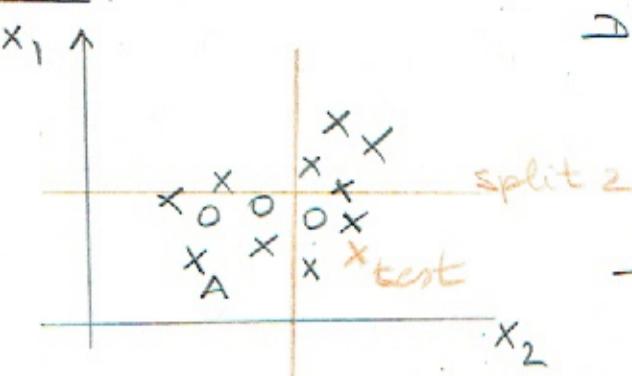


# LESSON 11: Decision Trees, Bagging, Boosting ADABOOST, Gradient Boosting, XGBoost Random Forests

## KD Trees



DATA SET

↳ KNN classification

↳ slow (calculate distances)

memory intensive

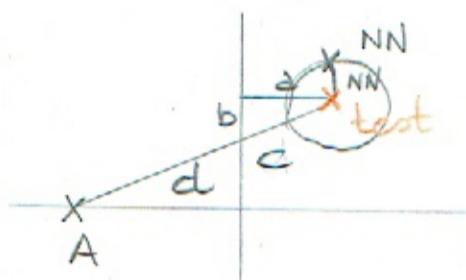
→ CAN we simplify?

↳ YES

we will split data along  
mean of axis with highest  
variance

→ apply KNN on reduced dataset

CLAIM: point A is further from  $x_{test}$  than  $x_{test}$ 's NN



$$c^2 = a^2 + b^2 \rightarrow c = \sqrt{a^2 + b^2}$$

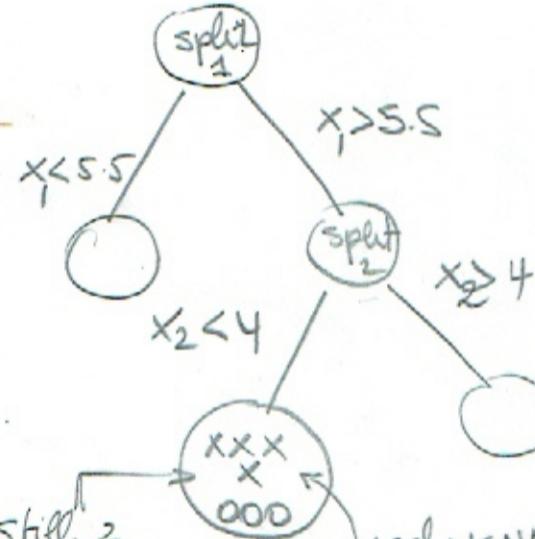
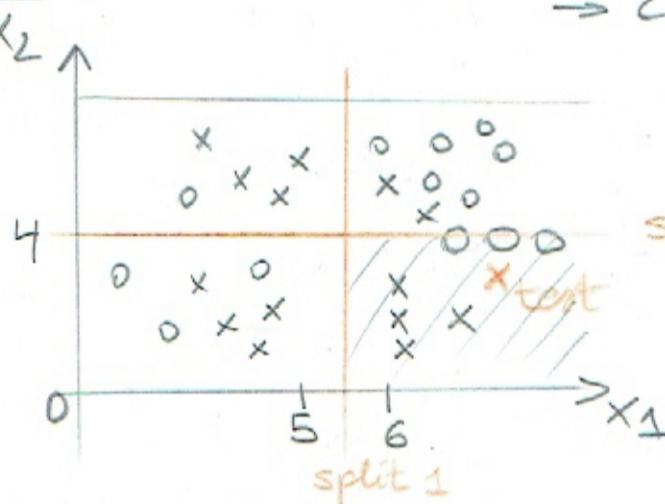
$$\rightarrow c+d = \sqrt{a^2+b^2} + d$$

$$\rightarrow c+d > \sqrt{a^2} + d$$

$$c+d \geq a+d \quad a \geq NN$$

$$\rightarrow c+d \geq a \rightarrow \boxed{c+d \geq NN}$$

Example



Define impurity

functions  $\Leftrightarrow$

Ideally we have pure leaves (one class)

still 2 classes

Apply KNN on leave

## IN PURITY ← GINI ← ENTROPY

↳ likelihood of classifying a point wrong  
Quality measure for a split

### GINI

$$\begin{array}{c} \text{xxx} \\ \text{x} \\ \text{ooo} \end{array} \rightarrow P(X_{\text{test}} = x) = 4/7 \\ P(X_{\text{test}} = o) = 3/7$$

$S$  is set of labels of dataset  $\rightarrow \{x, o\}$

$$S_k = \{(x, y) \in S \mid y = k\}$$

$$k=x \hookrightarrow S_x = \{(x, y) \in S \mid y = x\} \quad P_k = \frac{|S_k|}{|S|}$$

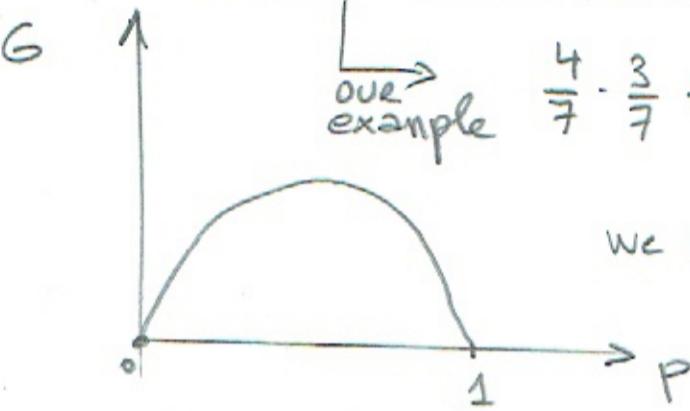
$$S = S_1 \cup S_2 \cup \dots \cup S_k$$

$$G(S) = \sum P_k (1 - P_k)$$

if  $k=2 \rightarrow$  probability of picking 1 class vs. another

$$\text{our example } \frac{4}{7} \cdot \frac{3}{7} + \frac{3}{7} \cdot \frac{4}{7} = 0.48 \quad \hookrightarrow \text{NOT VERY GOOD}$$

we want value close to 1 or 0



### ENTROPY

↳ worst possible distribution in a leaf is when all classes are equally likely to be picked!

$$d_1 = d_2 = \dots = d_K = 1/K \rightarrow \text{WE DONT WANT } d!!$$

KL-DIVERGENCE : measure of relative similarity of 2 distributions

$$\rightarrow KL(P \parallel d) = \sum_{k=1}^K P_k \log \frac{P_k}{d_k} \hookrightarrow 1/K$$

$$\rightarrow D_{KL}(p||d) = \sum p_k \log p_k + \underbrace{\sum p_k \log \frac{p_k}{c^k}}_{\log k \leq \sum_{k=1}^n p_k}$$

$$\rightarrow D_{KL}(p||d) = \sum p_k \log p_k + \log k$$

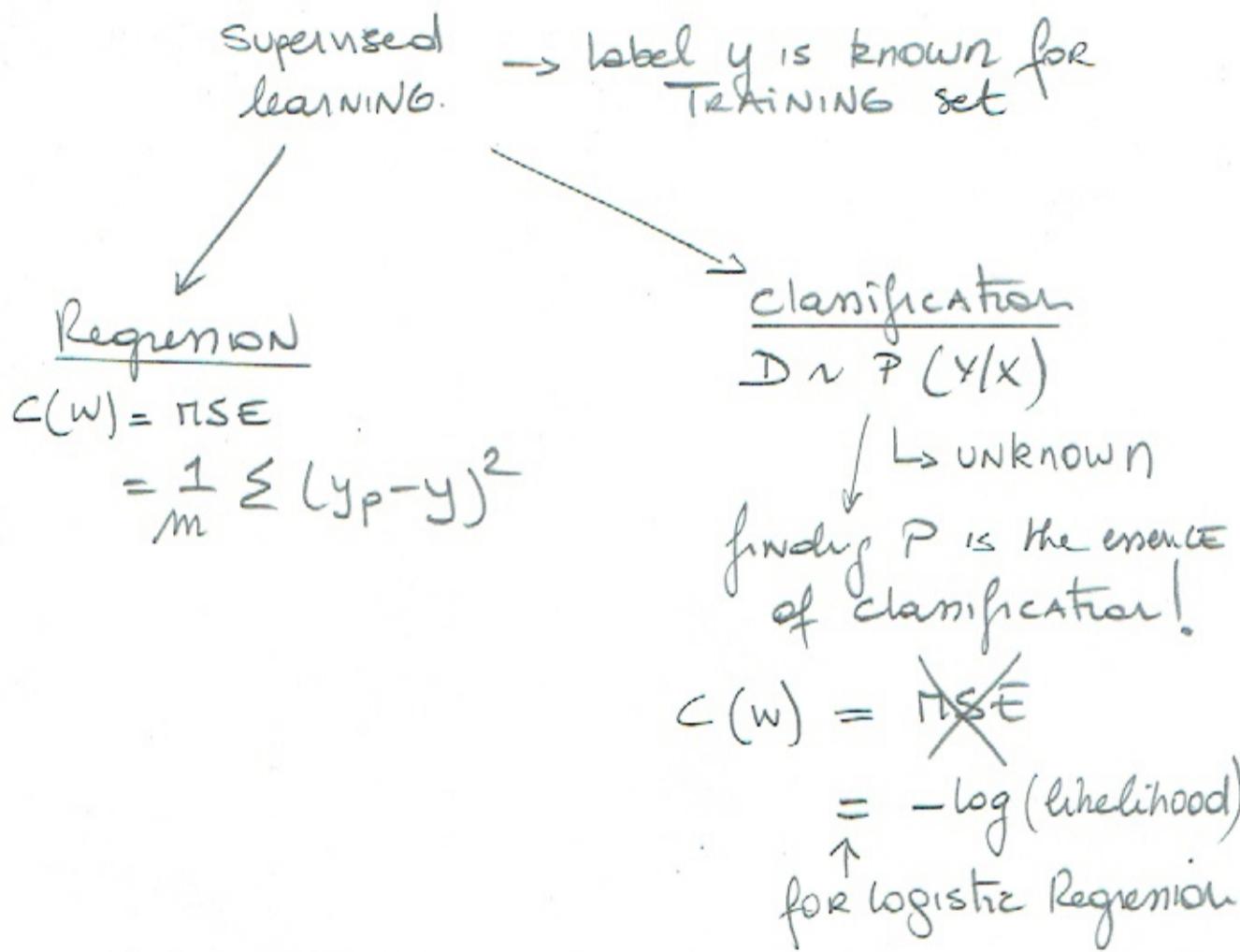
↳ drop when  
max/min as  $c^k$

$$\rightarrow \underset{p}{\text{MAX}} D_{KL}\left(\sum p_k \log p_k\right)$$

or  $\underset{p}{\text{MIN}} D_{KL}\left(-\sum p_k \log p_k\right)$

To be continued  
p6

Optional: ENTROPY - CROSS-ENTROPY - KL divergence



example: Random Variable  $X$   $\xrightarrow{\text{Cat dog}} \underbrace{A \ B \ C \ D}_{\substack{\text{possible outcome} \\ (\text{animal}) \\ \text{pics}}} \quad (4)$

