

x^1	x^2	y
1.2	1	a
1.4	0.3	a
0.5	-0.2	b
0.3	0	b

How to perform
the classification

tree?

Remember in regression
we do the computation
for one devision δ of $\Delta R(t_i, \delta)$

For the classification
we have to compute

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

? What is imp

$$\text{imp}(t) = \Pr(P(1|t), \dots, P(I,t))$$

↑
???

In general, the imp fact is the

sum function defined by

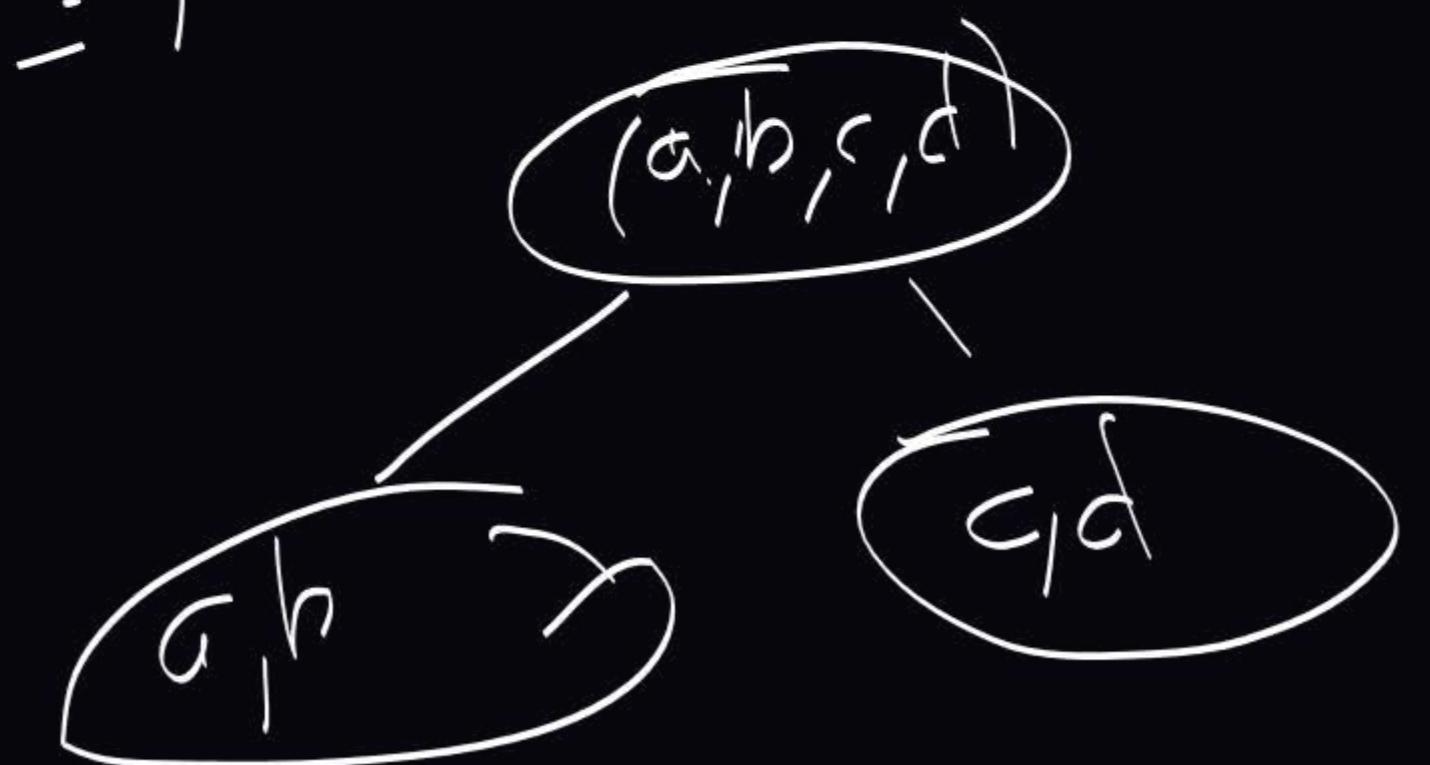
$$\text{imp}(t) = \sum_{i \neq j} P(i|t) \times P(j|t)$$

$$\Rightarrow \Pr(P_1, \dots, P_I) = \sum_{i \neq j} P_i P_j$$

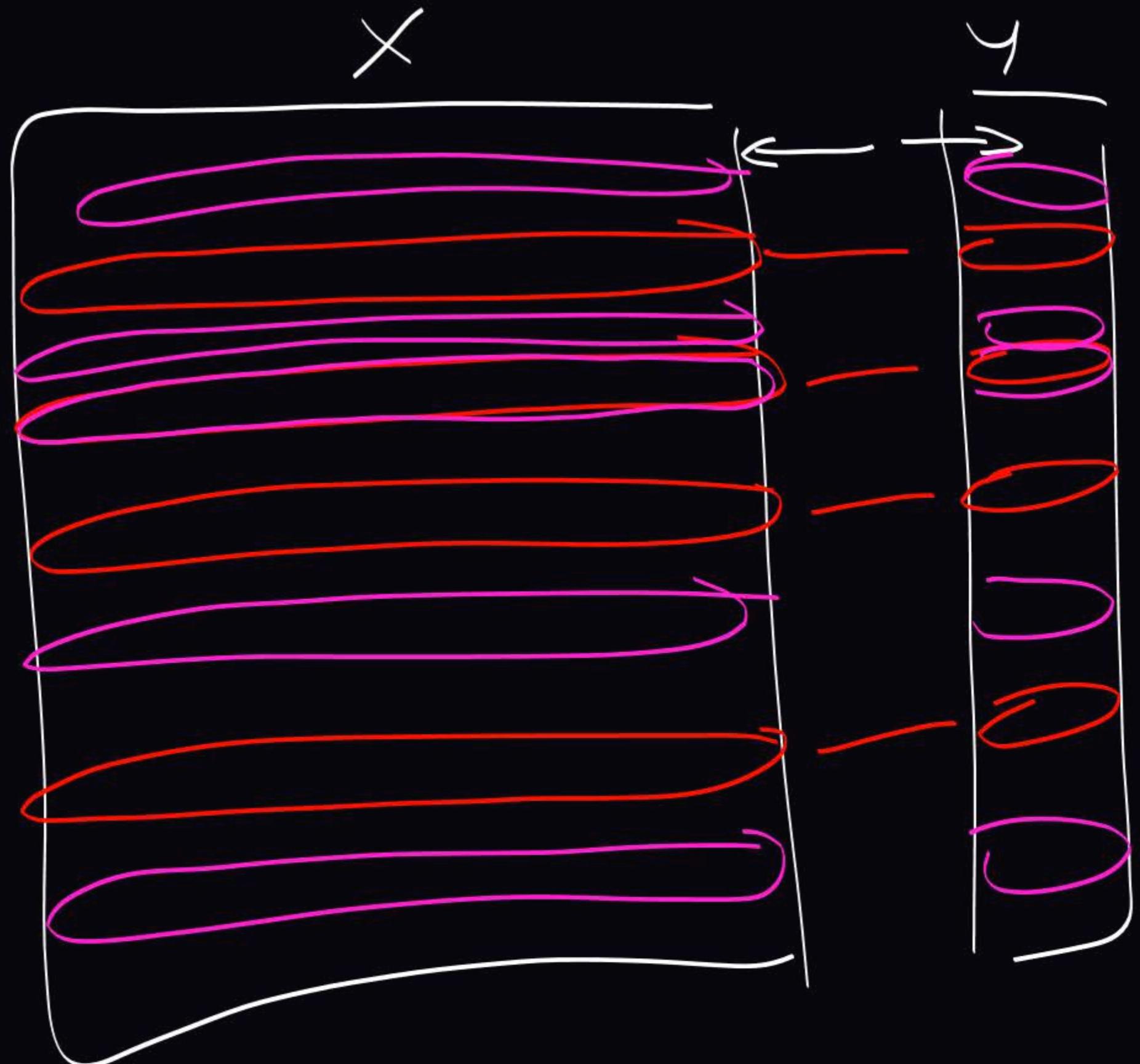
$$\text{imp}(t) = \sum_{i=1}^{|I|} \left(P(i|t) \times \sum_{j \neq i} P(j|t) \right)$$

