)$

definition 2

let t be a node .

let h be an heterogeneity function .

let define the imp. function by :

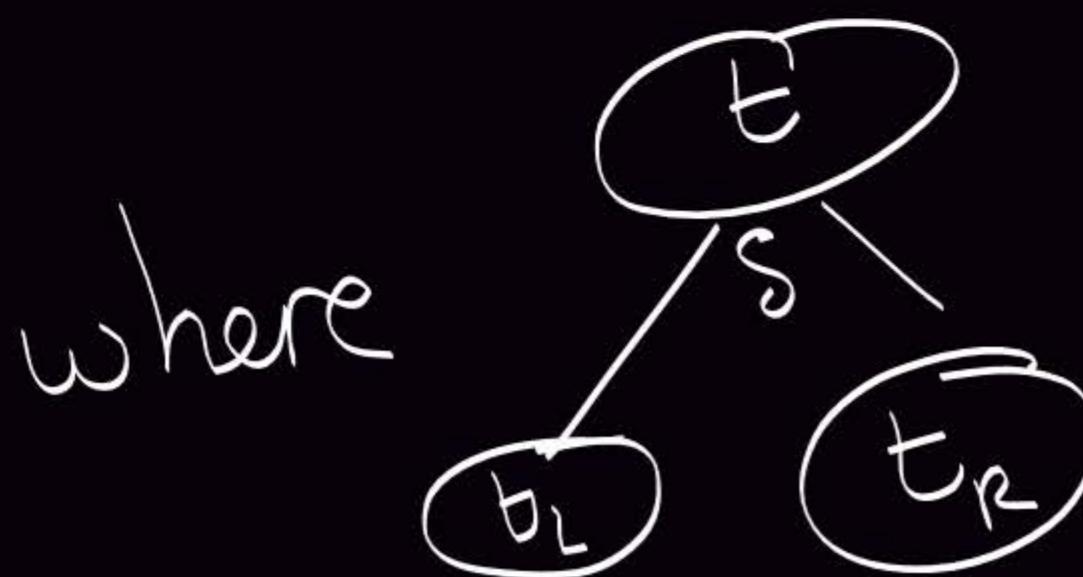
$$\text{imp}(t) = h\left(P(1|t), P(2|t), \dots, P(I|t)\right)$$

Proposition:

If the heterogeneity function is

concave, then we have

$$\Delta \text{imp}(t, s) = \text{imp}(t) - P_R \text{imp}(t_R) - P_L \text{imp}(t_L) \geq 0$$



where

and

$$P_R = P(t_R)/p(t)$$

The way to construct the model
free in the classification framework
is similar to the regression framework
except two things:
- the criterion $\Delta R(t, s)$ is replaced
by $\Delta \text{imp}(t, s)$

- \bar{Y}_t are replaced by $j(t)$ that
are defined by

$$j(t) = \operatorname{argmax}_{i \in \{1, \dots, I\}} P(i|t)$$

R.R.: When in the computer, we do not define
some prior probabilities, the computer we

the prior probabilities π_i .

In this case, $j(t)$ is also the class which is the most represented in t .

If you have other prior probabilities, be careful this other definition is not always true!

Rk: classification

Learning error:

$$\frac{1}{n} \sum_{i=1}^n \mathbb{I}_{y_i \neq \hat{y}}$$

$\text{!} = \text{predict}(\bar{x}; \text{type} = \text{'class'})$

$$1/\text{nrow}(x) * \text{sum}(\text{!})$$

Exercise

Create a matrix X where :

- the first column contains 1000 observations from $\mathcal{U}(2, 0.1)$
- the second one is associated to $\mathcal{E}(0.2)$
- the third one is associated to a $\mathcal{B}(5, 0.7)$
- the variable y $\mathcal{N}(-2, 0.05)$

- the variable S is associated to $\mathcal{B}(4, 0.2)$
 - 6 to $\mathcal{P}(-1, 0.1)$
 - 7 to $\mathcal{P}(1, 0.1)$
 - 8 to $\mathcal{B}(0.5)$

$$Y = \begin{cases} 2 - 5 * \sqrt{5} + 4 * \sqrt{2} + \varepsilon & \text{if } V8 = 0 \\ -3 + 2 * \sqrt{1} - 4 * \sqrt{3} + \varepsilon & \text{if } V8 = 1 \end{cases}$$

with $\varepsilon \sim \mathcal{P}(0, 1)$

Q1: run the CART algorithm
to explain γ thanks to all
the variables.
look at the variable importance
object.

Q2: create y_1 which is defined by:
 $y_1 = 1$ if $y < Q_1$ (1st quartile of y)
 $y_1 = 2$ if $Q_1 \leq y < Q_2$

$$y_1 = 3 \text{ if } Q_2 \leq y < Q_3$$

$$y_1 = 4 \text{ if } y \geq Q_3$$

Run the CA \overline{BUT} algorithm to explain y_1 thanks to V_1 until V_8 .

Q_3 : Repeat Q_1 and Q_2 but you transform V_3 and V_5 into qualitative variables.

Q_y: Repeat just Q₁ with the
initial variables, but create
S bootstrap samples (with
replacement)