1 TRAINING sample  
(picture of dog)  $\rightarrow [0 \ 0 \ 0 \ 1]$   
 $\hookrightarrow$  distribution  $p$

After TRAINING  
 $\rightarrow$  model  $h$ .  $\rightarrow [0.2 \ 0.1 \ 0.1 \ 0.6]$   
 $\hookrightarrow$  distribution  $q$

How close is  $q$  to  $p$ ?

Definition: INFORMATION content : IC

$IC = 0$  if info doesn't add value to what was known.

ex:  $p(X=D) = 1$  ON TRAINING sample

$$IC(X=D) \approx \frac{1}{p(X=D)} \rightarrow IC(X=D) = f(p(X=D))$$

Assume 2 RVs  $X$  and  $Y$   $\xrightarrow{\text{whether animal is found in SINGAPORE or not}}$

A	1
B	0
C	1
D	1

$$IC(X \wedge Y) = IC(X) + IC(Y) \quad \text{if } X, Y \text{ independent}$$

$$\downarrow \quad \downarrow$$

$$f(p(X=D)) \quad f(p(Y=1))$$

$$f(p(X=D) \cdot p(Y=1)) = f(p(X=D)) + f(p(Y=1))$$

$\Rightarrow f$  is log

$$\rightarrow IC(x=D) = \log \frac{1}{P(x=D)} = -\log P(x=D) \quad (b)$$

$$\rightarrow [IC p = -\log p] \text{ for any distribution } p$$

Expected VALUE  
of a random variable  $\rightarrow E[x] = \sum_x p(x=x) \cdot x$

$$E[IC] = -\log p(x=A) \cdot p(x=A) - \dots - \log p(x=D) \cdot p(x=D)$$

$$\text{ENTROPY } H(p) = - \sum_x p(x=x) \log p(x=x) \quad (1)$$

If certainty  $\rightarrow \log 1$   
 $\rightarrow 0$   
 $\rightarrow$  no contrib

### CROSS-ENTROPY

① Unknown distribution  $P_y(y|x) \xrightarrow{\text{distr.}} p$

② distribution  $q$  obtained through TRAINING

$$③ IC q = -\log q$$

$$(1) \rightarrow E[IC] = C(p||q) = - \sum_x p \log q(x=x)$$

$$\text{CROSS-ENTROPY } H(p||q) = - \sum_x p \log q$$

cost function  
 $\rightarrow$  minimize!  $\rightarrow \frac{\partial}{\partial q_i} = 0$  set  
 WITH constraint  $\sum q_i = 1$

$$\rightarrow \sum q_i - 1 = 0$$

$$\frac{\partial}{\partial q_i} E[\bar{I}_{\bar{C}}] = -\frac{p_i}{q_i} + \lambda \stackrel{\downarrow \text{set}}{\uparrow} = 0 \quad [Lagrange] \quad (6)$$

$$\rightarrow p_i = \lambda q_i \rightarrow \lambda = \frac{p_i}{q_i}$$

$$\rightarrow p_i \approx q_i \text{ if } \lambda \approx 1$$

KL-divergence → relative entropy

↑ measure  
for similarity  
of  $p$  and  $q$

$$D_{KL}(p \parallel q) = H(p \parallel q) - H(p)$$

↳ measure for comparing distributions  $p$  and  $q$

$$H(p) = \sum p_i \log \frac{1}{p_i}$$

$$= \sum p_i \log \frac{q_i}{p_i} \cdot \frac{1}{q_i}$$

$$= \sum p_i \left( \log \frac{q_i}{p_i} - \log \frac{1}{q_i} \right)$$

$$= - \underbrace{\sum p_i \log \frac{p_i}{q_i}}_{\text{relative entropy}} + \underbrace{\sum p_i \log \frac{1}{q_i}}_{\text{cross-entropy}}$$

→ describes mismatch between  $p$  and  $q$ .

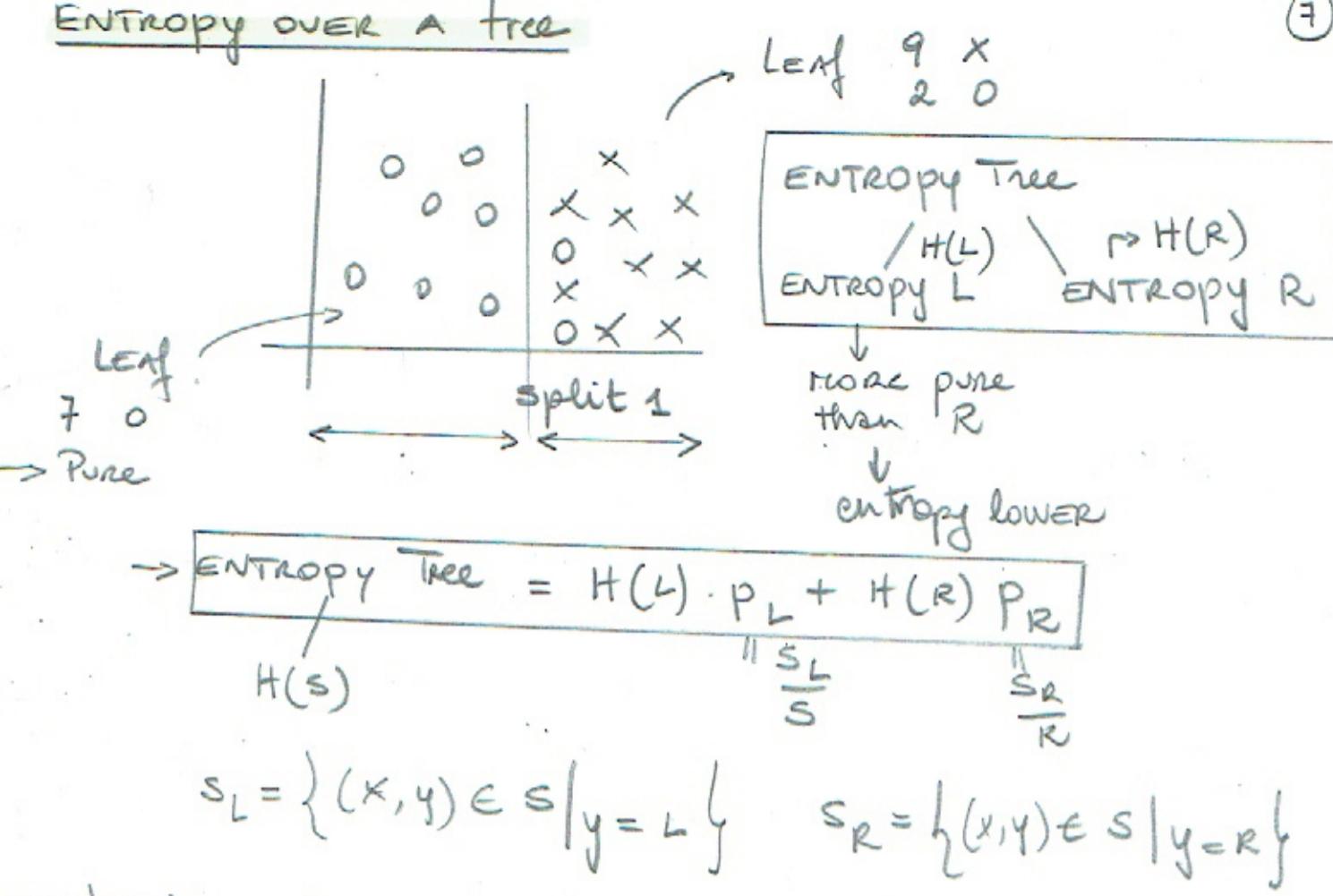
$$D_{KL}(p \parallel q)$$

now we continue lesson from page 3

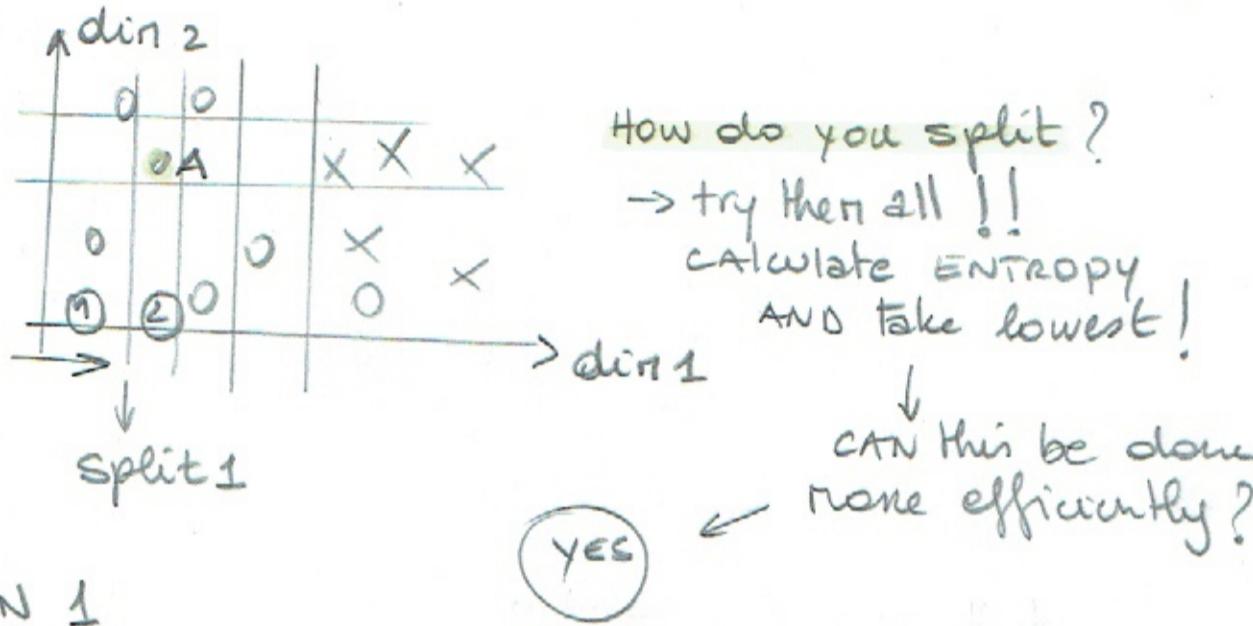
$$\min_p D_{KL} \rightarrow \min_p \left( - \sum p_k \log p_k \right)$$

