

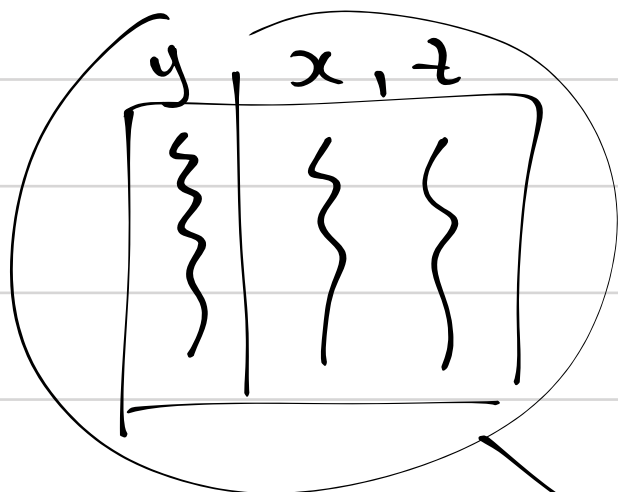
Hi!



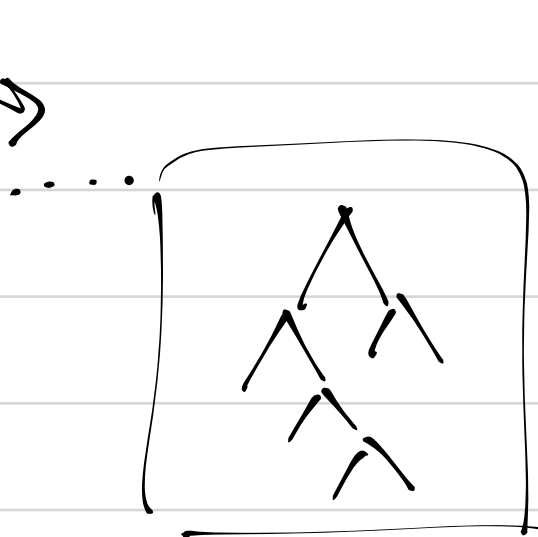
Trees!

→ classification + regression.