$$= \sum_{i=1}^{|I|} P(i|t) \times (1 - P(i|t))$$

Shanon

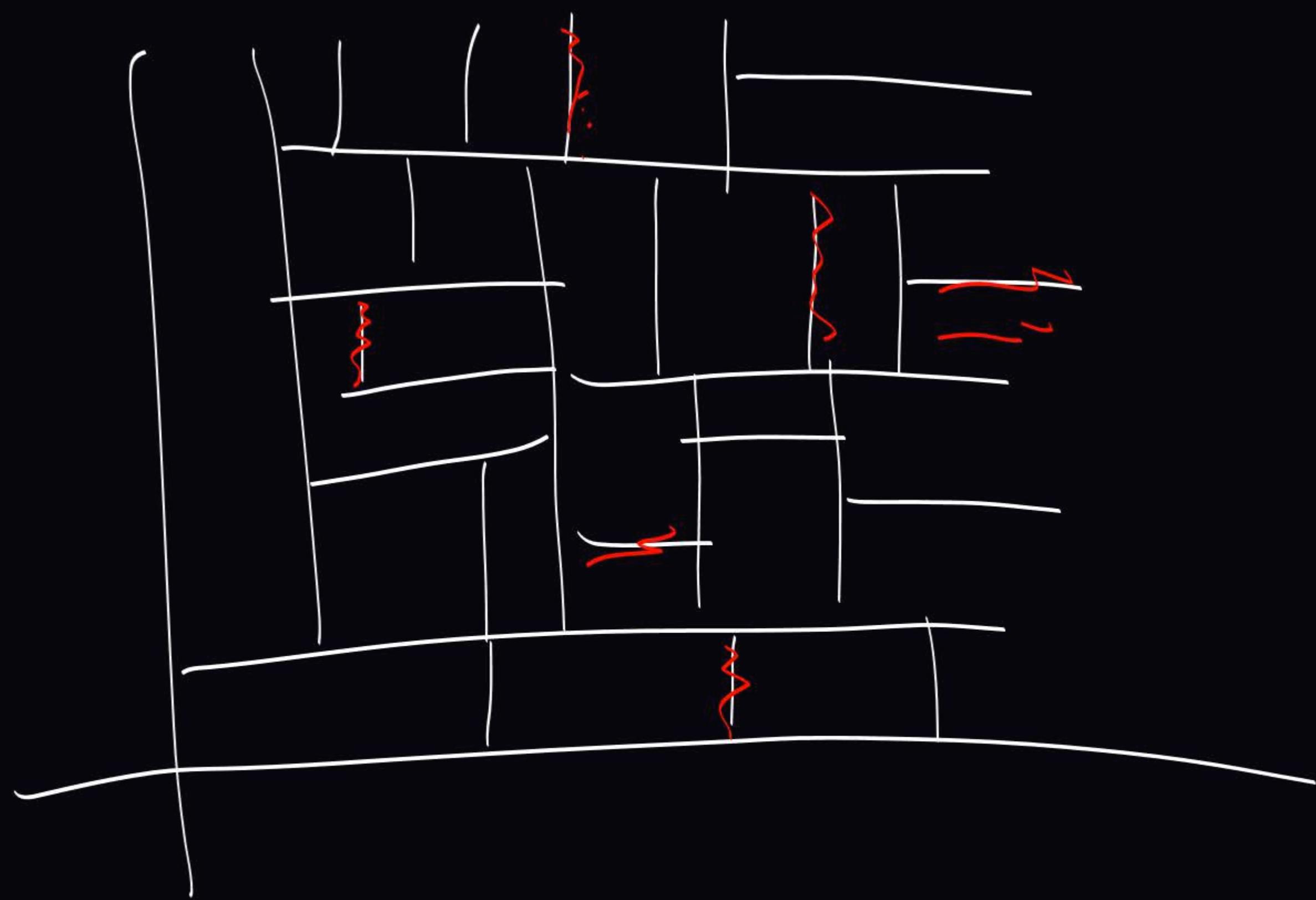


$$\sum_{i=1}^{|I|} -P_i \log P_i$$



At the end of the explained construction of the maximal tree

we get a tree that cannot be used in practice because it is too complex and it overfits the data.



The idea is thanks to the maximal
free, we want to create a smaller

goodness-of-fit

comple-
xity

Free !

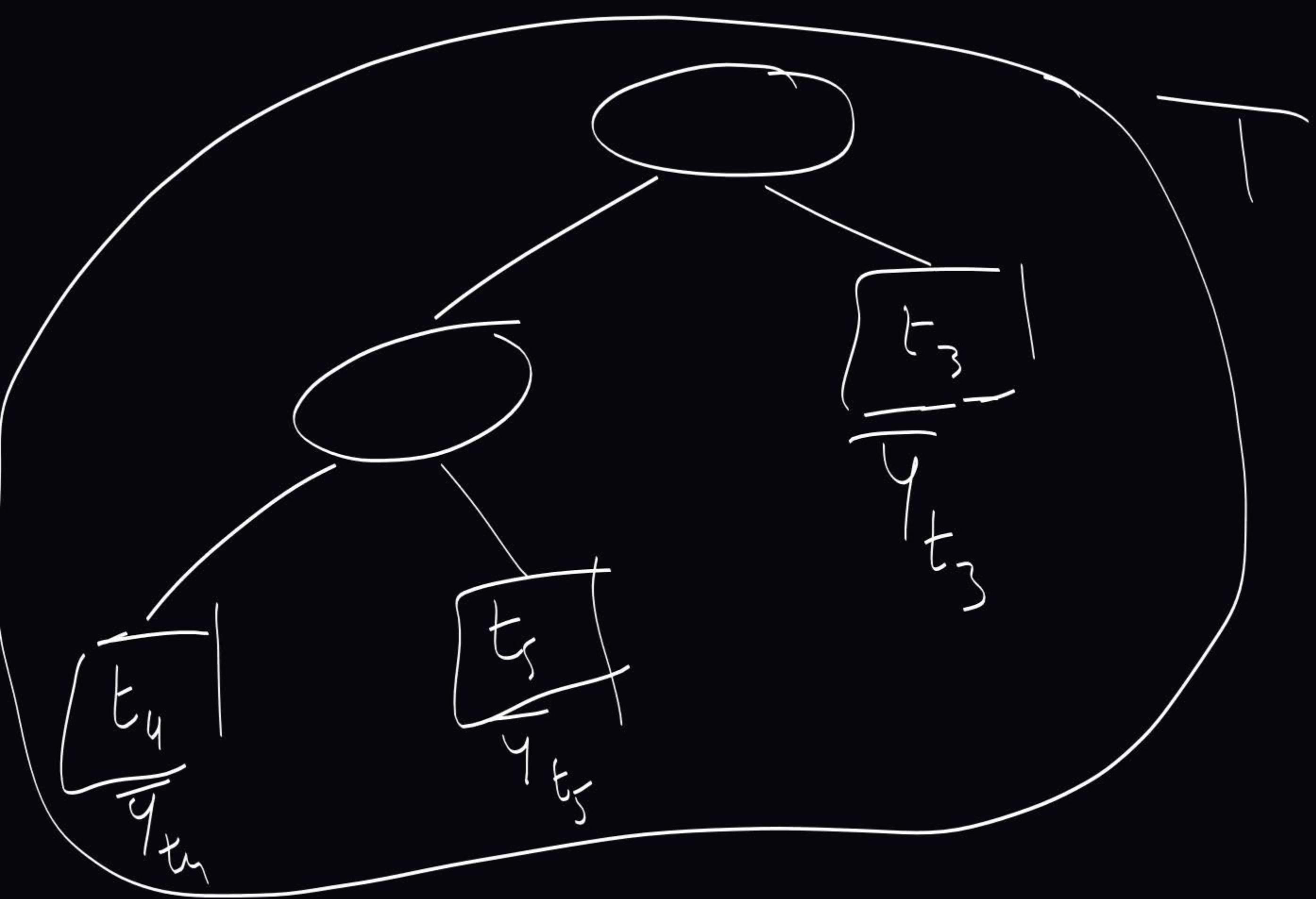
a criterion: $\text{crit}_\alpha(T) = \text{error}(T) + \frac{\alpha}{n} |T|^2$