$\underbrace{H(p)}_{\text{ENTROPY}}$

## Entropy over a tree



## Practical



### SITUATION 1

$$H(L) \cdot P_{L_1} = \frac{s_{L_1}}{m}$$

→ keep all calculations from previous split and just take into account new point A

## SITUATION 2

$$H(\zeta) \cdot P_{L_2} = \frac{S_{L_2}}{S}$$

If you add  $x$  to the left side  
 $\rightarrow H(L)$  will change

# ENSEMBLE LEARNING

BAGGING → BREIMAN '96

Boosting

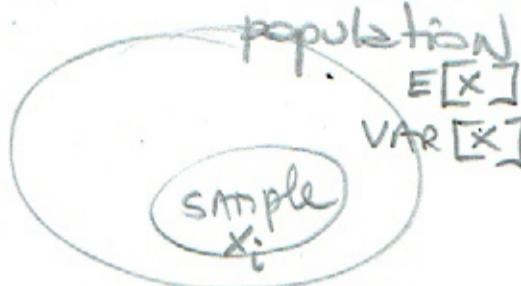
⑧

↳ group a set of learners in order to create a SUPER LEARNER!

→ F. GALTON → weight of ox study

## BAGGING : Bootstrap Aggregation

Can be done on any classifier WITH A VARIANCE problem



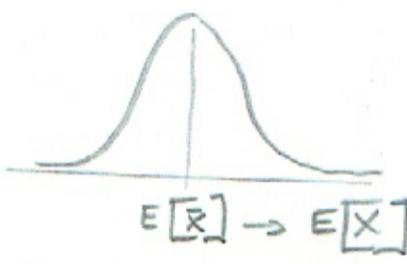
COIN TOSS →  $E[\bar{x}] = 50\%$

if  $n \rightarrow \infty$

$$E[\bar{x}] \rightarrow E[x]$$

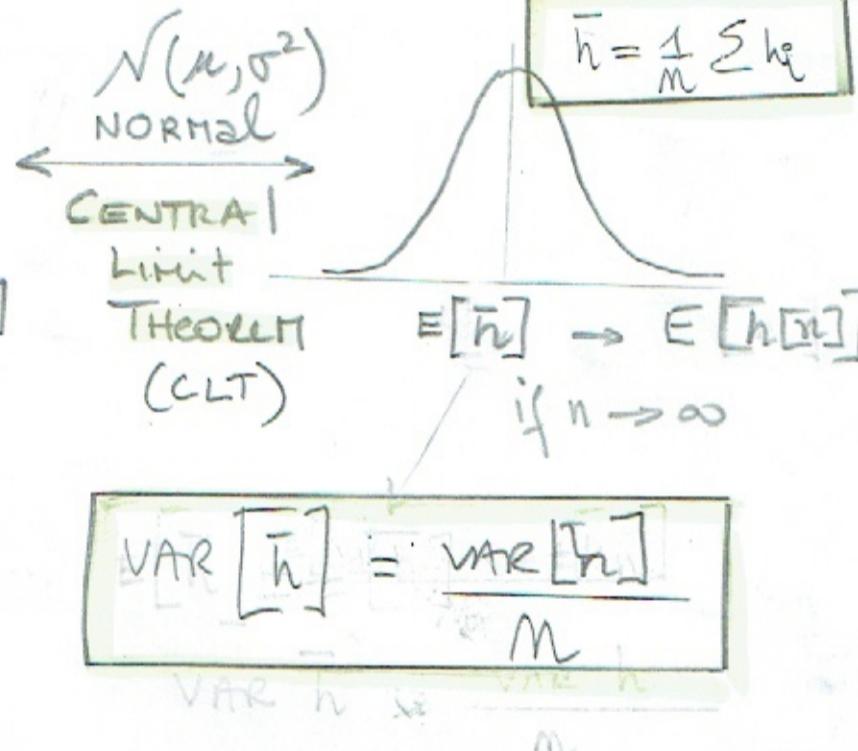
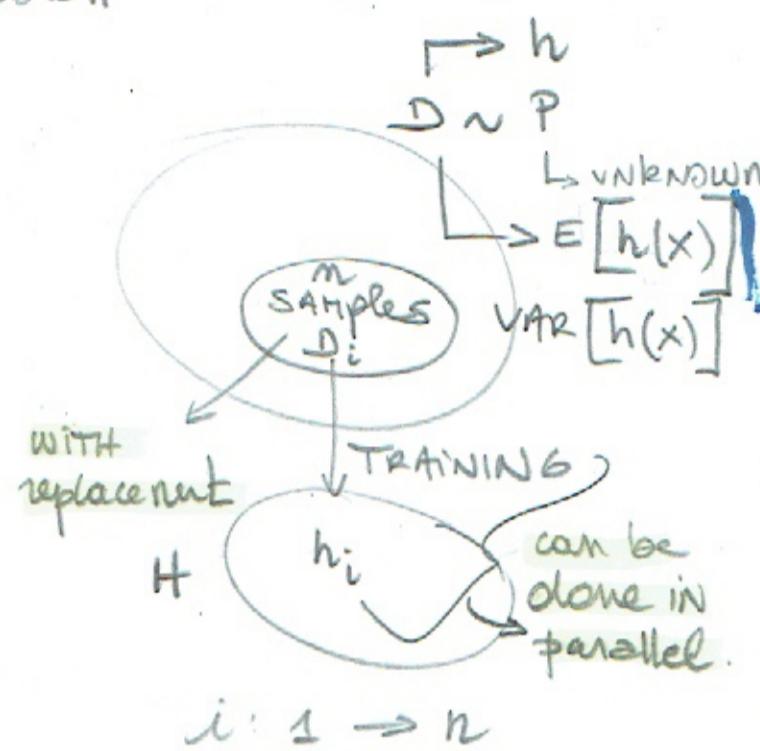
$$\bar{x} = \frac{x_1 + \dots + x_n}{n} \quad \text{WLLN}$$

$x_i$ : 100 COIN TOSSES



$$\text{VAR}[\bar{x}] = \frac{\text{VAR}[x]}{n}$$

This is how you resolve VARIANCE problem!





WLLN  
CLT

assumes i.i.d

(9)

① identical ✓  
 $D_i \sim P$

② independent X

↳ not fulfilled !!

nevertheless: BAGGING works well!

cannot use  
WLLN & CLT

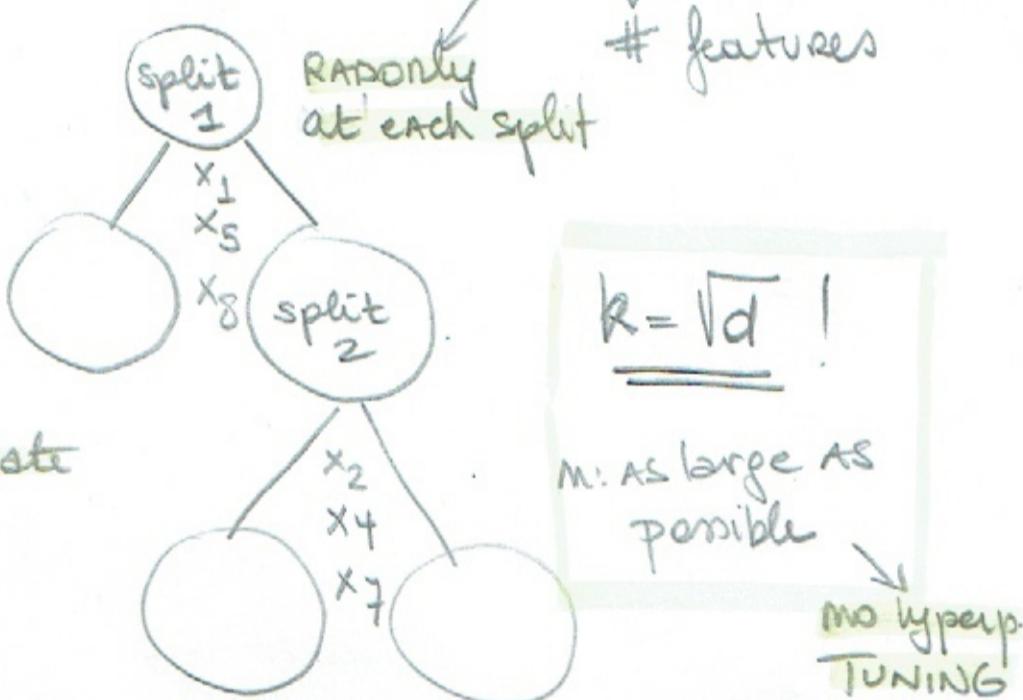


## RANDOM FORESTS

↳ BAGGING ON Trees!  
↳ BREIMAN '99

- ①  $D \sim P$  → run through all trees until all leaves are pure  
 $D_1 D_2 D_n$  → overfitting  
→ split on  $k < d$  features

example:  $d = 9$   
 $k = 3$



→ you can evaluate importance of certain features by splitting on them and calculate IMPURITY!

