

Lecture FYS-  
STK3155/4155,  
November 11, 2024

# FYS-STK3155/4155

november 11

---

Why Decision trees ?

- simple
- building block in Ensemble methods
  - Bagging
  - Random Forest
  - Boosting
  - Gradient Boosting

# A simple tree

Binary selection

Root node

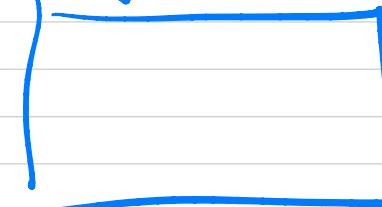
TRUE

Essential feature

FALSE



TRUE  
FALSE



Interior node

TRUE  
FALSE



LEAF NODE

= Pure case  
only one feature

what is an impurity?  
when all elements of a  
branch (= leaf) contain  
only one class, we say we  
have a Pure case

How can we define a  
measure for when a  
class is pure?

Two ways to measure:  
gini factor      g  
entropy            s

$$g, s \in [0, 1]$$

$g = 0$ , pure case

$s = 0$  — 2

Example: binary case

True has probability =  $P$

$$P \in [0, 1]$$

False :  $1 - P$

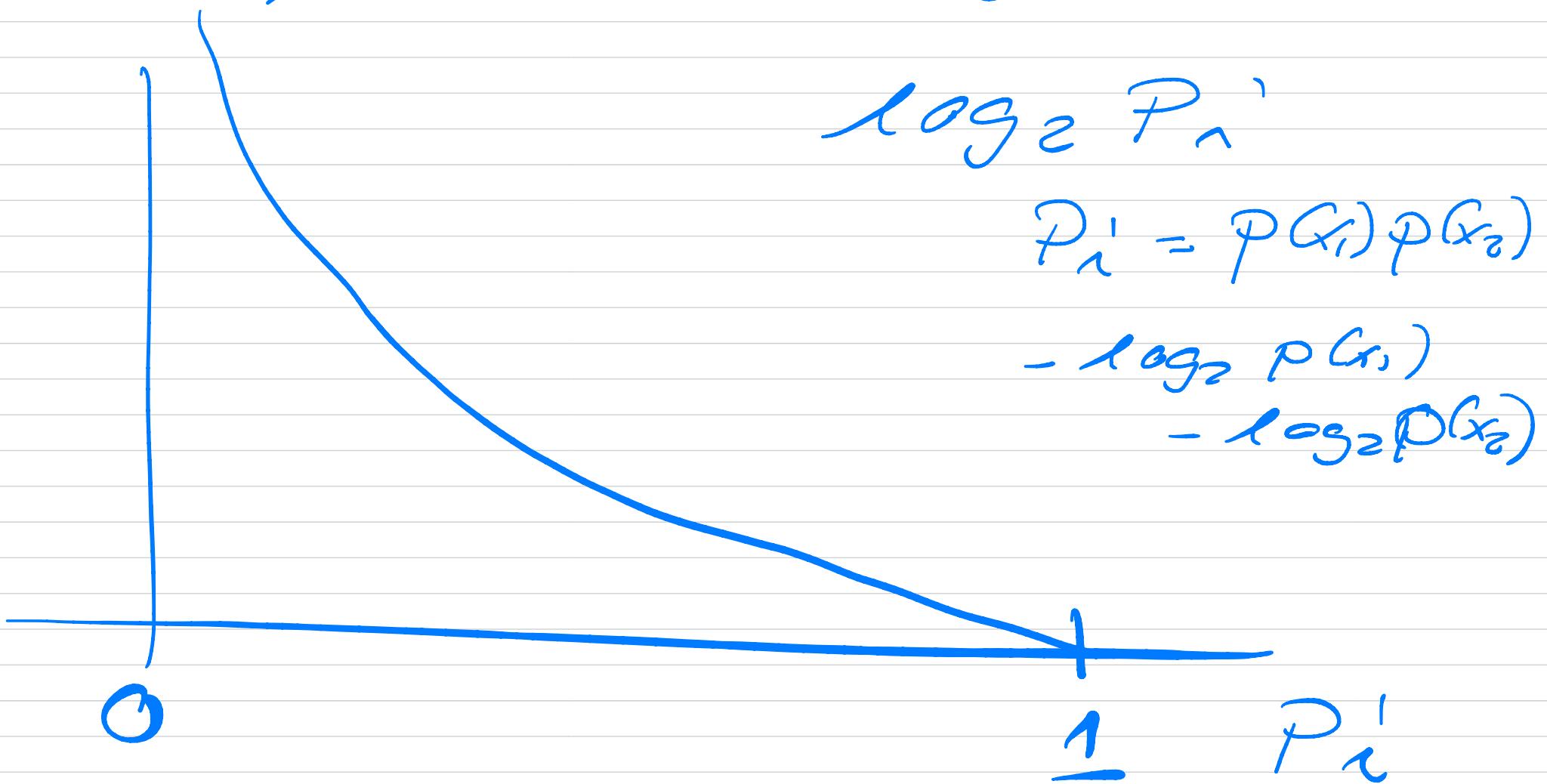
$$\sum_{i=1}^2 P_i = (1-P) + P = 1$$

Criterions for  $I$  or  $S$   
Measure of information

- additive
- continuous
- non-negative
- high for unlikely events

entropy :

$$S_i = - P_i \log_2 P_i$$



$$S = - \sum_{i=1}^n p_i \log_2 p_i$$

$$(0 \log_2 0 = 0)$$

$$-\frac{p \log_2 p}{\text{True}} - \frac{(1-p) \log_2 (1-p)}{\text{False}}$$

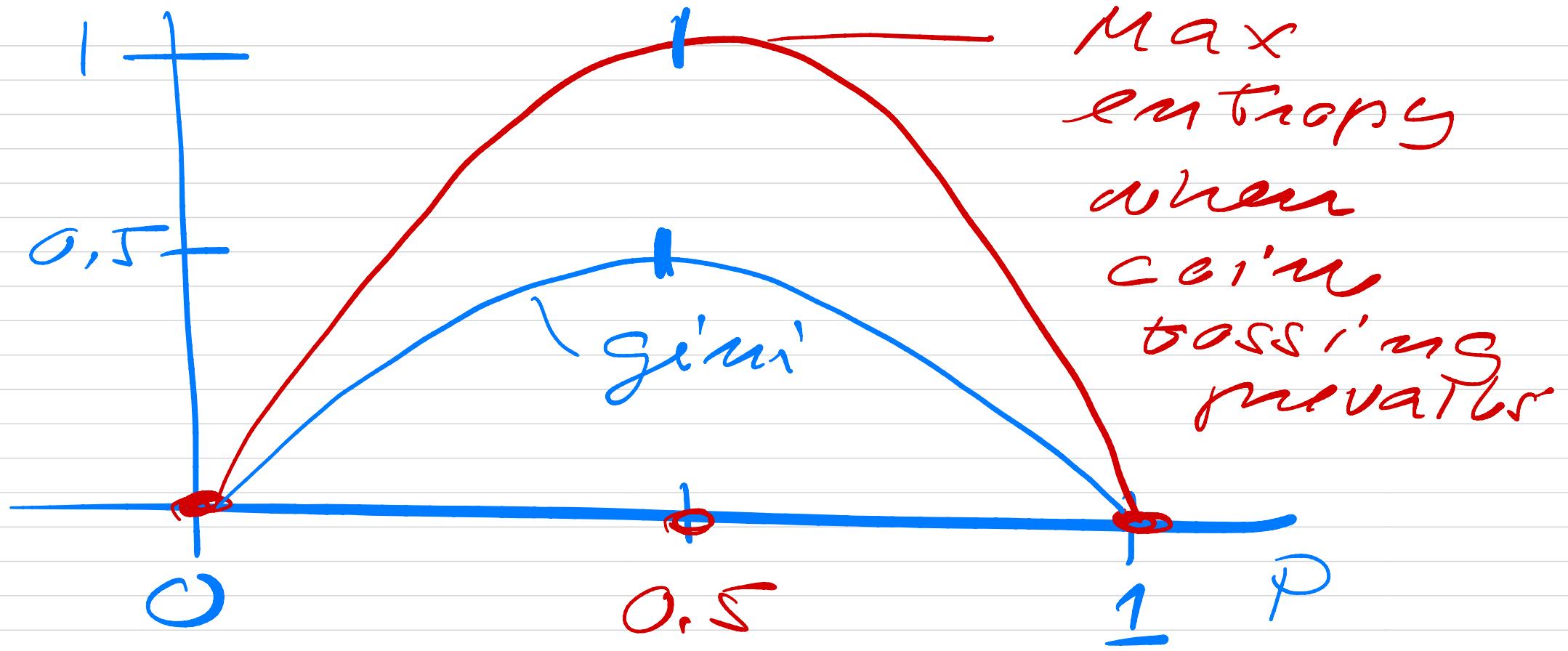
$$S = 0 \quad \text{if } p = 0 \quad \text{FALSE}$$