with $\alpha > 0$, $|T|$ the number of leaves of α
the tree T .

in R, this is
proportional to CP

$$\text{error}(\tau) = \begin{cases} \frac{1}{n} \sum_{i=1}^n (\gamma_i - \hat{\gamma}_{\tau,i})^2 & \text{in regression} \\ = \frac{1}{n} \sum_{t \in T} \sum_{x_i \in t} (\gamma_i - \bar{\gamma}_t)^2 \\ \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{\gamma_i \neq \hat{\gamma}_{\tau,i}} & \text{in classification} \end{cases}$$

where $\hat{\gamma}_{\tau,i}$ is the prediction thanks
to the τ for the individual i



The best tree associated to α is defined by:

$$\bar{T}_\alpha = \begin{cases} \text{argmax } \text{crit}_\alpha(\bar{T}) \\ \bar{T} \not\in T \\ \text{if } \text{crit}_\alpha(\bar{T}_\alpha) = \text{crit}_\alpha(\bar{T}) \\ \text{then } \bar{T}_\alpha \not\in T \end{cases}$$

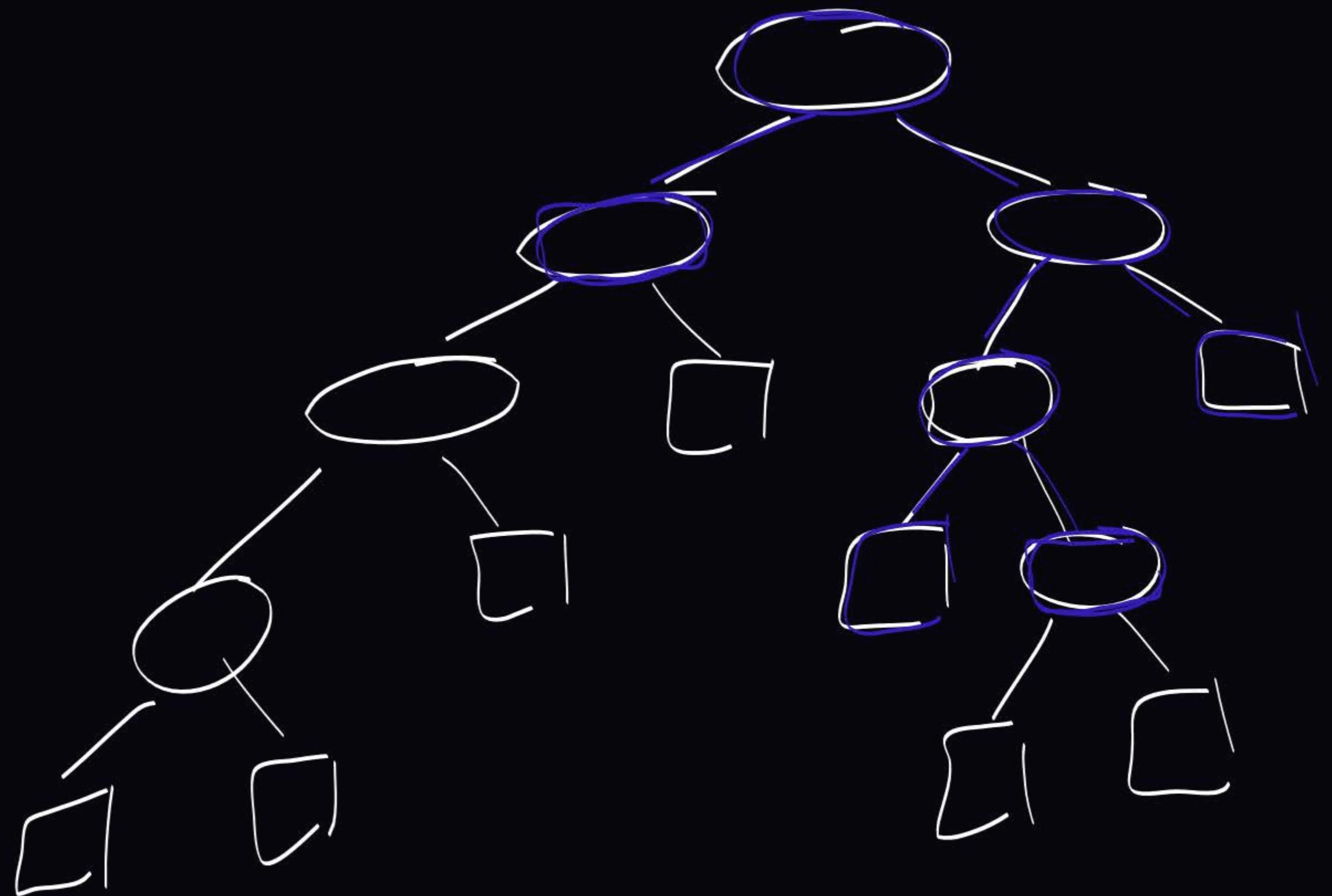
$$\text{crit}_\alpha(\bar{T}) = \min_{T' \subset T} \text{crit}_\alpha(T')$$

and \bar{T}_α is the smallest sub-tree
of T that minimizes crit_α .

where $\overline{T}' \not\subset \overline{T}$ means

\overline{T}' is a sub-tree of \overline{T}

\overline{T}' sub-tree of \overline{T} means that \overline{T}' is obtained from \overline{T} by suppressing nodes of \overline{T} .



Bk: for $\alpha = 0$, T_α is the
maximal tree.

Then we increase α . Thus
the complexity of T_α becomes
smaller.