Step 1.
Build the tree

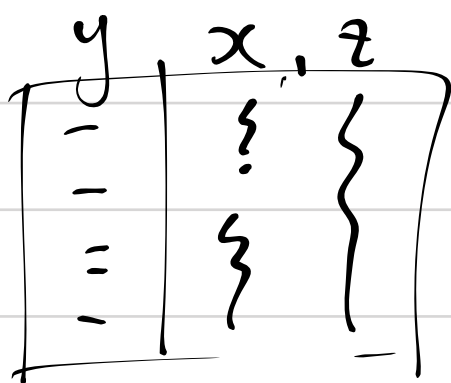


y - response variable
x, z [...] - predictors

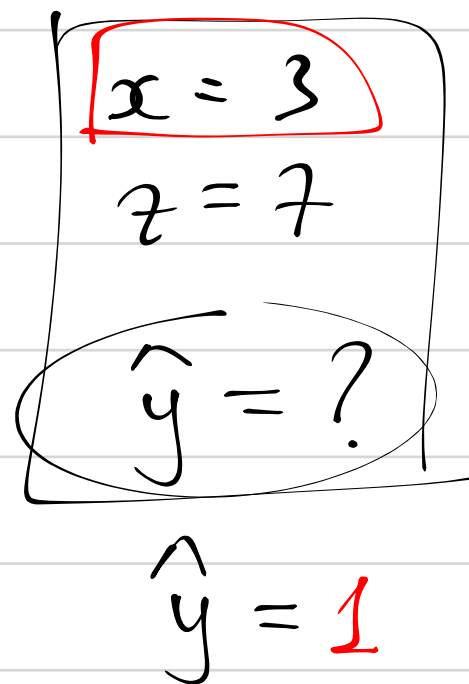
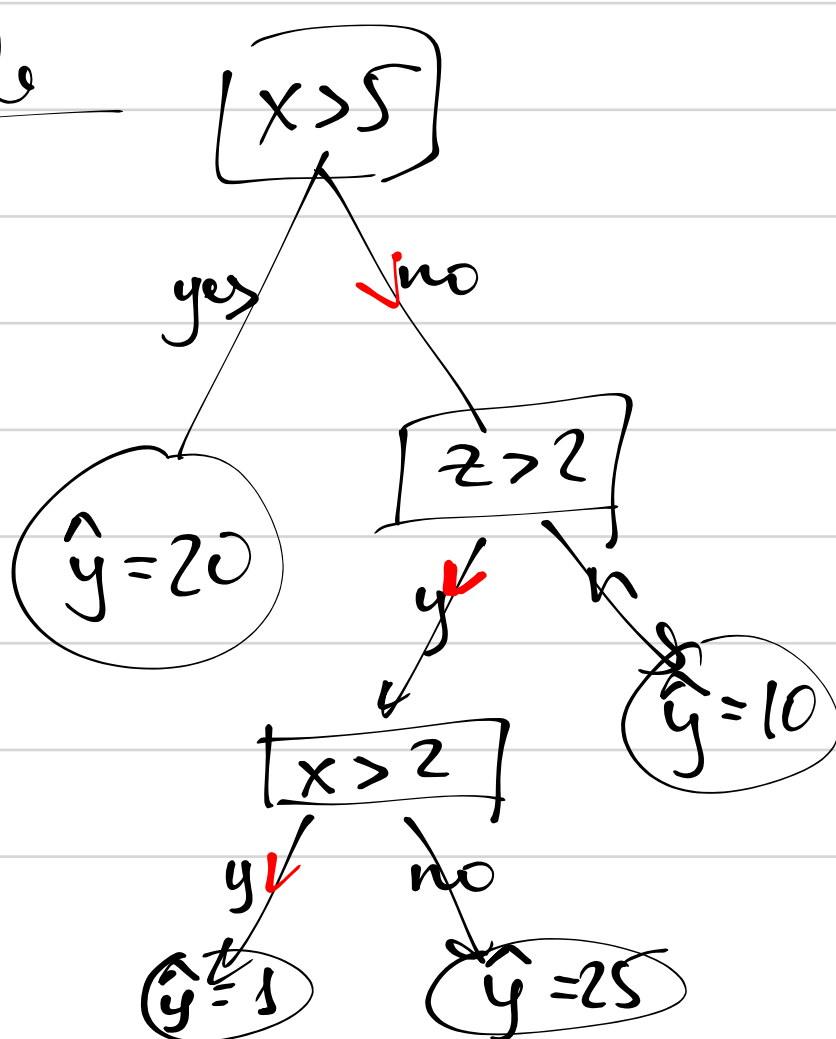


Step 2

Use the tree to forecast y
(the response variable)