(pure)

$$S = 0 \quad \text{if } p = 1 \quad \text{TRUE}$$

(pure)

$$p = 0.5 \Rightarrow S = 1$$



$P = 0.5$ , max entropy is  
the probability of an  
imperfect classification  
Both cases are equally  
probable.

Gini's factor

$$g = \sum_{i=1}^M P_i(1-P_i)$$
$$= \underbrace{\sum_i P_i}_{=1} - \sum_i P_i^2 =$$

$$1 - \sum_{n=1}^N p_n^2$$

$$g = [\bar{0}, 1]$$

For our simple example

$$g = 1 - \sum_{n=1}^N p_n^2$$

$$= 1 - p^2 - (1-p)^2 = 2p(1-p)$$

$$p=0 \Rightarrow g=0$$

$$p=1 \Rightarrow g=0$$

$$p=0.5 \Rightarrow g=0.5$$

# Classification Example

- Features/Attributes
  - Grade trend (above/below)
  - # of hours studied  
(high/low)
  - # of hours slept  
(high/low)
- output  
grade above or  
below average

Track	Study	Sleep	Grade
Above	low	H	A
Below	high	L	B
A	L	H	A
A	H	H	A
B	L	H	B
A	L	L	B
B	H	H	B
B	L	H	B
A	L	L	B
A	H	H	A

Gini factor in order to make a tree

$$P(\text{Trend} = A) = 6/10$$

$$P(\text{Trend} = B) = 4/10$$

if Trend = above & grade

= below)

$$P = 2/6 = 1/3$$

if Trend = above & grade  
Above)

$$P = 4/6$$

$$\text{Simi} = 1 - \left( \left( \frac{4}{5} \right)^2 + \left( \frac{2}{8} \right)^2 \right) = 0.45$$

if (Trend = below & grade above)  
 $P = 0$

if (Trend = below & grade below)  
 $P = 4/4 = 1$

$$\text{Simi index } 1 - (1^2 + 0^2) = 0$$

weighted sum for Trend:

$$6/10 \times 0.45 + 4/10 \cdot 0 = \boxed{0.27}$$

Simi factor for hours slept

$$g = \boxed{0.47} \quad (\text{weighted})$$

weighted factor for hours  
stressed

$$g = \boxed{0.34}$$

Roct

Grade  
Trend  
 $g = 0,27$

TREND  
Above  
(6)

?

LHS

TREND  
Below  
(4)

?

RHS

LHS

		Sleep study G				
		L	H	A		
A		L	H	A		
A	A	H	H	A		
A	A	L	L	A		
A	A	L	L	B		
A	A	H	H	B		
				A		

$$P(\text{Sleep} = H) = 2/6$$

$$P(\text{SLEEP} = L) = 4/6$$

if (H & grade  
above)

$$P = 1$$

if (H & grade  
below)

$$P = 0$$

:

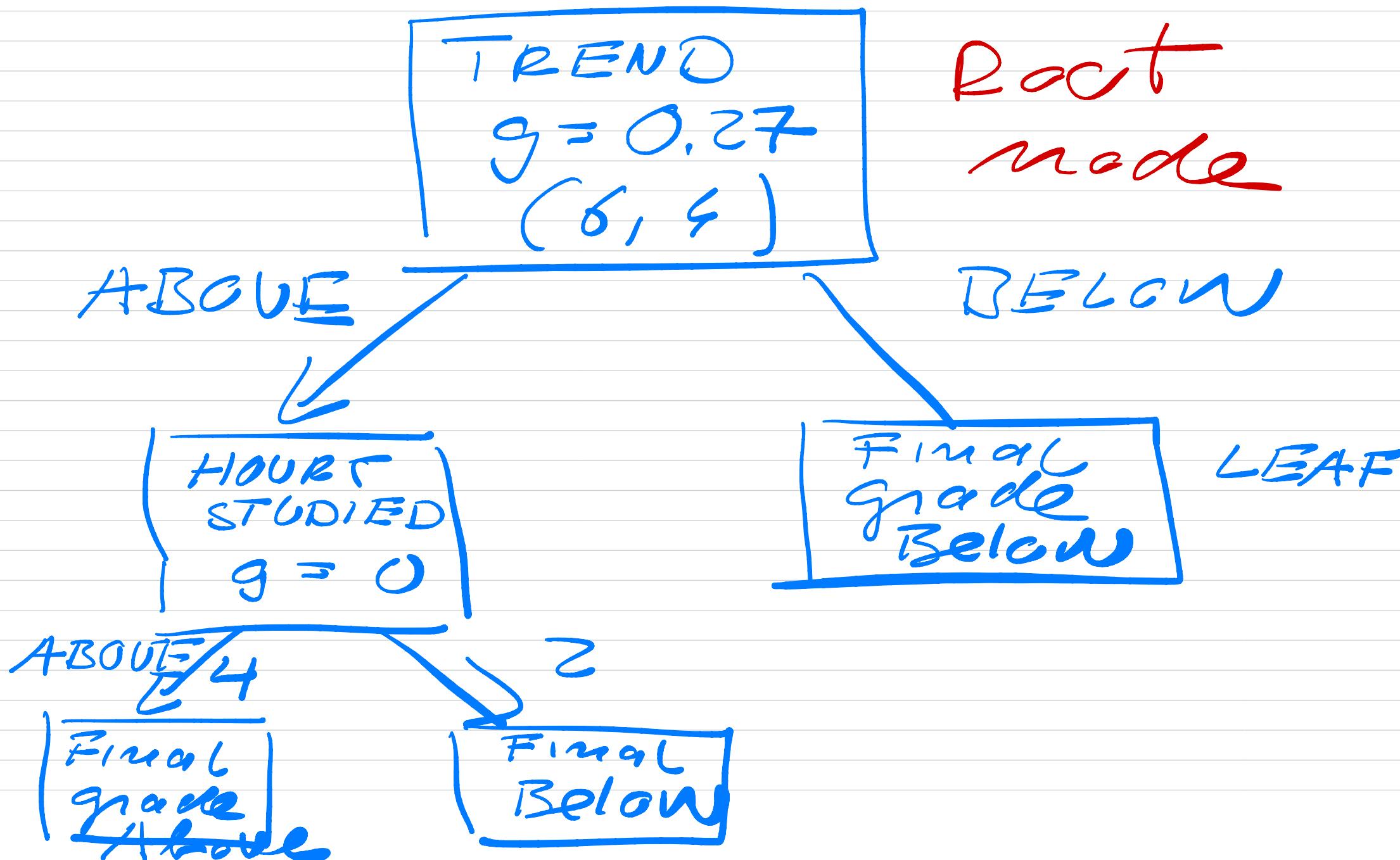
$$\begin{aligned} \text{weighted sum} \\ \text{for hours slept} \\ = 0.53 \end{aligned}$$