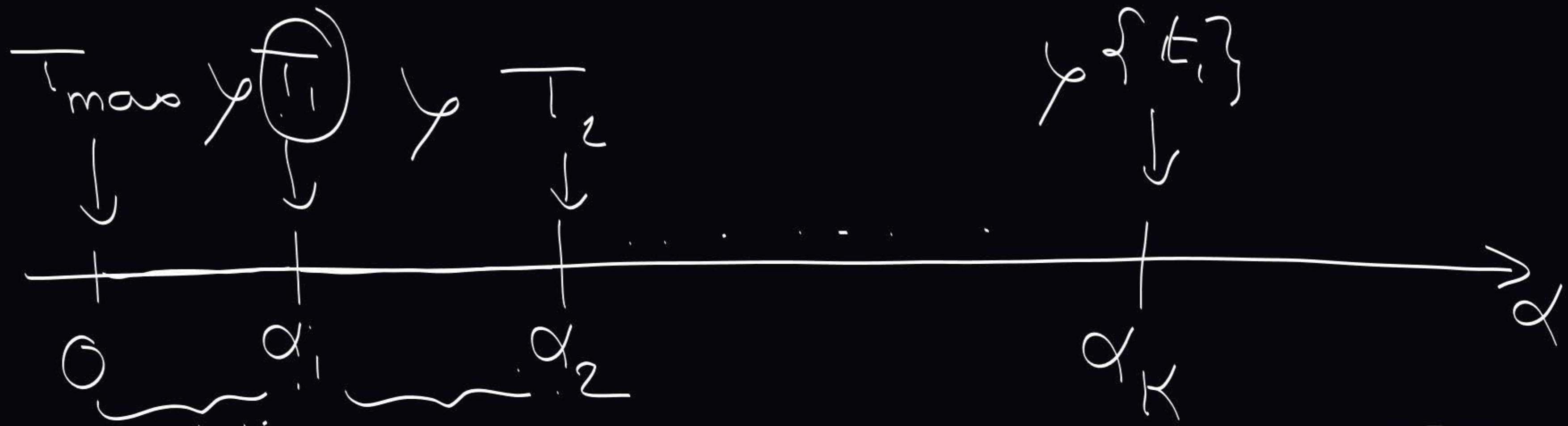
Theorem:

There exists a sequence (α_k) $k \in \{0, \dots, K\}$

which is data-driven such that:

- $\alpha_0 = 0$ $\alpha_1 < \alpha_2 < \dots < \alpha_K$
- $\bar{T}_0 = \bar{T}_{\max}$ $\bar{T}_1, \bar{T}_2, \dots, \bar{T}_K = \{t_i\}$

- $\forall k \in \{0, \dots, K-1\}, \forall a \in [\alpha_k, \alpha_{k+1}[$, $T_a = \bar{T}_k$



The solution
of the minimi-
zation p
is T_{\max}

The solution
of the
minimization
 p is T_1

The solution
of the
minimization
 p is $\{t_i\}$

PR₂:

The sequence of sub-trees that is created is a sequence of nested sub-trees!

• If $|T_{\max}| = m$, your sequence of sub-trees induced by the minimization problem does not necessary contains sub-tree for all the complexity between 1 and m.

⇒ the creation of the sequence
of sub-trees is the second
step of the CART algorithm

The name of this second step is the
pruning step!

P^b: at the end of the second step, we have $(k+1)$ estimators.

How to chose one?

This is the goal of step 3 which is called final selection.

Final selection: 2 way

- 1) Test sample
- 2) Cross-validation

test sample: We assume that we have individual
that have not been used in the
construction of T_{max} , and of T_1, T_2 .

We compute, for an individual j

of the test sample, the prediction

induced by $\overline{T}_{\max}, \overline{T}_1, \dots, \overline{T}_K$

→ For the individual j we have $(K+1)$ predictions

We denote them $\hat{y}_{j,k}$ where $\hat{y}_{j,k}$ is the

prediction for individual j by T_k

We compute the error in test for

all the sub-trees T_0, \dots, T_K by

$$\text{error}_k = \begin{cases} \frac{1}{m} \sum_{j=1}^m (y_j - \hat{y}_{j,k})^2 & \text{regression} \\ \frac{1}{m} \sum_{j=1}^m 1 | y_j \neq \hat{y}_{j,k} & \text{classification} \end{cases}$$

(m : n^b of individuals
in the test sample)

The final selection is the sub-tree

associated to the smallest test

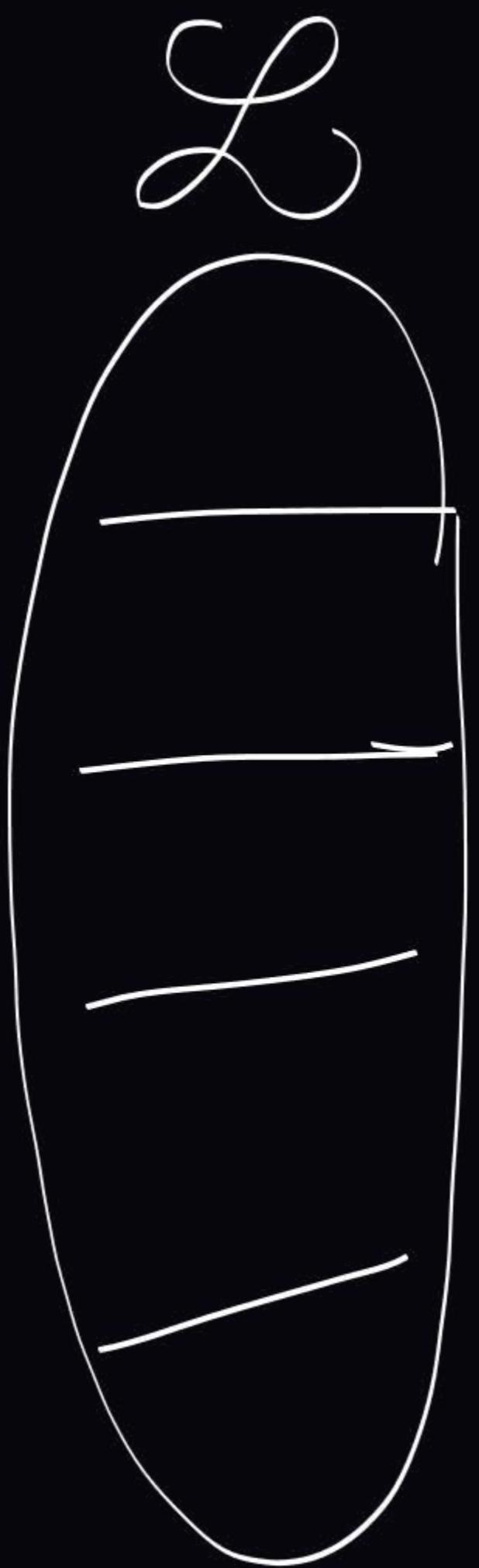
error

Cross - Validation

$$\mathcal{L} = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$$
 training sample

We apply a method A to construct a decision rule \hat{f} to estimate f .

$$\mathcal{L} \xrightarrow{A} \hat{f}$$



a partition
of L

\rightarrow in V elements
created at random.

Each element of the

partition has quite
the same number of observations ($\approx \frac{n}{V}$)

L_1

L_2

\vdots

L_V

Step 1:

$$\text{Consider } \mathcal{L}_2 \cup \dots \cup \mathcal{L}_V = \mathcal{L}_{-1} \quad \frac{\mathbb{R}\mathbb{L} \cdot \#\mathcal{L}_{-1}}{V} \approx \frac{(V-T)}{V} n$$

to construct an estimator $\hat{\varphi}_{-1}$

by application of the method A.