Example

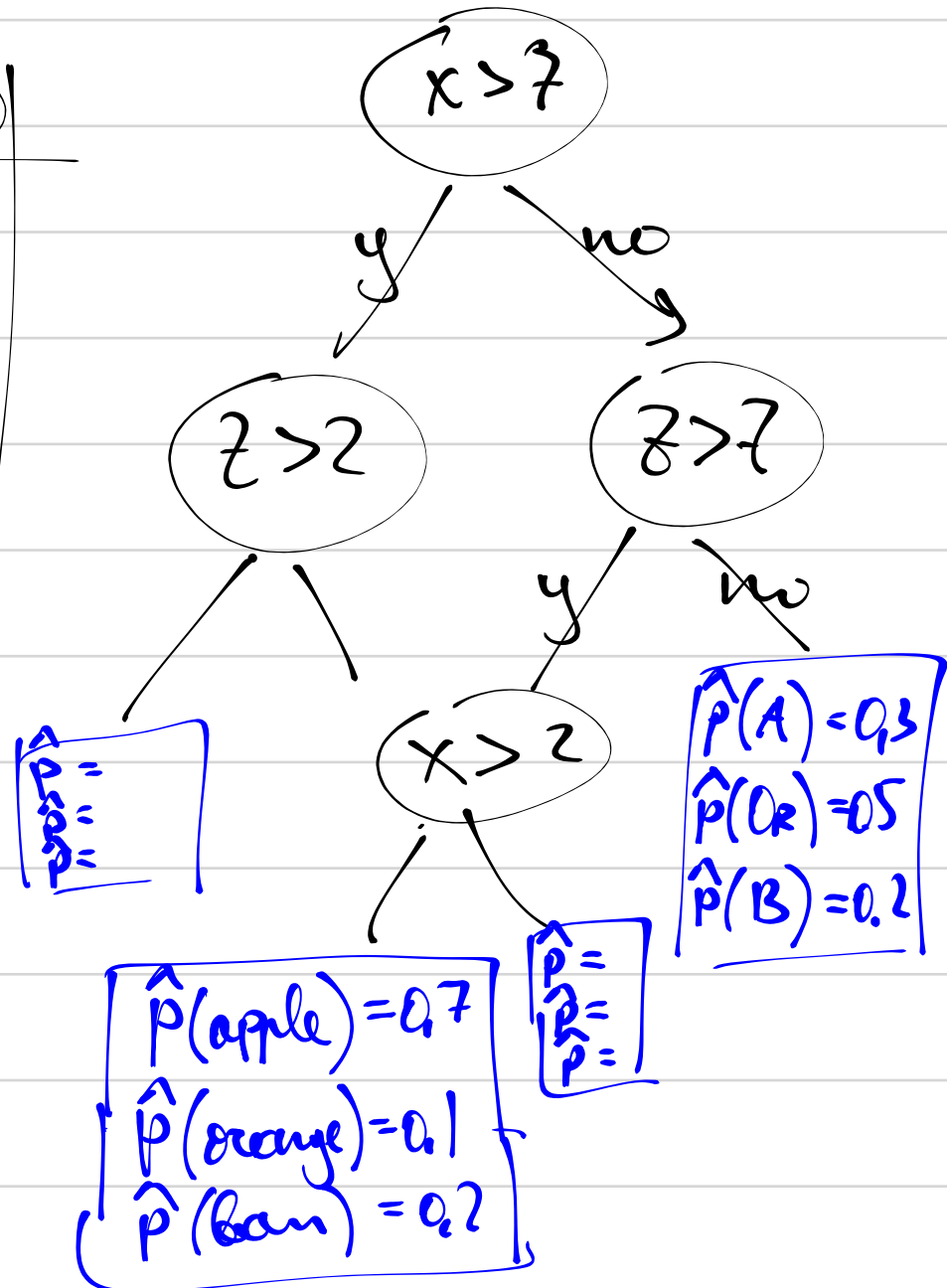


Step 1. How to construct a tree?

Classification task. response variable y takes finitely many values (non-numerical values)

response: predictors:

y	x	z
apple	20	5
orange	10	1
banana	20	3
banana	10	4
apple	-5	5



! There are many variants of the procedure!

Q1. How the node is split in two nodes?

Which function is optimized?

Q2. When we should stop the splitting process?

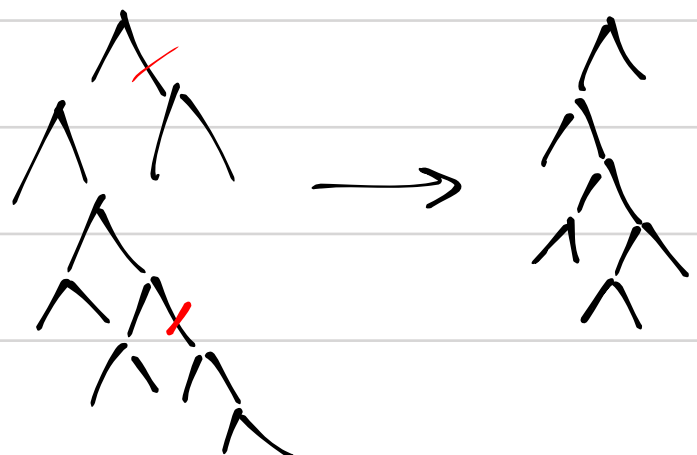
Q3. How the tree processes missing values of predictors?

y	x	z
orange	20	-

Q4. Should we cut the tree after construction? And how?

1. grow a tree

2. cut the tree.