weighted grade for hours  
studied  $g = 0$

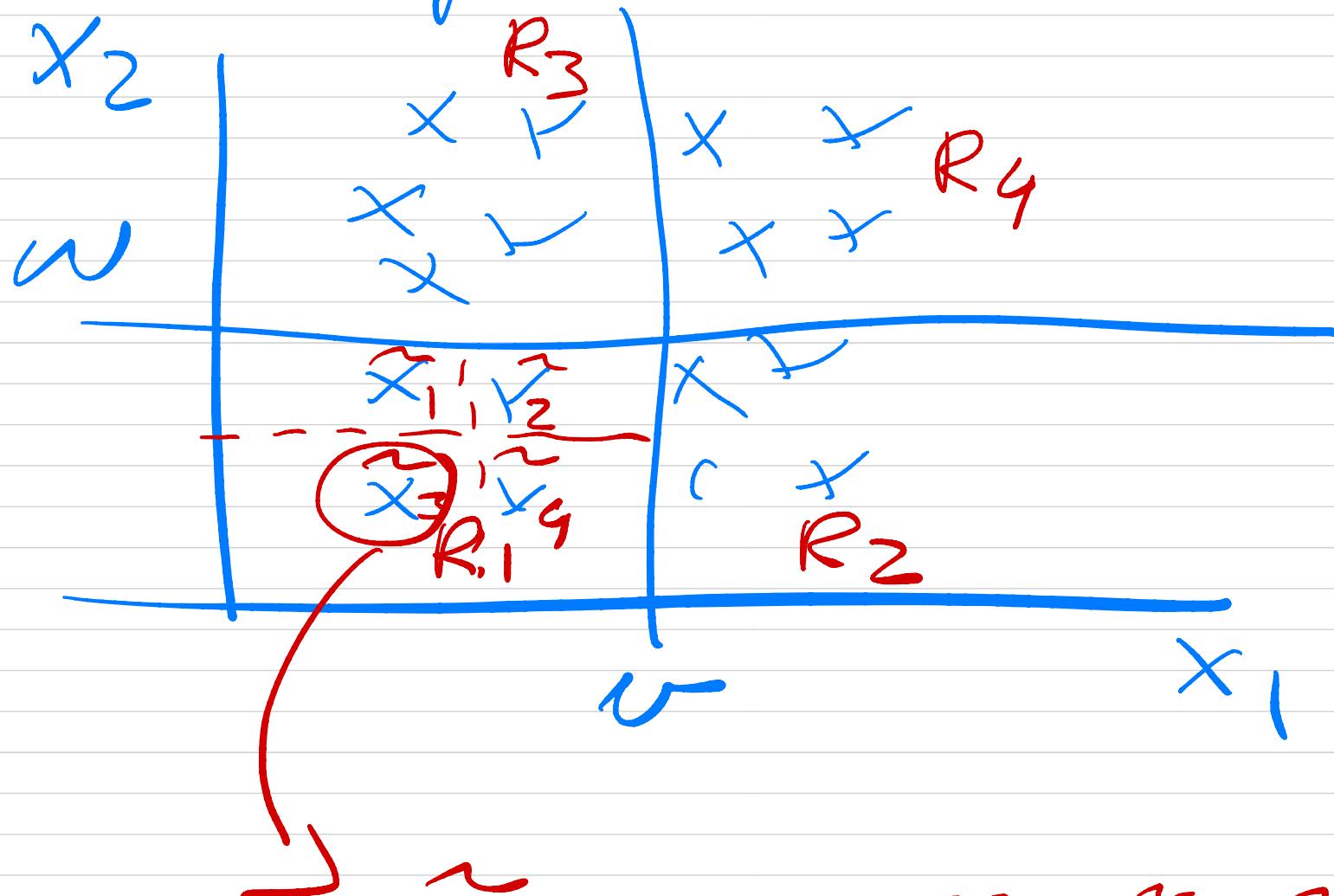
RHS

T	Sleep	Study	Grade
B	H	L	B
B	L	H	B
B	H	H	B
B	L	H	B

# TREE



# Regression Tree



$$\tilde{x}_{R_1} = \frac{1}{4} (\tilde{x}_1 + \tilde{x}_2 + \tilde{x}_3 + \tilde{x}_4)$$

TRUE

$\boxed{\text{if } x_1 < w}$

$R_1 \text{ or}$   
 $R_3$

$\boxed{\text{if } (x_2 < w)}$

TRUE

FALSE

$\boxed{R_1}$

$\boxed{R_3}$

$\boxed{R_2}$

$\boxed{R_4}$

Regression prediction =  
average value in Region  $R_i$