We compute the test error of $\hat{\varphi}_{-1}$ thanks to \mathcal{L}_1
 $\underbrace{\text{err}_{-1}}$

Example:

linear regression

$$X = \begin{pmatrix} 1 & 1.2 & 2.5 \\ 1 & 1.1 & 3.4 \\ 1 & 0.5 & 1.4 \\ 1 & 2.3 & 0.7 \\ 1 & 0.7 & -0.1 \end{pmatrix}$$

$$Y = \begin{pmatrix} 0.5 \\ 1.2 \\ -0.1 \\ 0.3 \\ 1.5 \end{pmatrix}$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

$$\hat{\beta}(T) = T \hat{\beta}$$

$$\begin{aligned}\hat{f}(x_1, x_2) &= \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 \\ &= (1 \ x_1 \ x_2) \times \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix}\end{aligned}$$

$$V = 5$$

$$\mathcal{L}_1 : X_1 = \begin{pmatrix} 1 & 1.2 & 2.4 \end{pmatrix} Y_1 = 0.5$$

$$X_{-1} = \begin{pmatrix} 1 & 1.1 & 3.7 \\ 0.5 & 1.4 & 0.7 \\ 2.3 & 0.7 & -0.1 \\ 0.7 & -0.1 & \end{pmatrix}$$

$$Y_{-1} = \begin{pmatrix} 1.2 \\ 0.1 \\ 0.3 \\ 1.5 \end{pmatrix}$$

$$\text{error}_{-1} = (X_1 \hat{\beta}_{-1} - Y_1)^2$$

$$\hat{\beta}_{-1} = X_{-1}^{-1} Y_{-1}$$

with

$$\hat{\beta}_{-1} = \hat{\beta}_{-1}^T$$

$$\hat{\beta}_{-1}^T(\top) = \top \hat{\beta}_{-1}$$

Step 2 :

We use $\mathcal{L}_{-2} = \bigcup_{\substack{v \in \{1, 2\} \\ v \neq 2}} \mathcal{L}_v$

to construct thanks to the method A
an estimator \hat{f}_{-2} .
We compute the test error (err_{-2}) of \hat{f}_{-2} by
using \mathcal{L}_2 .

We continue by step 3, Step 4, ...
until step V.

At the end, we have

$$\text{err}_{-1}, \text{err}_{-2}, \dots, \text{err}_{-V}$$

The cross validation error of \hat{f} is given by:
$$\text{err}_W = \frac{1}{n} \sum_{i=1}^n n_i \text{err}_{-i}$$
 with $n_i = \# \mathcal{L}_i$

Pk: if $V = n$ this is the
leave-one-out cross validation.

for the application to CART

algorithm

$$\mathcal{L} \xrightarrow{\text{CART}} \max_{\bar{T}_1, \dots, \bar{T}_K} \left(\alpha_0 = 0, \alpha_1, \dots, \alpha_K \right)$$

$$\mathcal{L} \xrightarrow{\text{partition}} \mathcal{L}_1, \mathcal{L}_V$$

Step 1

$$\mathcal{L}_{-1} = \bigcup_{i \in \{2, \dots, N\}} \mathcal{L}_i$$

We use \mathcal{L}_{-1} to construct $\overline{\mathbf{T}}_{\max, -1}, \overline{\mathbf{l}}_{1, -1}, \dots, \overline{\mathbf{T}}_{K(-1), -1}$

We use \mathcal{L}_1 to compute the test error associated

to this sequence.

$$\alpha'_0 = \sqrt{\alpha_0 \alpha_1}, \quad \alpha'_1 = \sqrt{\alpha_1 \alpha_2}, \quad \alpha'_2 = \sqrt{\alpha_2 \alpha_3}, \quad \dots, \quad \alpha'_K = \sqrt{\alpha_K (\alpha_{K+1})}$$

for $i \in \{0, \dots, k\}$
err _{α_{i-1}} the test error using \mathcal{L}_i

$$g^T \overline{\alpha}_{i-1}' = \overline{\alpha}_{i-1}^T \overline{\alpha}_{i-1}$$
$$\overline{\alpha}_{0,-1} \quad \overline{\alpha}_{1,-1} \quad \overline{\alpha}_i \quad \overline{\alpha}_{2,-1} \quad \dots \quad \overline{\alpha}_{(k-1),-1}$$
$$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \dots \quad \downarrow$$
$$\alpha_{0,-1} \quad \alpha_{1,-1} \quad \alpha_i' \quad \alpha_{2,-1} \quad \dots \quad \alpha_{(k-1),-1}$$

Example :

$$\alpha_0' = 0 \quad \alpha_1' = 0.2 \quad \alpha_2' = 0.3$$

Step 1 $\bar{t}_{0,-1}$

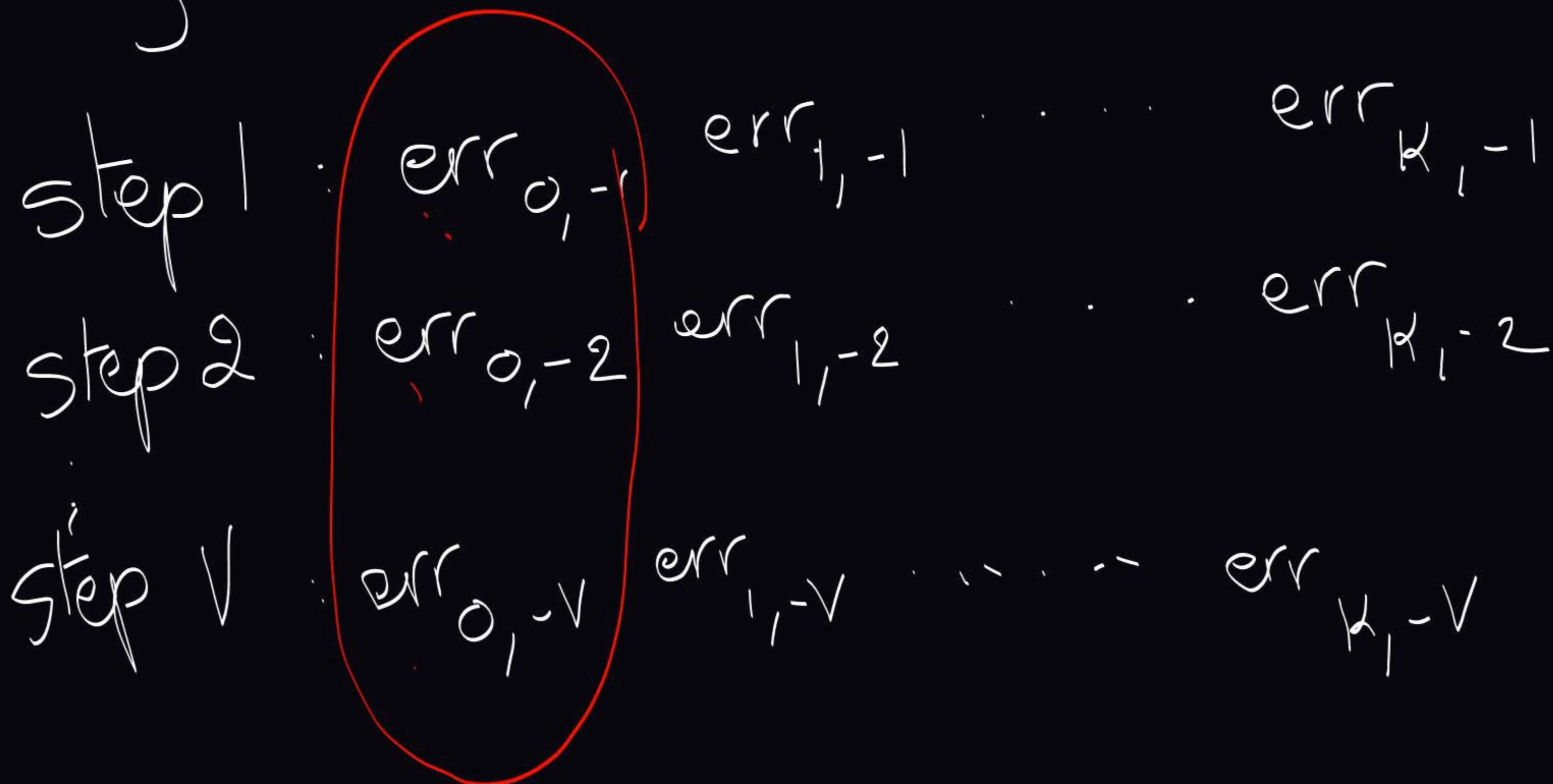
$$\alpha_{0,-1} = 0 \quad \alpha_{1,-1}' = 0.05 \quad \alpha_{2,-1}' = 0.1, \quad \alpha_{3,-1}' = 0.4$$

$\text{err}_{1,-1}$ = test error of $\bar{t}_{1,-1}$

$\text{err}_{2,-1} = \underline{\hspace{10em}}$

We repeat the same idea

for step 2 until step V.