Q1. How to split a node in two nodes?

two criteria: entropy, Gini impurity index

Ask a question

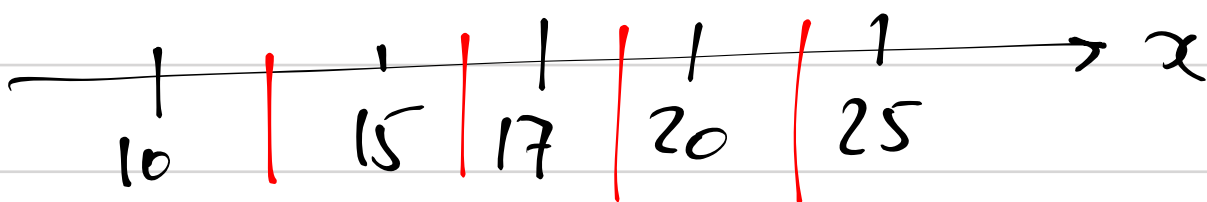
predictor > split.

$x > 7$ $z > 3$ $x > 2.5$

$z > 1.5$

y	x	z
A	20	1
A	10	2
B	15	3
B	17	4
B	20	100
B	25	200

$$H(y) = - \left[\frac{2}{6} \cdot \ln \frac{2}{6} + \frac{4}{6} \cdot \ln \frac{4}{6} \right]$$



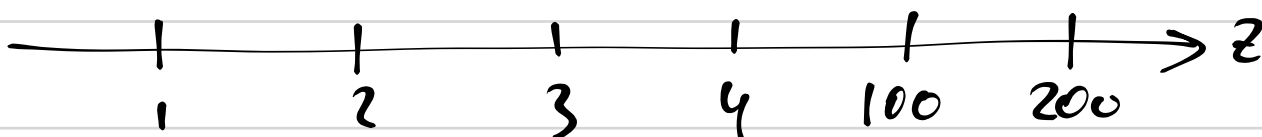
Possible questions are

$x > 2.5$

$x > 16$

$x > 12.5$

$x > 22.5$



$z > 1.5$

$z > 2.5$

$z > 3.5$

$z > 52$

$z > 150$

Calculate all possible conditional entropies $H(y|q)$

y	x	z
A	20	1
A	10	2
B	15	3
B	17	4
B	20	100
B	25	200

$$Q = \{z > 52\}$$

$$H(Y|Q) = ?$$

$$z > 52$$

y		z	
yes	B	20	100
	B	25	200
no	A	20	1
	A	10	2
	B	15	3

$$H(Y|Q=Yes) = -1 \cdot \ln 1 = 0$$

$$H(Y|Q=No) = -\left[\frac{2}{4} \cdot \ln \frac{2}{4} + \frac{2}{4} \cdot \ln \frac{2}{4}\right] = -\ln \frac{1}{2} \approx 0.69.$$

$$\underline{H(Y|Q)} = \frac{2}{6} \cdot 0 + \frac{4}{6} \cdot \ln 2 = \frac{2}{3} \cdot \ln 2 \approx 0.46.$$

Q, Possible questions	$H(Y Q)$
$Q = \{z > 52\}$	0.46
$\{x > 12.5\}$	\vdots
$\{z > 1.5\}$	\vdots
\vdots	\vdots

3 possible questions

lowest cond-l entropy is chosen.

$$\min_Q H(Y|Q) \Leftrightarrow$$

$$\max_Q [H(Y) - H(Y|Q)]$$

Second criterion: joint impurity index

y	
A	
A	
B	
B	
B	
B	

$\mathcal{I}(Y)$ = the probability that two players

will choose different values of Y if they choose indep - by one observation each

$$\begin{aligned}\mathcal{I}(Y) &= P(n_1 \text{ ch } A, n_2 \text{ ch } B) + P(n_1 \text{ ch } B, n_2 \text{ ch } A) = \\ &= \frac{2}{6} \cdot \frac{4}{6} + \frac{4}{6} \cdot \frac{2}{6} = 2 \cdot \frac{2}{6} \cdot \frac{4}{6} = \\ &= 1 - \left(\frac{2}{6}\right)^2 - \left(\frac{4}{6}\right)^2\end{aligned}$$

Q possible quest-s
 $\{z > 52\}$

$\mathcal{I}(Y|Q)$
 $\frac{1}{3}$

the question with lowest $\mathcal{I}(Y|Q)$ is chosen.

$$\begin{aligned}\min_Q \mathcal{I}(Y|Q) \\ \max_Q [\mathcal{I}(Y) - \mathcal{I}(Y|Q)]\end{aligned}$$