# Boosting

Freund & Shapire '96  
ADAPTIVE Boosting

## ADABoost

↳ classification only!

IDEA: combine a set of 'weak' learners  
to make a super learner!

sequential buildup

of classifier  $H(x)$  popular at Kaggle competition  
OVER A few rounds

① XG BOOST 2014 Tianqi Chen

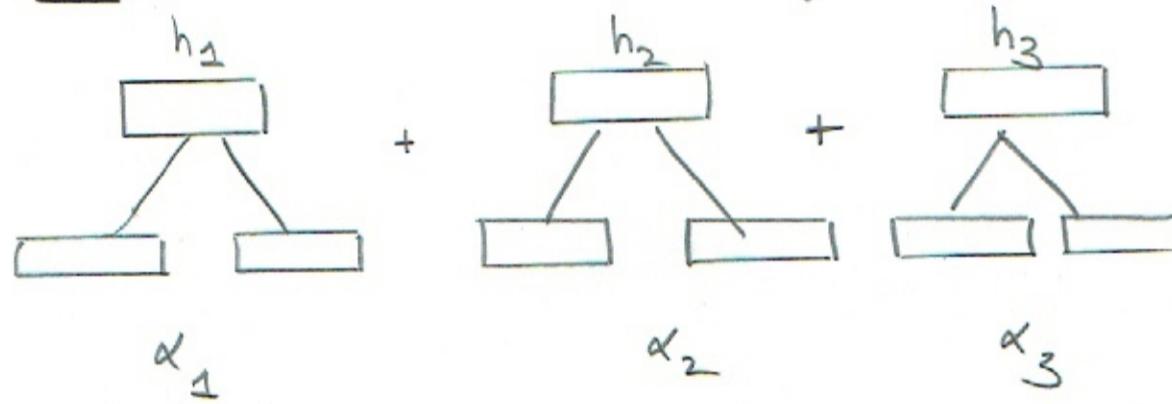
② LIGHT GBM (Microsoft) 2017

high Bias

Fast  
Scalable  
Accurate

ADABoost → stumps (Tree of depth 1)

IDEA:  $h_i \rightarrow$  weak classifiers



$\alpha_1$

$\alpha_2$

$\alpha_3$

$$\alpha_i \sim \frac{1}{\epsilon_i}$$

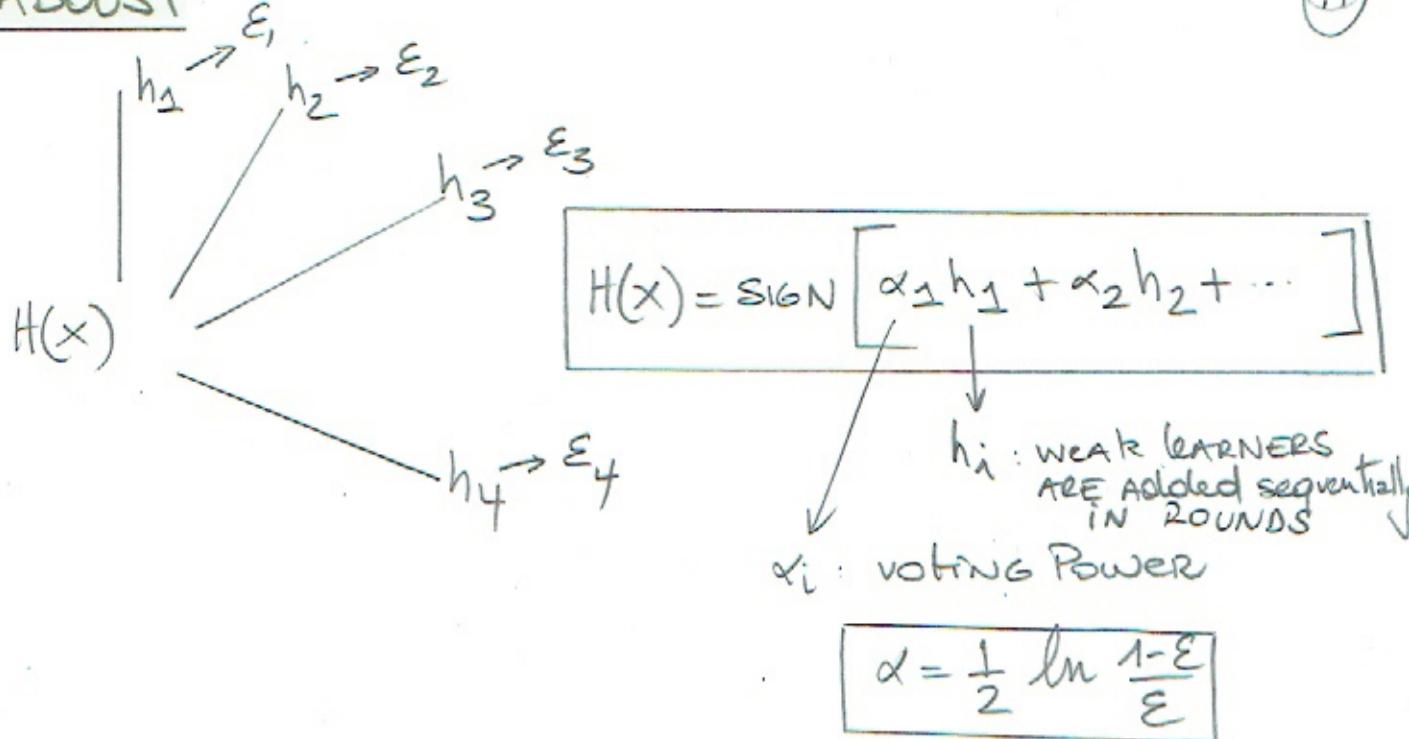
VOTING power ERROR of  $h_i$

$$H(x) = \text{sign} [\alpha_1 h_1 + \alpha_2 h_2 + \dots]$$

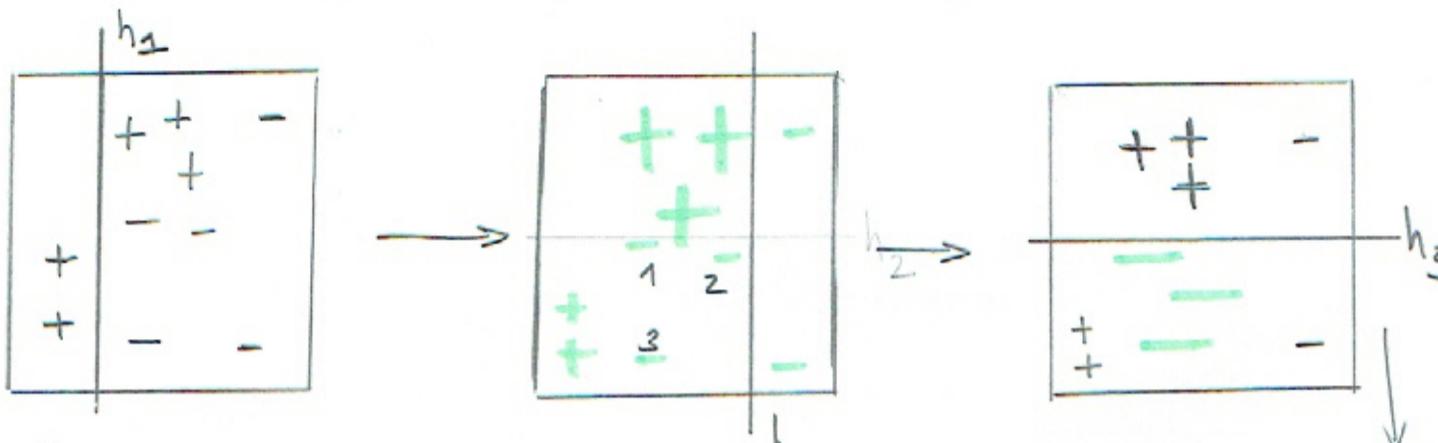
↑  
Stop when target  
reached or when all  
points correctly classi

# ADABOOST

(1)



Weak Learners are stumps (simple trees)



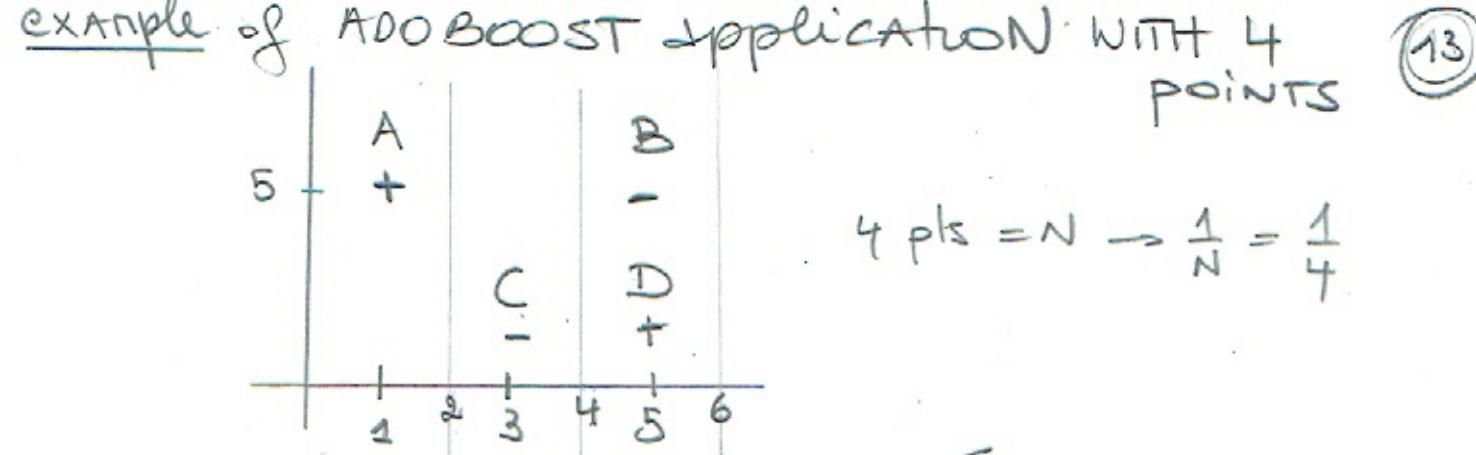
- ① INITIALISE weights  
 $w_i = \frac{1}{n}$   $n$  is # points.
- ② pick best weak learner  
 $h_1$  (lowest error)
- ③  $h_1 \rightarrow 3$  points misclassified
  - ↳ INCREASE weight
  - ↳ reduce weight other points
  - Assign voting power to  $h_1 \rightarrow \alpha_1 \rightarrow \alpha_1 = 0.41$
  - if  $\epsilon$  is large we want  $\alpha$  to be small
  - ↳  $\sum$  weights of misclassified points
- ④ Append  $\alpha_i h_i$  to  $H(x)$ .