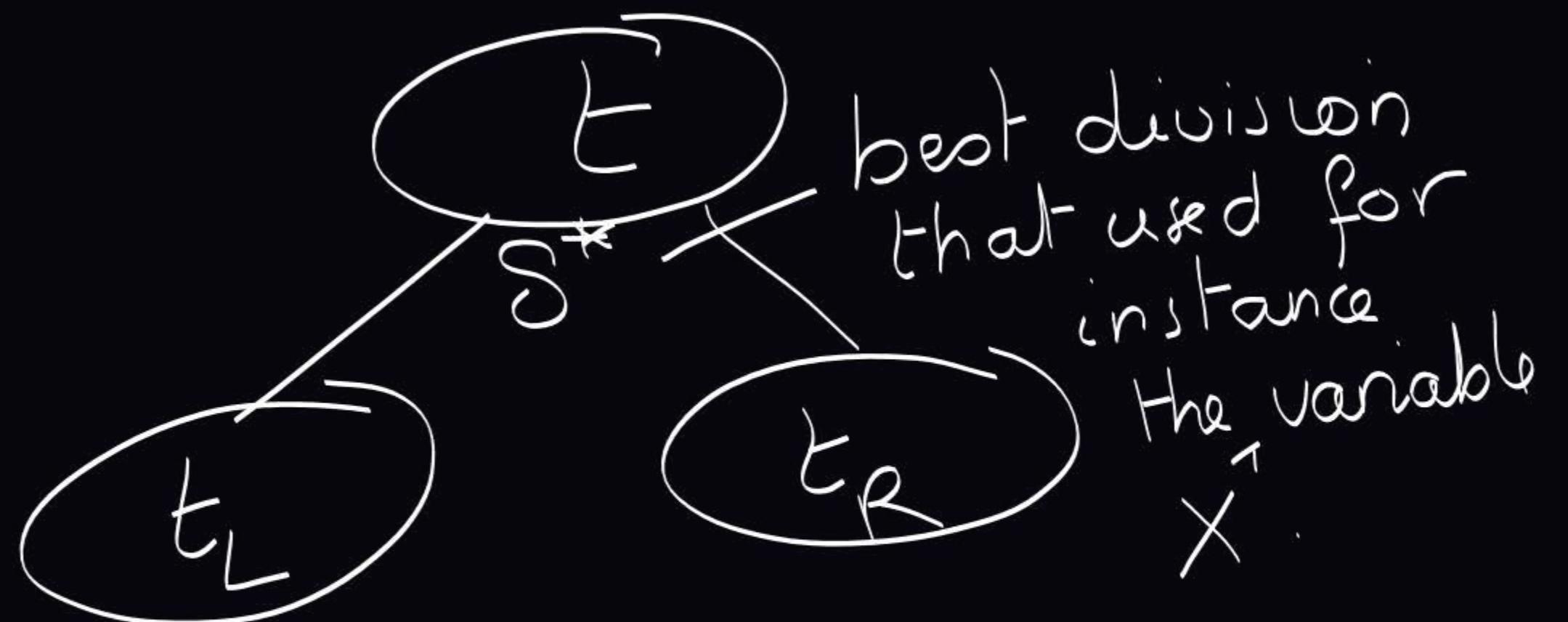
cross validation error of \bar{T}_i is given by:

$$e_{cv}(\bar{T}_i) = \sum_{j=1}^J \text{err}_{i-j} \times \frac{n_j}{n}$$

with $n_j = \# \mathcal{L}_j$

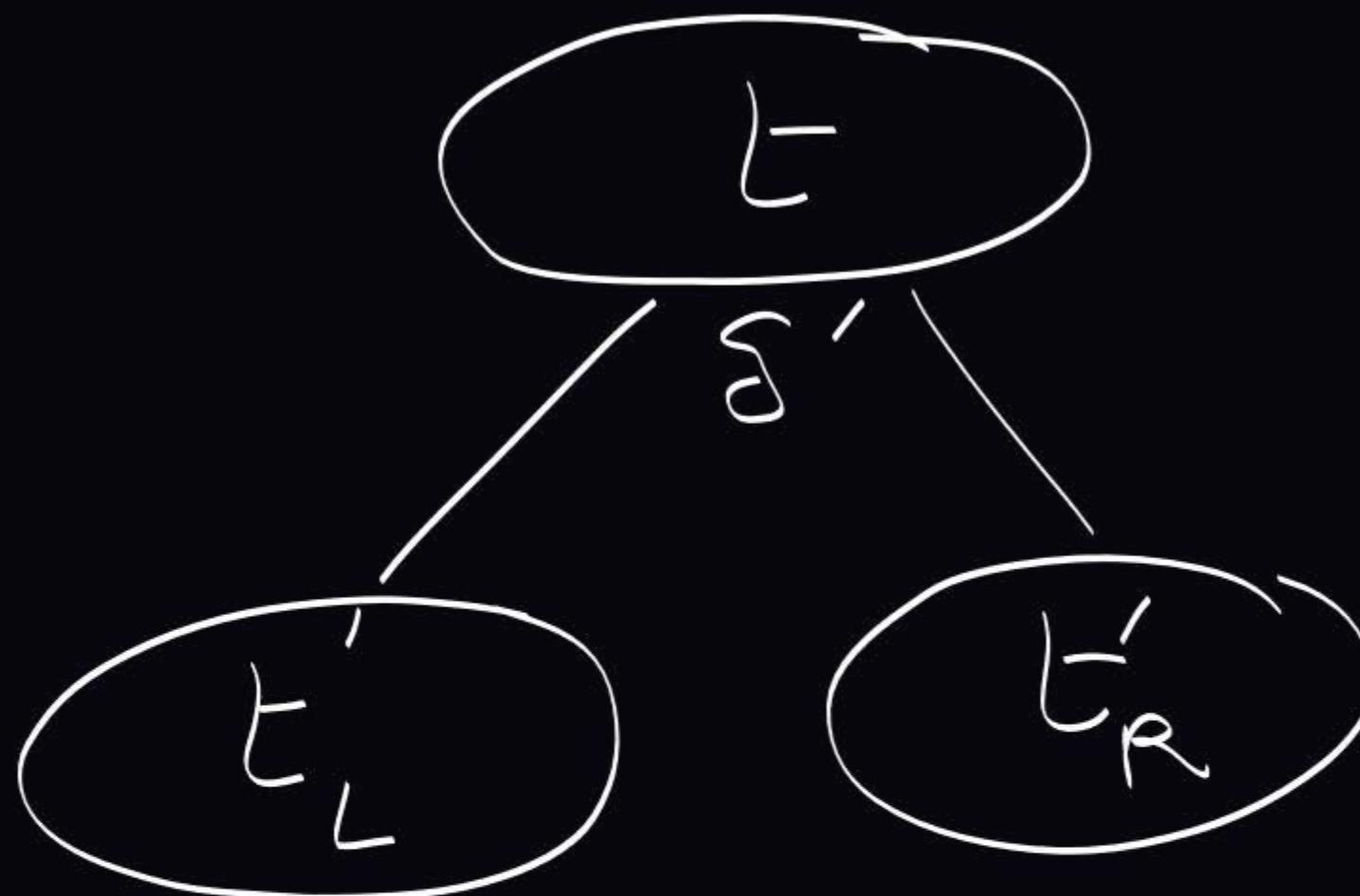
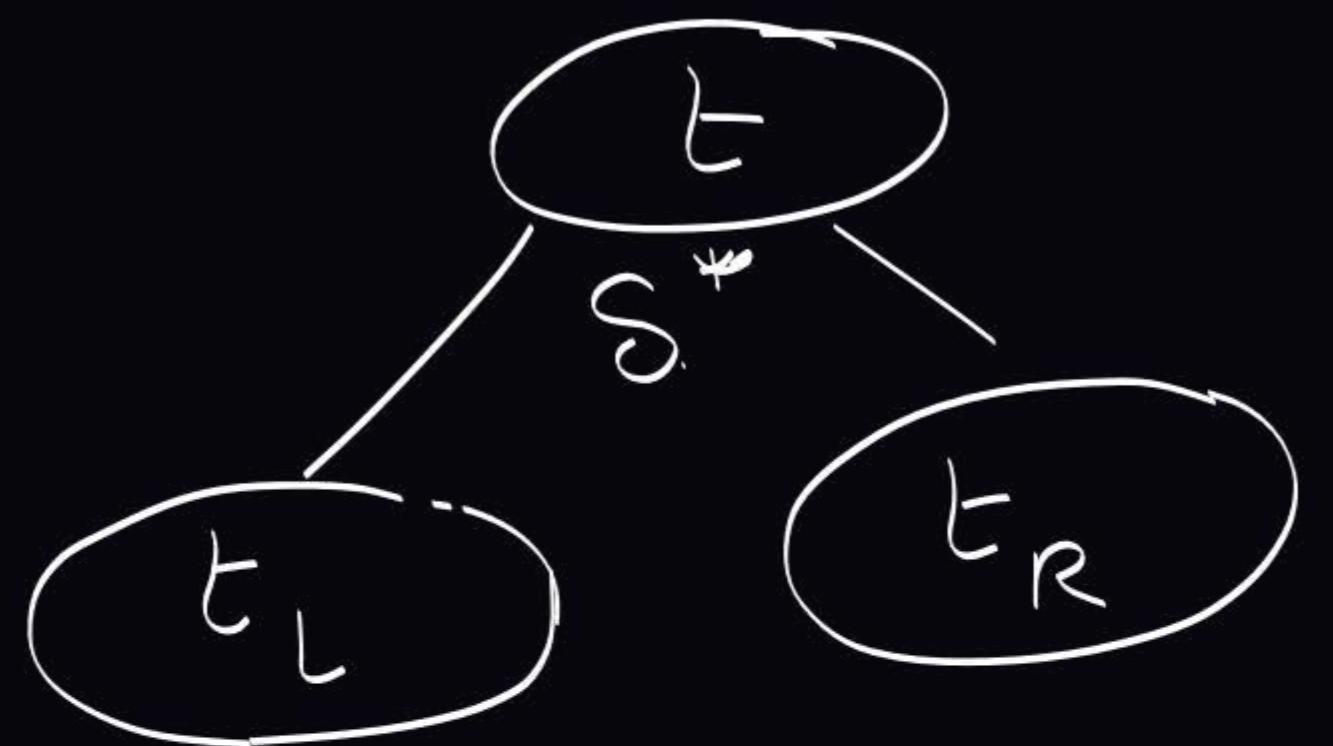
\Rightarrow To select the best sub-tree, we keep the one associated to the minimal value of the cross-validation error!

Surrogate splits



We are going to define the best division or X^2 or X^P or X^- or X^+ which will be "similar" to s^* .

“Similar”



P_{LL}' = the proportion of observations of t that
are in t_L and in t'_L

P_{RR}' also

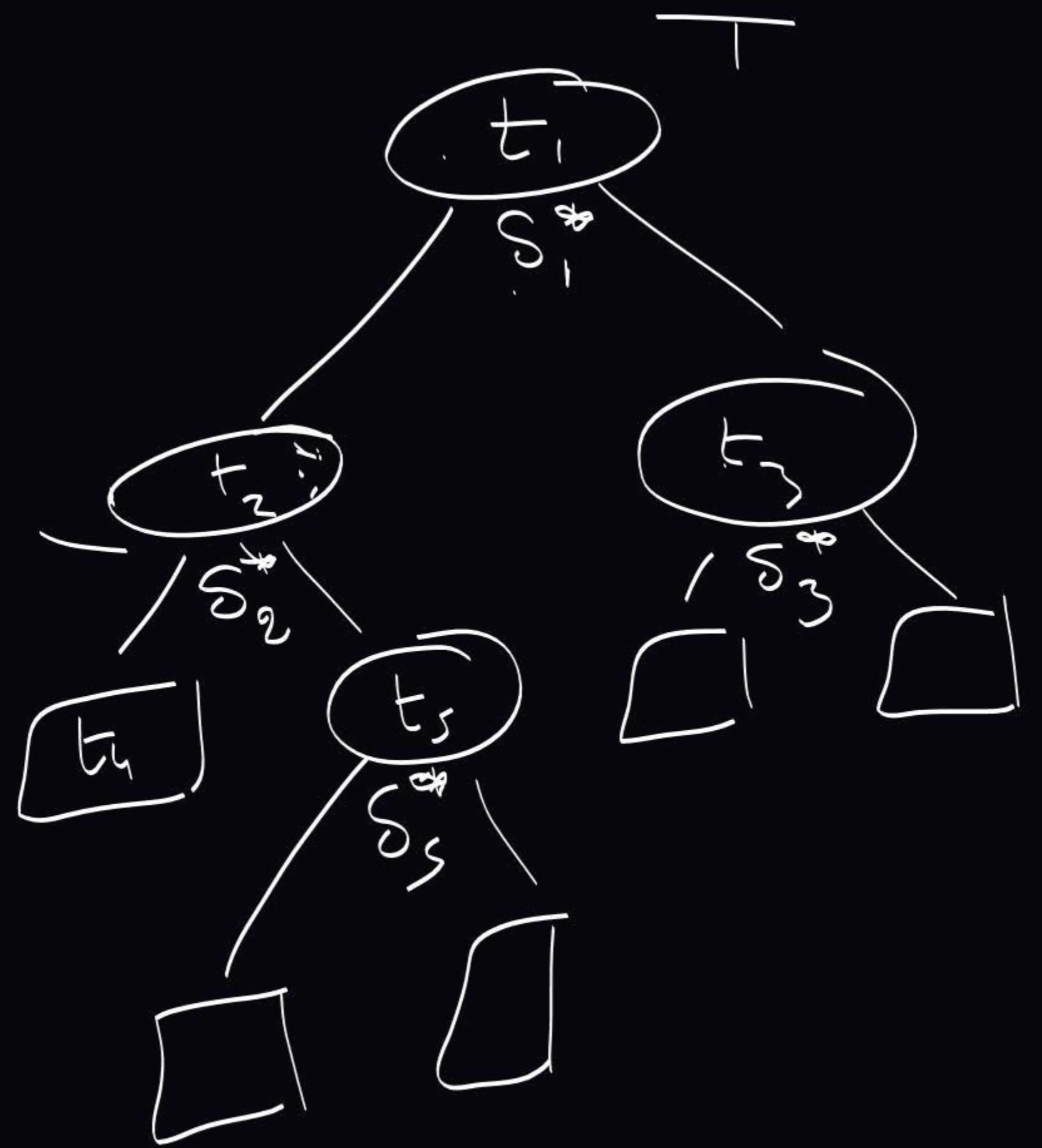
We would like to maximize

$$P_{LL'} + P_{RR'}$$

Variable importance

$$VI(X^i) = \begin{cases} \sum_{t \in \bar{T} - \tilde{T}} \Delta_{imp}(t) & \text{classification} \\ \sum_{t \in \tilde{T}} \Delta_R(t) & \text{regression} \end{cases}$$