$z > 52$

y	x	z
B	20	100
B	25	200

$$\mathcal{I}(Y|Q=\text{Yes}) = 0$$

y	x	z
A	20	1
A	10	2
B	15	3
B	17	4

$$\mathcal{I}(Y|Q=\text{No}) = \frac{2}{4} \cdot \frac{2}{4} + \frac{2}{4} \cdot \frac{2}{4} = \frac{1}{2}$$

$$\begin{aligned}\mathcal{I}(Y|Q) &= \frac{n_Y}{n} \cdot \mathcal{I}(Y|Q=\text{Yes}) + \frac{n_{\text{No}}}{n} \cdot \mathcal{I}(Y|Q=\text{No}) = \\ &= \frac{2}{6} \cdot 0 + \frac{4}{6} \cdot \frac{1}{2} = \frac{1}{3}\end{aligned}$$

For regression problem

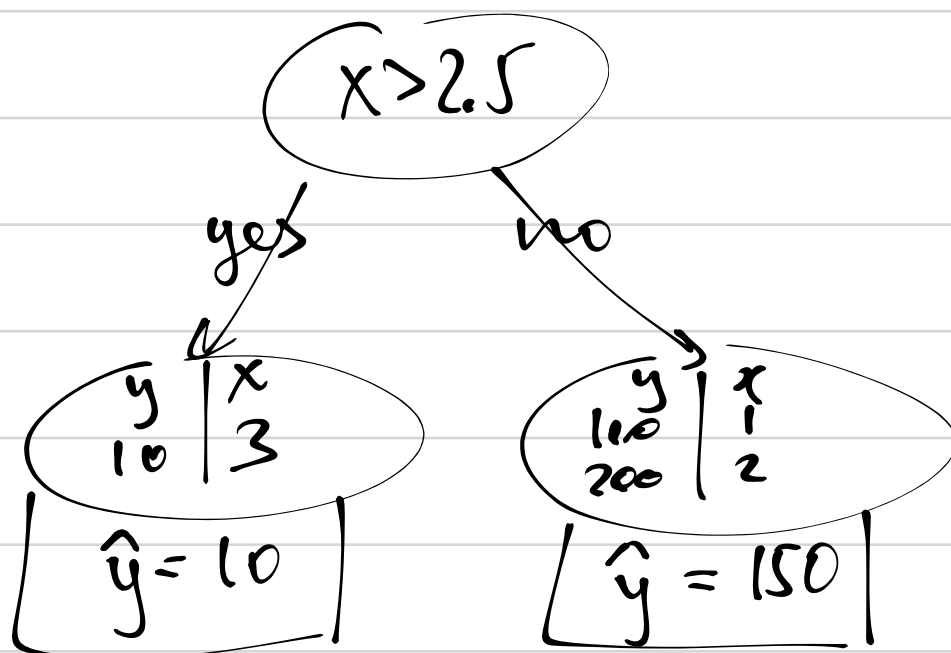
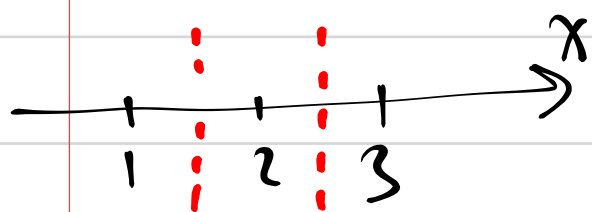
Residual sum of squares \rightarrow min.
Q

y	x
100	1
200	2
10	3

two possible questions.