12

	+	+	-
	+	-	-
+	-	-	-
+	-	-	-

$$H(x) = \text{sign} \left[ 0.42 \begin{array}{|c|} \hline \text{---} \\ \hline \end{array} + 0.66 \begin{array}{|c|} \hline \diagup \\ \hline \end{array} + 0.93 \begin{array}{|c|} \hline \diagdown \\ \hline \end{array} \right]$$

+++      --      ---  
 we are done  
 as all pts correctly  
 classified!



$$4 \text{ pts} = N \rightarrow \frac{1}{N} = \frac{1}{4}$$

STUMPS of KD Tree  $\sum w_i$  → wrong

	Misclass	RND1	RND2	RND3	int	RND2	RND3	
$x < 2$	D	1/4	3/6	9/24	$w_A$	1/4	1/6	3/24
$x < 4$	C, D	2/4	4/6	15/24	$w_B$	1/4	1/6	6/24
$x < 6$	B, C	2/4	2/6	12/24	$w_C$	1/4	1/6	6/24
$x > 2$	A, B, C	3/4	3/6	15/24	$w_D$	1/4	3/6	9/24
$x > 4$	B, A	2/4	2/6	9/24				
$x > 6$	A, D	2/4	4/6	12/24				

② Assign weight  $\rightarrow \frac{1}{N}$  at start

③ Calculate E

$$\alpha = \frac{1}{2} \ln \frac{1-E}{E}$$

④ Pick best h

	RND1	RND2	RND3	Final
$h$	$x < 2$	$x < 6$	$x > 4$	$x < 2$
$E$	1/4	1/3	3/8	3/10
$\alpha$	$\frac{1}{2} \ln 3$	$\frac{1}{2} \ln 2$	$\frac{1}{2} \ln \frac{5}{3}$	$\frac{1}{2} \ln 3$

$D$  voted correctly  
by (1)  $\alpha(2)$

⑤ Append H(x) =  $\text{sgn} \left[ \frac{1}{2} \ln 3 (x < 2) + \frac{1}{2} \ln 2 (x < 6) + \frac{1}{2} \ln \frac{5}{3} (x > 4) \right]$

$$= \frac{1}{2} \ln 3 (x < 2)$$

Are we done? → no  $\rightarrow$  D is misclassified!

(14)

↓  
update weights → ROUND 2

$$w_{\text{new}} = \begin{cases} w_{\text{old}} \cdot \frac{1}{2} \frac{1}{1-\varepsilon} & \text{if correctly classified} \\ w_{\text{old}} \cdot \frac{1}{2} \frac{1}{\varepsilon} & \text{if wrongly classified} \end{cases}$$

$$w_A = 1/4 \cdot \frac{1}{2} \cdot \frac{1}{3} = 1/6$$

$$w_B = 1/4 \cdot \frac{1}{2} \cdot \frac{4}{3} = 1/6$$

$$w_C = 1/4 \cdot \frac{1}{2} \cdot \frac{4}{3} = 1/6$$

$$w_D = 1/4 \cdot \frac{1}{2} \cdot 4 = 1/2 \quad \checkmark$$

Are we done? → B, C still misclassified!

→ ROUND 3 → we need it to have the necessary weight + classify D correctly

$$w_A = \frac{1}{6} \cdot \frac{1}{2} \cdot \frac{3}{2} = \frac{3}{24}$$

$$w_B = \frac{1}{6} \cdot \frac{1}{2} \cdot 3 = \frac{3}{12} = \frac{6}{24} \quad \checkmark$$

$$w_C = \frac{1}{6} \cdot \frac{1}{2} \cdot 3 = \frac{3}{12} = \frac{6}{24} \quad \checkmark$$

$$w_D = \frac{3}{6} \cdot \frac{1}{2} \cdot \frac{3}{2} = \frac{9}{24}$$

misclassified weights  
 $\leq$  to  $1/2$   
and some for  
correct classified weight

Are we done? → yes as no more points misclassified.

$$w_A = \frac{3}{24} \cdot \frac{1}{2} \cdot \frac{3}{2} = \frac{9}{48}$$

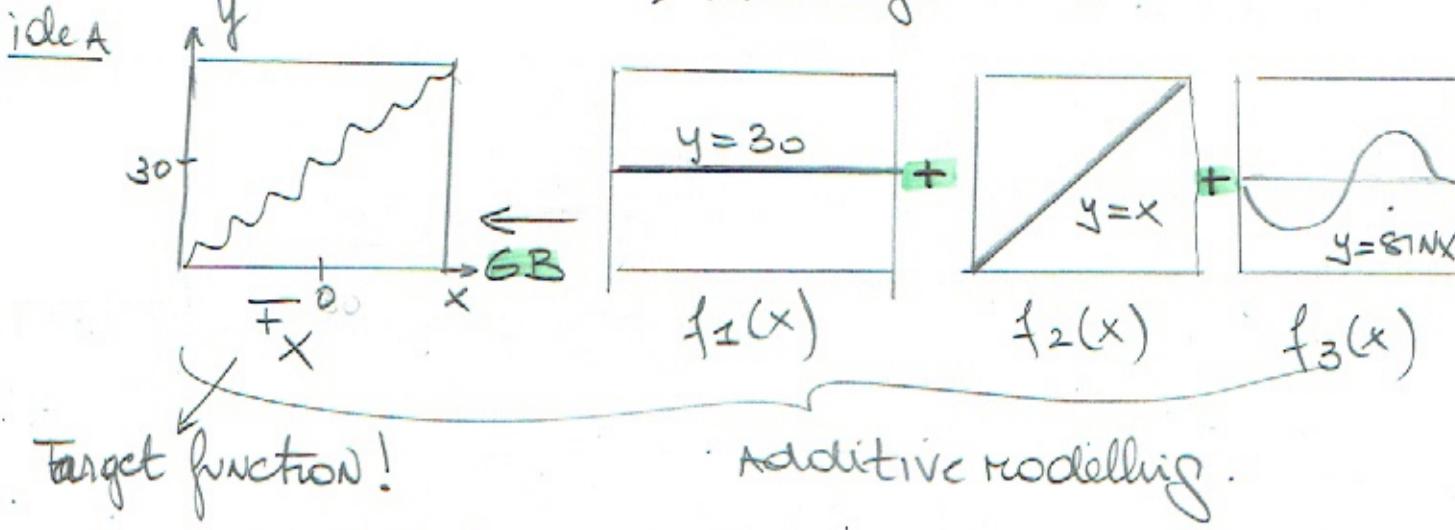
$$w_B = \frac{3}{24} \cdot \frac{1}{2} \cdot \frac{3}{2} = \frac{9}{48}$$

$$w_C = \frac{3}{24} \cdot \frac{1}{2} \cdot \frac{3}{2} = \frac{9}{48} = \frac{3}{16}$$

$$w_D = \frac{3}{24} \cdot \frac{1}{2} \cdot \frac{3}{2} = \frac{9}{48}$$

Gradient Boosting → Algorithm used by search engines  
Trees  $\leftarrow \begin{cases} -\text{xGBoost} \\ -\text{lightGBM} \end{cases}$

→ KAGGLE's favorite!



$$F(x) = f_1(x) + f_2(x) + f_3(x)$$

# functions : M

hyperparameter

$$\rightarrow F_M(x) = \sum_{m=1}^M f_m(x) \quad (1)$$

in machine learning : function = model

$$\parallel y_p \\ h$$

$$\text{OR } F_M(x) = F_{m-1}(x) + f_m(x)$$

for intuition we will use a regression Tree  
 ↪ fitting function through  $(x_i, y_i)$

→ Classification can be found on YouTube

→ StatQuest : GRADIENT BOOST Part 4 :  
 classification details.

MATH is challenging

# EXAMPLE : house price prediction!

(16)

i	size (m <sup>2</sup> ) $x^{(i)}$	price (k\$) $y^{(i)}$	$f_0(x^{(i)})$	$\leftarrow$ simple : 1 feature : size $y - f_0 = r_{11}$	$f_1(x^{(i)})$	$y - f_1 = r_{12}$ house
1	50	60	398	-338	+378	-318
2	100	180	398	-218	+378	-198
3	150	250	398	-148	+378	-128
4	200	300	398	-98	+378	-78
5	250	700	398	302	+438	262
6	300	900	398	502	+438	462
	350	?				