all the nodes of \bar{T}
that are not a leaf



Assume that we
want to compute

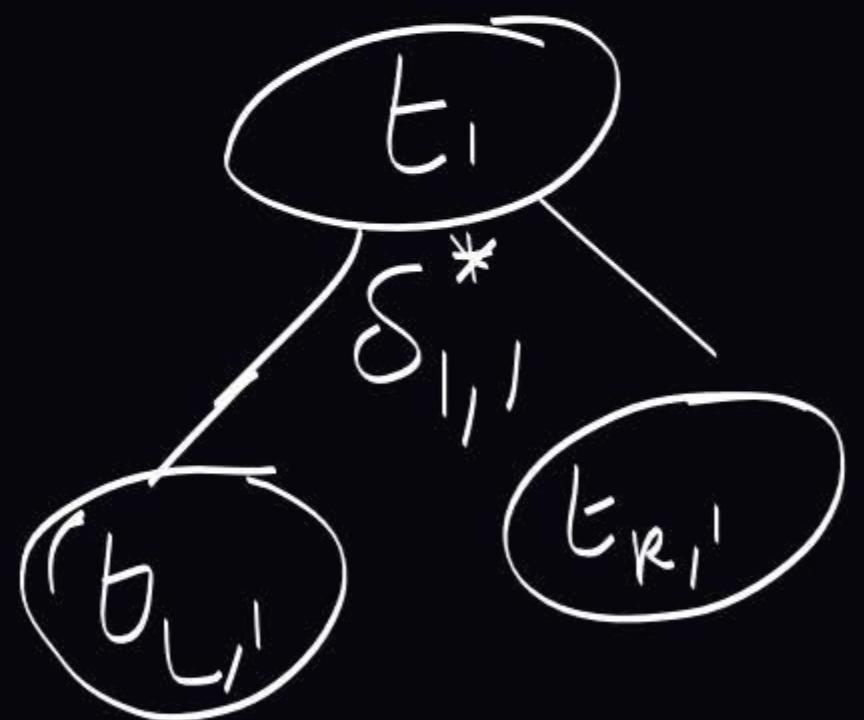
The variable importance

δ_X

If s^* likes X' we are
just compute $\Delta \text{imp}(t_1)$
 $= \Delta \text{imp}(t_1)$

If s_i^* does not use x'

We consider $s_{i,i}^*$ which is
the best surrogate split of t_i



that uses x'

$$\Delta' \text{imp}(t_i) = \text{imp}(t_i) - P_R \text{imp}(t_{R,i}) - P_L \text{imp}(t_{L,i})$$

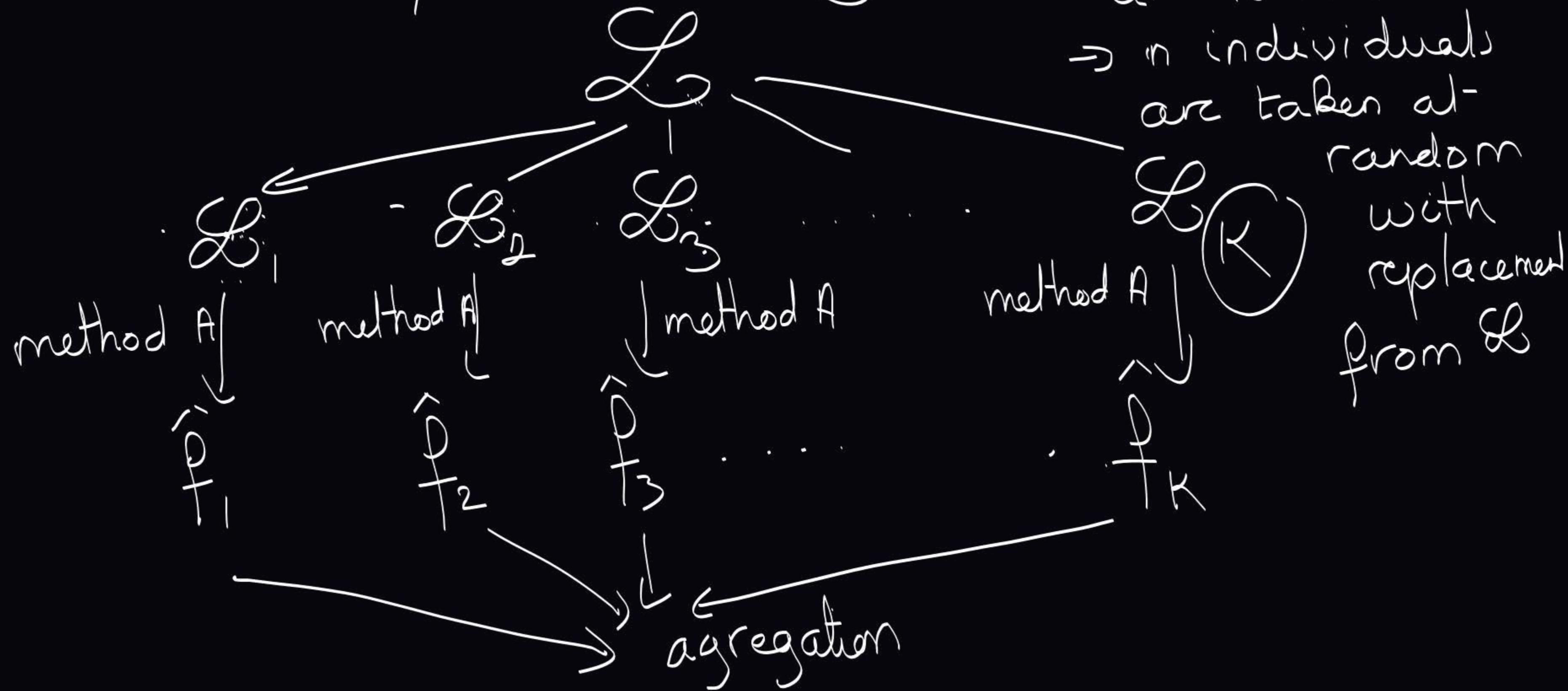
oftens:

$$VI(x^i) = 100 \times \frac{VI(x^i) - \min_i VI(x^i)}{\max_i (x^i)}$$

P^b: A CART tree is unstable.

So if we change just a little
the data, we may have something
different
 \rightarrow stabilization

1st way: bagging



L_1, \dots, L_K have
all n individuals
⇒ n individuals
are taken at-
random
with
replacement
from L

aggregation of $\hat{f}_1, \dots, \hat{f}_K$:

- in regression:

$$\hat{P}(x) = \frac{1}{K} \sum_{i=1}^K \hat{f}_i(x)$$

- in classification:

The response variable takes values in $\{1, \dots, I\}$.

Let x be an observation

$$\forall i \in \{1, \dots, T\}, n_i(x) = \#\left\{j \in \{1, \dots, K\} \text{ s.t. } \hat{f}_j(x) = i\right\}$$

$$\hat{f}(x) = \operatorname{argmax}_{i \in \{1, \dots, T\}} n_i(x)$$

CART

\mathcal{L} : all your data

1) $\downarrow \bar{T}_{\max}$

2) T_0, T_1, \dots, T_K

3) \bar{T}_i by cross validation