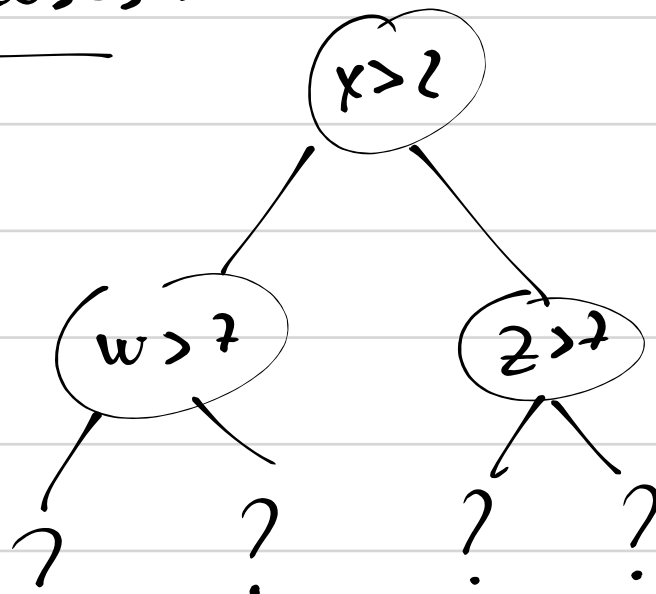
Q	Res SS
$x > 1.5$	

$x > 2.5$? \nearrow 5000



$$\text{Res SS} = \sum_i (y_i - \hat{y}_i)^2 = (10 - 10)^2 + (100 - 150)^2 + (200 - 150)^2 = 2 \cdot 50^2 = 5000$$

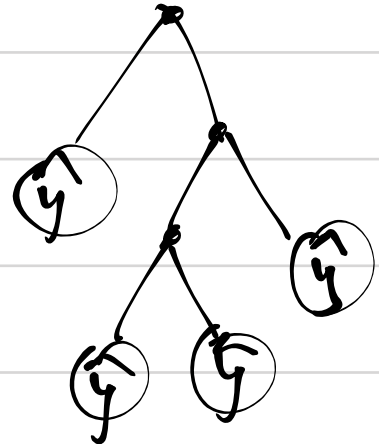
(Q1 b.) Where we should split if there are many cases?



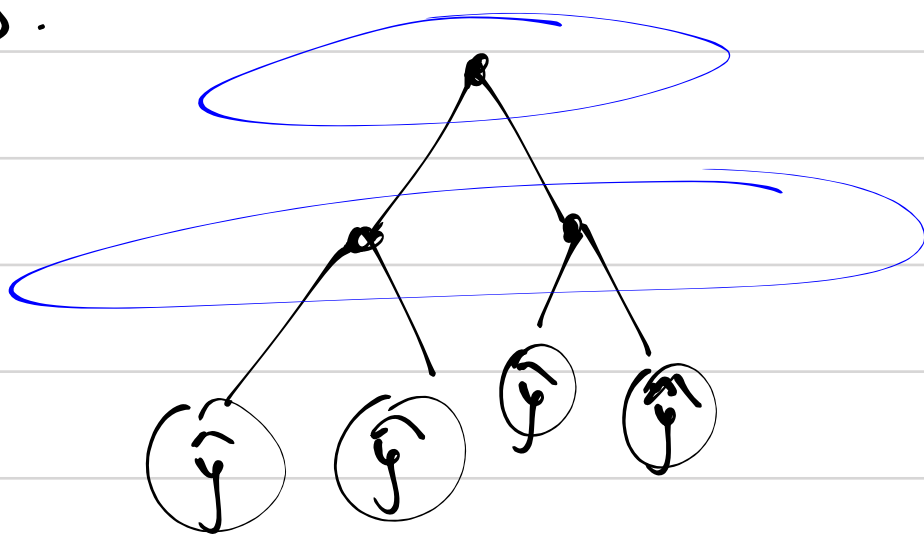
Consider all possible split places and choose the place where the entropy drop is max!

Q2. Where should we stop?

Possible answers: \rightarrow Stop after (3) splits.



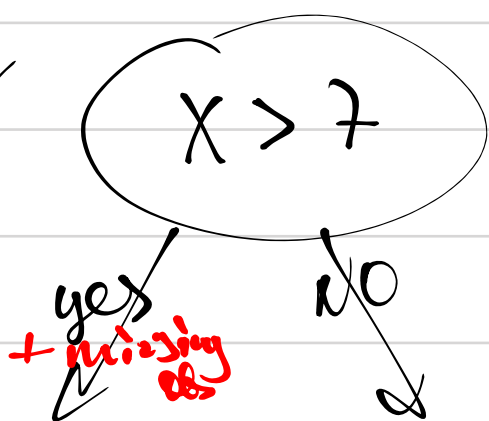
\rightarrow Stop when the tree has 2 levels of questions.



\rightarrow Split the node if there more than (5) observations inside.