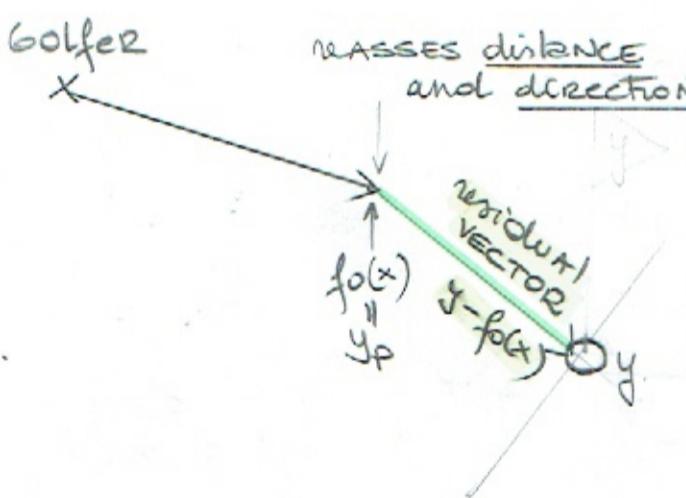
↑  
first prediction residual  
 $r_{11}$

↑  
2nd prediction  
 $r_{12}$

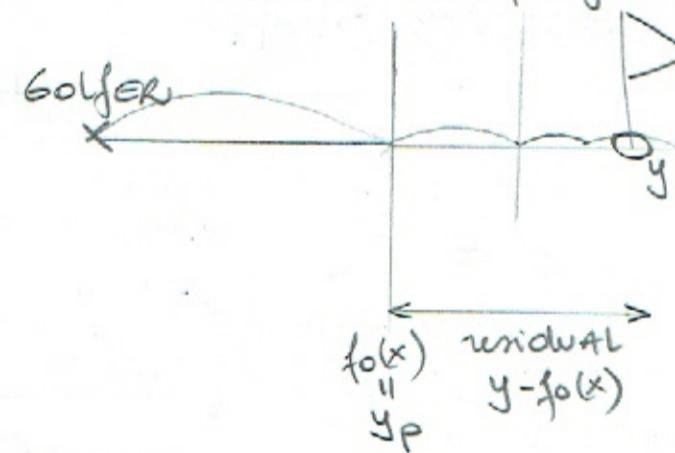
$i: 1 \rightarrow 6$  "samples" → fitting function that best connects  $(x^{(i)}, y^{(i)})$

## Analogy: Golf

### Top view projection



### side View projection



(1)  
p15

$$\rightarrow y_p = \sum_{m=1}^M f_m(x) \quad (2)$$

(17)

WE ARE IN FUNCTION SPACE!

Cost function

$$\frac{1}{2} (y - y_p)^2 = c_1 \quad (3) \rightarrow \text{must be differentiable!}$$

$$\begin{aligned} \frac{dc}{dy_p} &= \frac{1}{2} \cdot 2(y - y_p) \cdot (-1) \rightarrow c(y, y_p) \\ &= -\underbrace{(y - y_p)}_{\text{residual}} \end{aligned}$$

now: back to housing price exercise

$$\textcircled{1} \text{ Initialise model } \rightarrow y_p = F_0(x)$$

$$\textcircled{2} \text{ p17 } \rightarrow F_0(x) = \underset{y_p}{\text{ARGMIN}} \sum_{i=1}^6 c(y^{(i)}, y_p)$$

Initial prediction ↑ target we need to find  $y_p$  that minimizes this sum!

$$\textcircled{3} \rightarrow \frac{1}{2} (60 - y_p)^2 + \dots + \frac{1}{2} (900 - y_p)^2$$

$$\frac{\partial c}{\partial y_p} = -(60 - y_p) - (180 - y_p) - \dots - (900 - y_p)$$

$$= 6y_p - 2390 \stackrel{\text{set}}{=} 0 \rightarrow y_p = \frac{2390}{6} = \underline{\underline{398}}$$

Average of  $y^{(i)}$ 's

$$\Rightarrow F_0(x) = 398$$

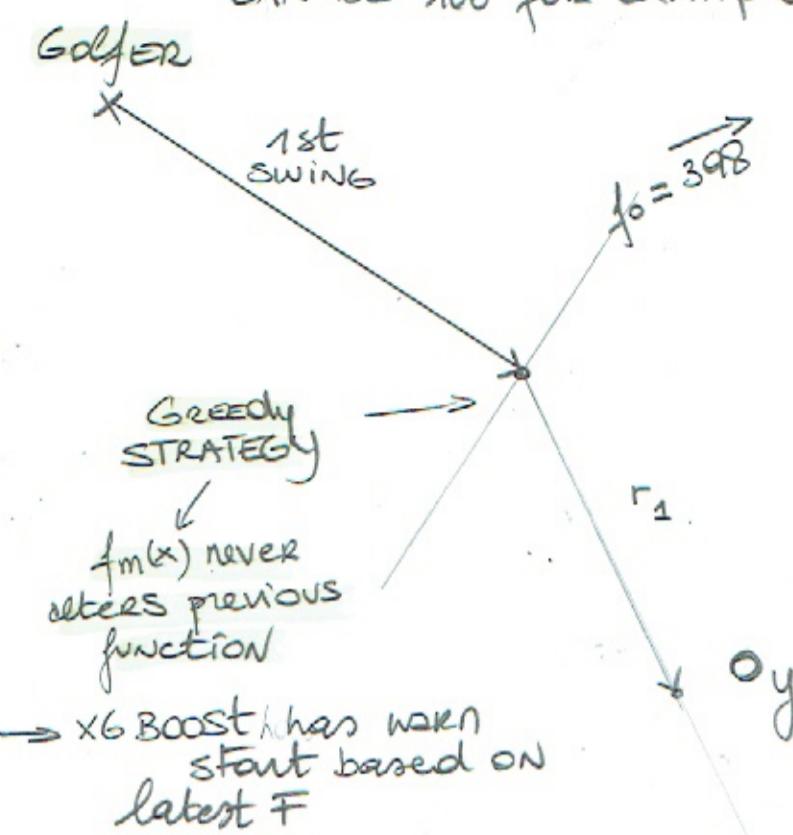
↳ is just a leaf. → ? is predicted to be 398k \$!

② → build  $M$  trees to predict the residuals!

instead of price of the house

CAN be 100 for example

$M: 1 \rightarrow 100$



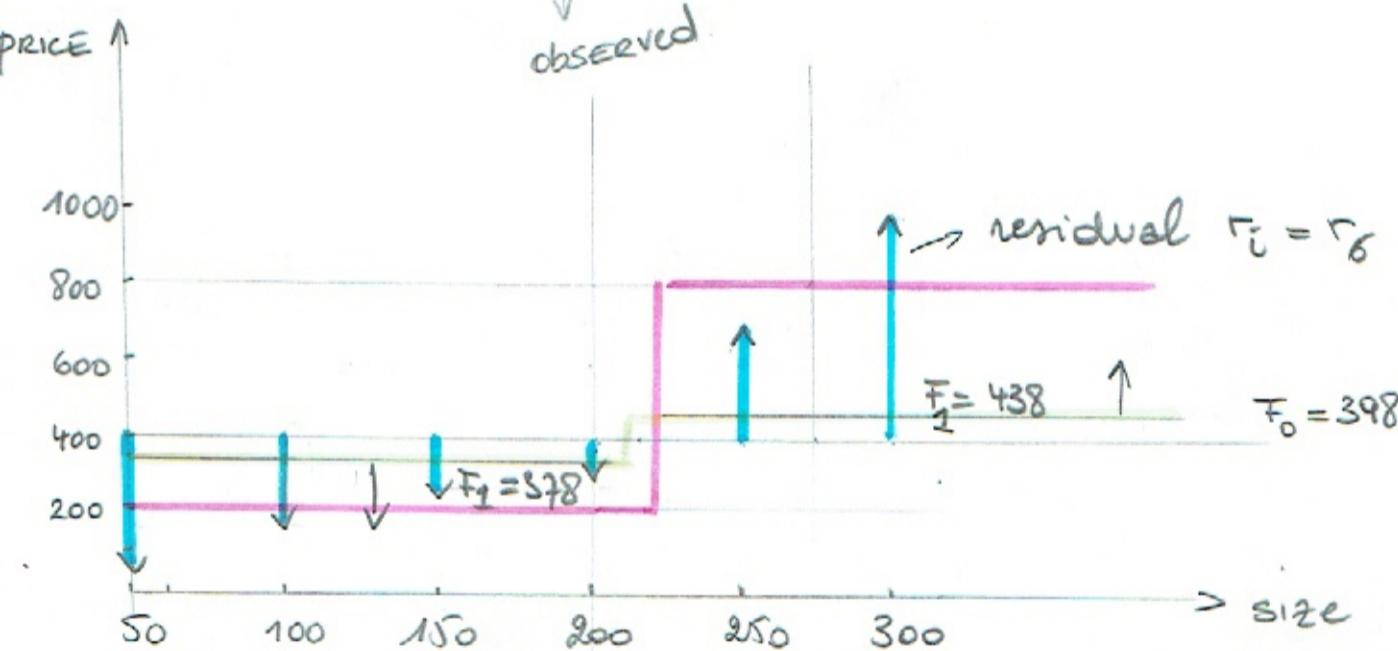
use size of house to predict  $r_{i1}$  etc...

// WE TRAIN ON THE RESIDUALS

$M=1$  → compute residual  $r$

$$r_{i1} = (y^{(i)} - y^{(i)}_{F_0}) \rightarrow \text{fill in in column } r_{i1}$$

1st sample  $\downarrow$       1st tree  $\downarrow$       observed  $\downarrow$   
 $F_0(x)$       Predicted.