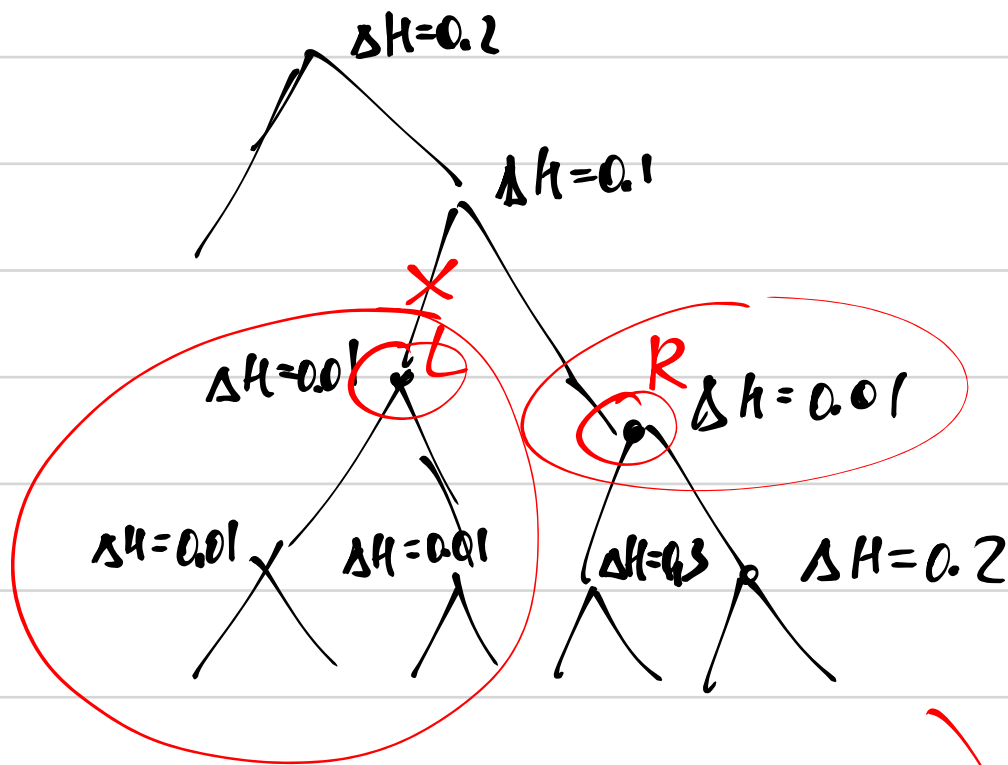
Q3. What to do with missing values in predictors?

- A_1 : Choose the best side (Yes/No)
- A_2 : choose (Y/No) side randomly.
- A_3 : Remove missing obs-s.



Q4. Should we cut the tree?
If yes then how?

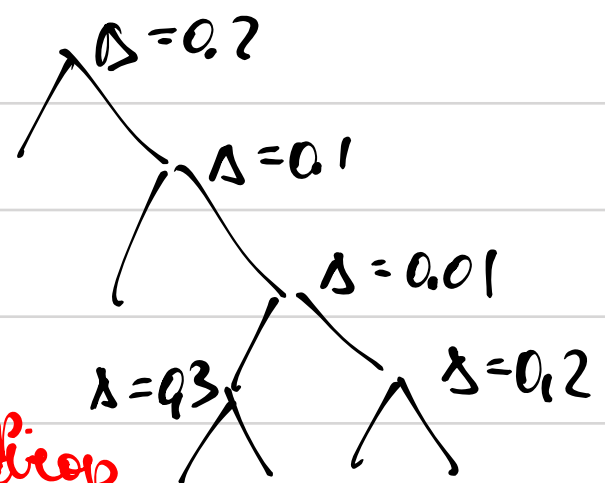
Stage one.



Stage two.
Cut "bad" branches.

price of quest =
0.05

cut the bad part of the tree



classif: ^{max} entropy drop / ^{max} gini imp. index drop
regression: ^{max} res SS drop

Q1A: How to split a node in two?

Q1B: Which node should we split?

Q2: When should we stop?

Q3: How to process missing values?

Q4: How to cut the tree?

→ split the node if it contains more than 5 obs.

→ see answer to Q1A / Q1B

→ compare the entropy drop with a branch price.