A perfect  $r_1$  would be  $F_1(x) = F_0(x) + [y - F_0(x)] \quad (19)$

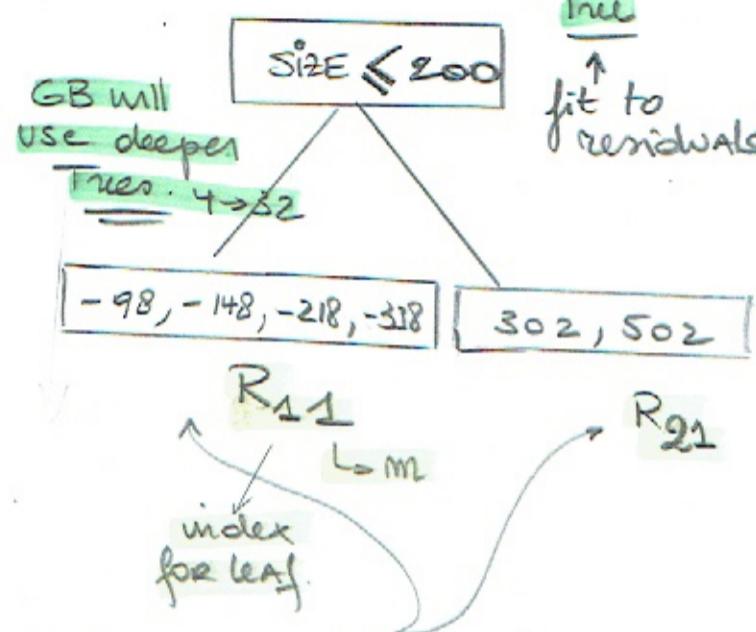
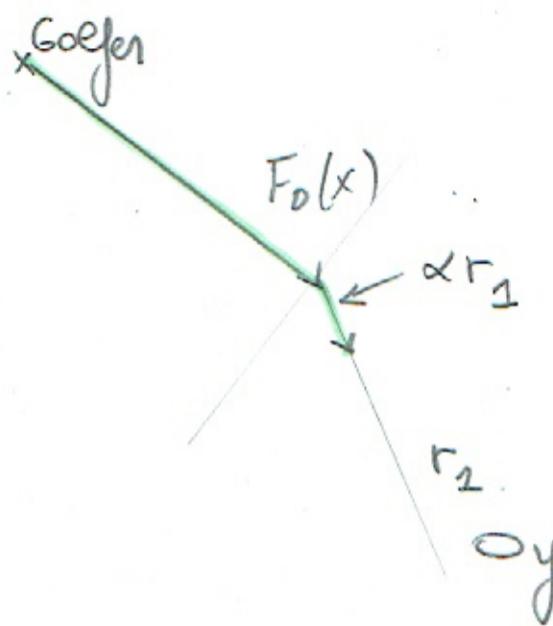
$$= y \quad \begin{matrix} \swarrow \\ \text{node is NOT perfect} \end{matrix}$$

let's define

$$F_m(x) = F_{m-1}(x) + \alpha r_1$$

$r_1$  vector does not end in  $y$ !

learning Rate hyperparameter



Now need to create terminal Regions:  $R_{j^M}$

& calculate outputs

$$R_{11} \rightarrow y_{p_{j1}} = \underset{y_p}{\operatorname{argmin}} \sum_{\substack{x^{(i)} \\ \in R_{11}}} C(y^{(i)}, \underbrace{F_0(x^{(i)})}_{\equiv} + y_p) \quad (1)$$

$\hookrightarrow$  you take previous prediction into account

$$(1) \rightarrow y_{p_{j1}} = \underset{y_p}{\operatorname{argmin}} \sum \frac{1}{2} (y^{(i)} - (F_{m-1}(x^{(i)}) + y_p))^2$$

$$\rightarrow \frac{\partial C}{\partial y_p} = \underbrace{\frac{1}{2} \cdot 2 [60 - (398 + y_p)]}_{\text{Term 1}} + \frac{1}{2} 2 [ ] \dots$$

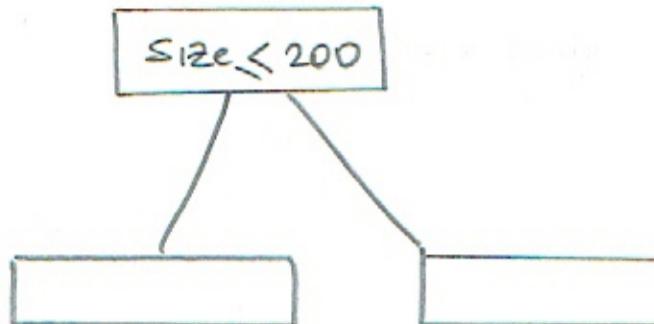
20

$$\begin{aligned}
 \text{TERM 1 : } & (-338 - y_p)(-1) \\
 2 : & (-218 - y_p)(-1) \\
 3 : & (-148 - y_p)(-1) \\
 4 : & (-98 - y_p)(-1)
 \end{aligned}
 \left. \right\} \rightarrow \sum \Rightarrow +4y_p + 802 = 0$$

↓  
 $y_{p_{21}} = 200.5$

Average of rewards  
that ended in leaf  
1

$R_{21}$  →  $y_{p_{21}} = +402$



$$y_{p_{11}} = -200.5 \quad y_{p_{21}} = 402$$

Now update  $F$  →  $\bar{F}_1(x) = \bar{F}_0(x) + \alpha \text{Tree}$

↳ NEW Prediction!  
for each sample

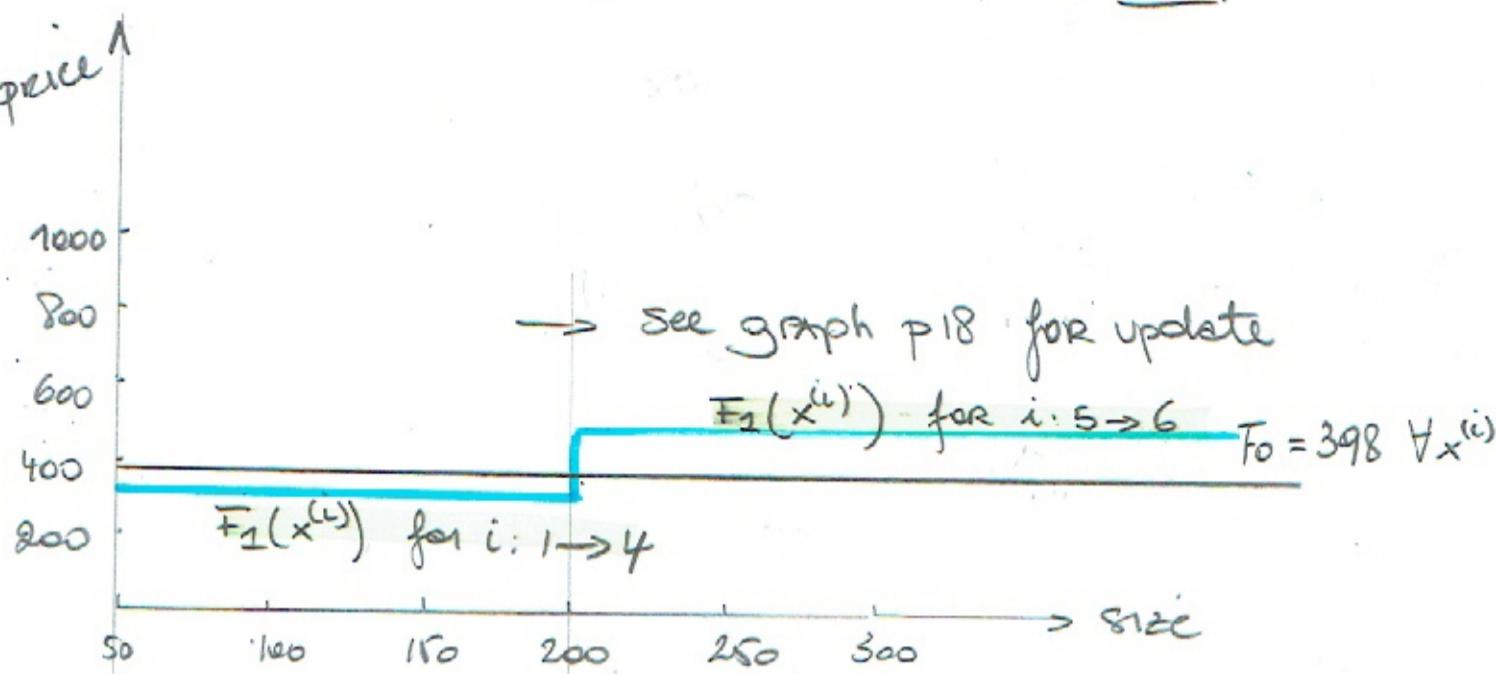
$\alpha \in [0, 1]$

$\bar{F}_1(x) = \bar{F}_0(x) + \alpha \sum_{j=1}^{m \rightarrow \# \text{leaves}} y_{p_j}$

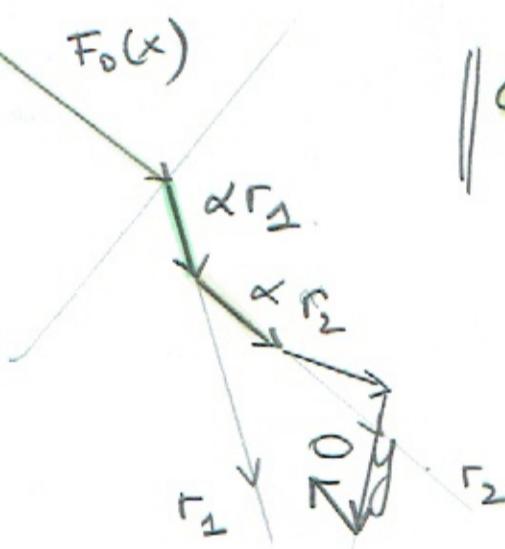
$$F_1(x^{(4)}) = +398 + 0.1(-200.5) = \boxed{378} \quad \checkmark$$

T<sub>2</sub> → same for  $x^{(2)}, x^{(3)}, x^{(4)}$

$$F_1(x^{(5)}) = F_1(x^{(6)}) = +398 + 40 = \boxed{438} \quad \checkmark$$



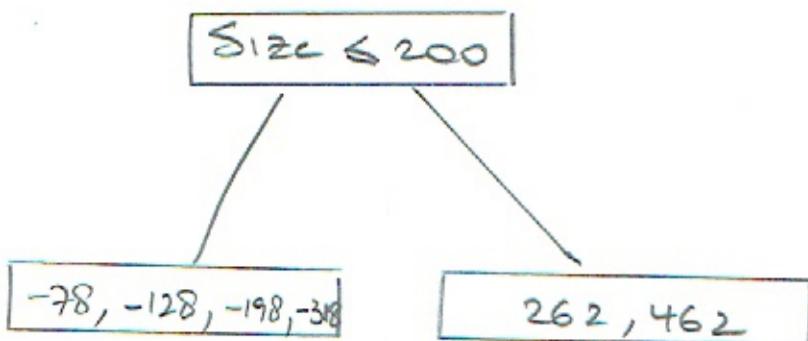
Golfer



// GB is Gradient Descent  
IN function space

↳ you could overshoot if  $\alpha$  too large

if  $\alpha$  too small → slow

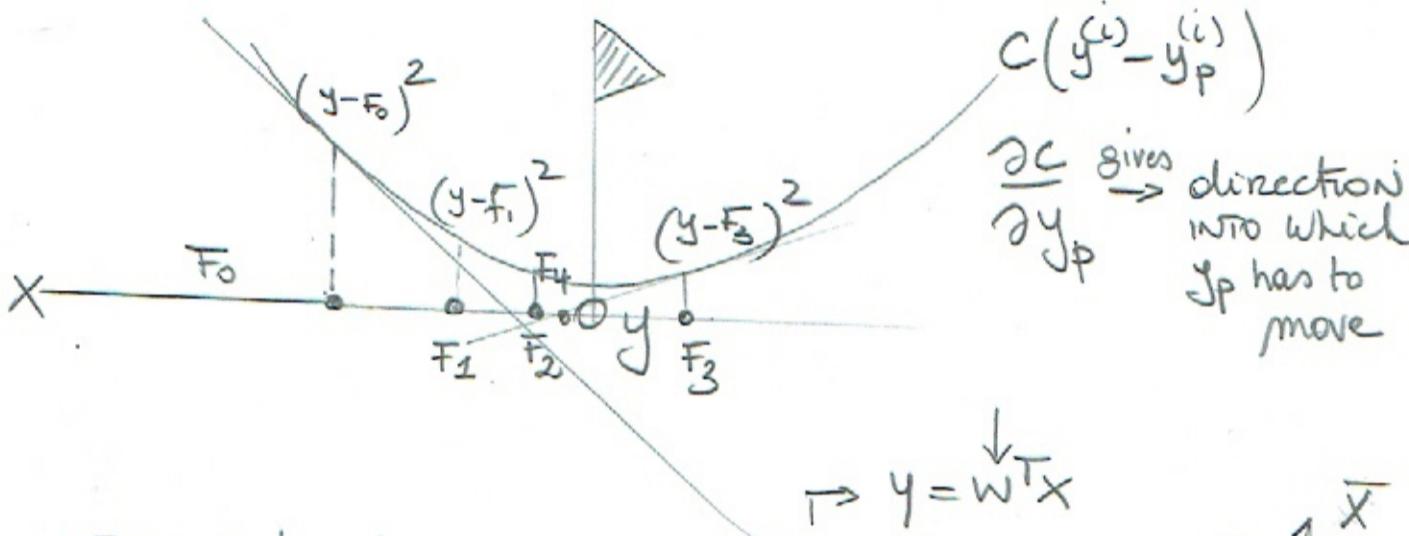
$M=2$ Tree

→ Repeat  $m$  times  $1 \rightarrow M \rightarrow$  gradually move to  $y$

hyperparameters  $\leftarrow M \rightarrow \# \text{ trees } (4 \rightarrow 32)$   
 $\alpha$ : learning rate

Gradient Descent? INTUITION Golfer

side view



GDescent: tweaking parameters/weights of features in order to minimize  $(y - y_p)$

GBoosting: Adding models together sequentially with the goal to optimize the composite model prediction  $y_p$  so as to bring it as close as possible to  $y$ .

boost model output

③ Production ASSURE  $M=2$  (23)

$$F_2(x) = \text{Final output}$$

→ new data  $x^{(7)} = 350 \text{ m}^2$

$\rightarrow \text{OUTPUT Tree 2 : } 362$   
 $\rightarrow \text{OUTPUT Tree 1 : } 402$

$x > 200$

$$F_2(x) = 398 + 0.1$$

Tree 1                      Tree 2  
└───┘ + 0.1            └───┘  
  └──┘   └──┘  
  └──┘   └──┘

$$= 398 + 40 + 36 = 474$$

## Classification

↪ same method but different Cost Function

pseudo  
Abnormal

- log (likelihood) (LR)

## XG Boost

Xtreme Gradient Boosting

ADABoost  
GRADIENT Boosting



BIAIS

XG  
Boost  
library.

2014 Tiengichen VAR

BAGGING  
Random forests

PARAMETER TUNING → Hyperparameters ①  
→ you as data scientist  
will need to choose

① GridSearchCV:  $\alpha$  0.1 0.5 0.7  
 $\beta$  0.3 2 4

	0.1	0.5	0.7
0.3	•	•	•
2	•	•	•
4	•	•	•

works when:  
① # parameters is small  
② # samples is relatively small

↳ 2 parameters each with 3 values!

k-fold cross-validation → 9 SCENARIOS to train & cross-validate  
AVERAGE the performance AT END

TRAINING	CROSS VAL.	TESTING	DATASET
----------	------------	---------	---------

↑ used to tune parameters!

→ Brute Force → takes long time  
↳ does not learn from previous rounds

② Random Search CV

•	•	•
•	•	•

→ A bit more efficient but still does not learn from previous rounds.

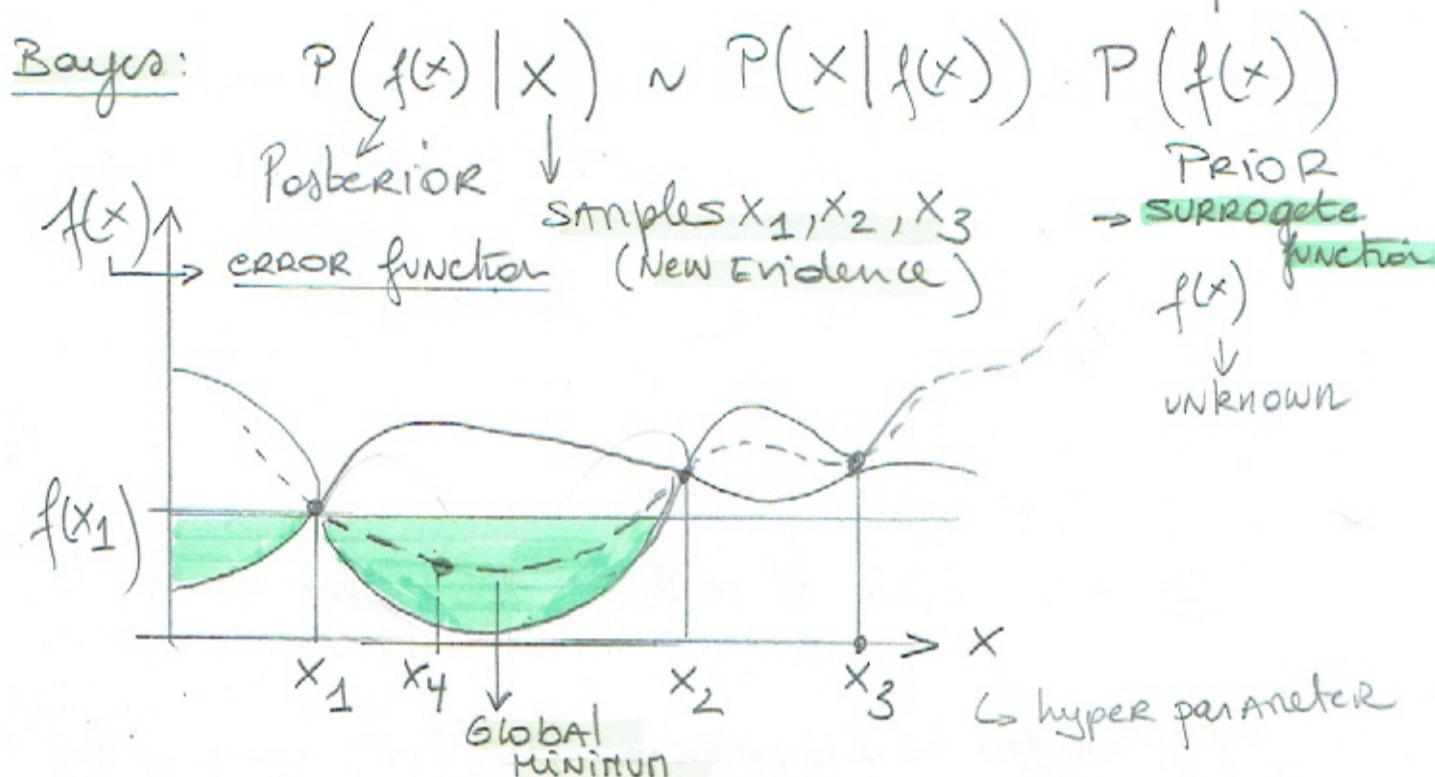
→ RANDOMLY sample from grid!

Conclusion: WE DO TRIAL & ERROR ENGINEERING!  
↓  
BLACK BOX defaults will not lead you best possible results

## Bayesian Optimization

- consider hyper parameters as random variables with a distribution that we don't know
- optimization problem for finding the minimum of a function without calculating gradients  $f(x)$

- 3 options
- ① Gaussian Processes  $\rightarrow P(f(x)|x) \sim N(\mu, \sigma)$
  - ② Random Forests
  - ③ Tree Parzen Estimators



How to find correct zone to sample  $x_4$  from?

- look for zones where the expected improvement is MAXIMUM and where  $f(x) \leq f(x_1)$

$$EI = E \left[ \max(f(x) - f(x_1)), 0 \right] \rightarrow \text{Acquisition function}$$

→ take  $x_4$  in green zone

## parameters XGBoost

③

① learning rate  $\alpha$

$$0.1 \rightarrow 1$$
$$F_1 = F_0 + \alpha r_1$$

② n\_estimators  $\rightarrow F_M = F_{M-1} + \alpha r_M$

$\hookrightarrow M$   
 $(25 \rightarrow 1000)$   $\hookrightarrow \# \text{ trees } M$

③ gamma  $(0 \rightarrow 5)$   $\rightarrow$  regularization factor

④ min\_child\_weight  $\rightarrow$  min # of items in leaf before you continue to split  
 $(2 \rightarrow 10)$

⑤ max\_depth  $\rightarrow$  against overfitting  
 $(1 \rightarrow 32)$

⑥ subsample  $\leftarrow D \begin{cases} D_1 \\ D_2 \\ D_3 \end{cases}$   
⑦ colsample  $(0.5 \rightarrow 1)$   $\downarrow$  analogues with Bagging

split  $k < d$ .