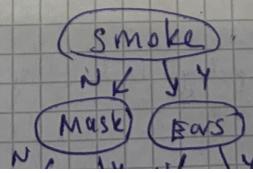


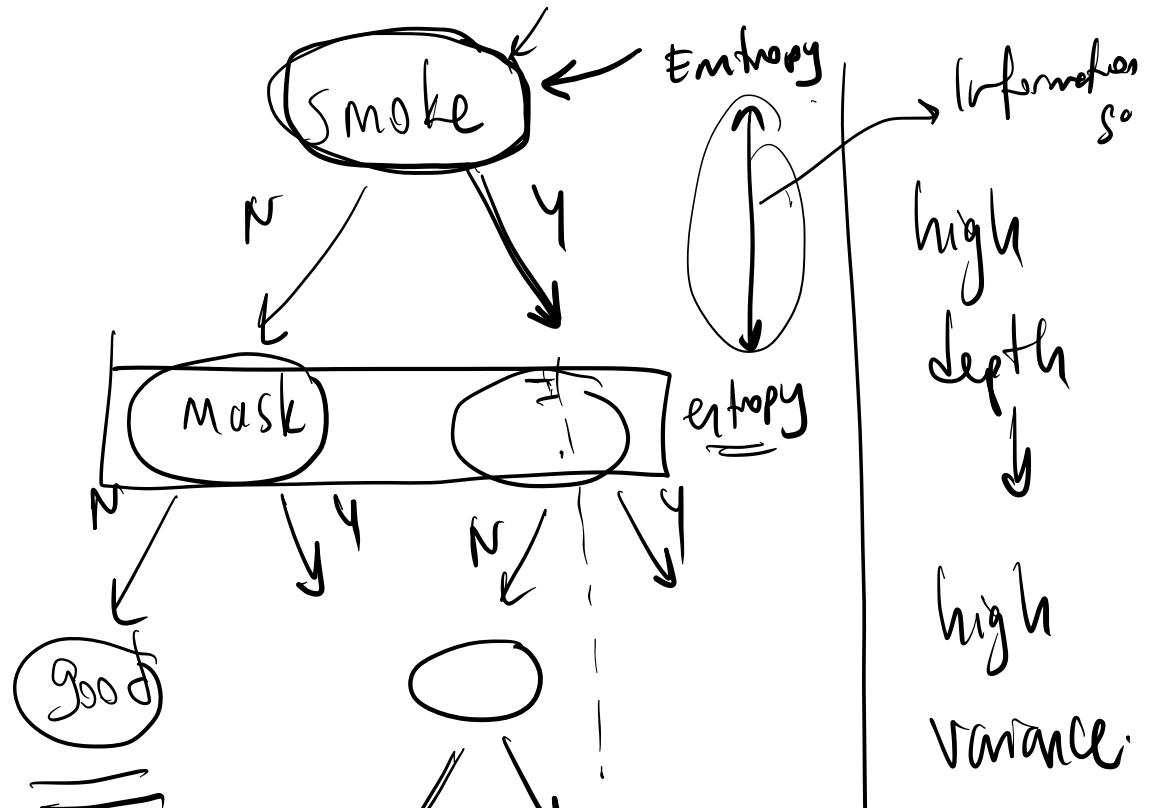
Decision Trees

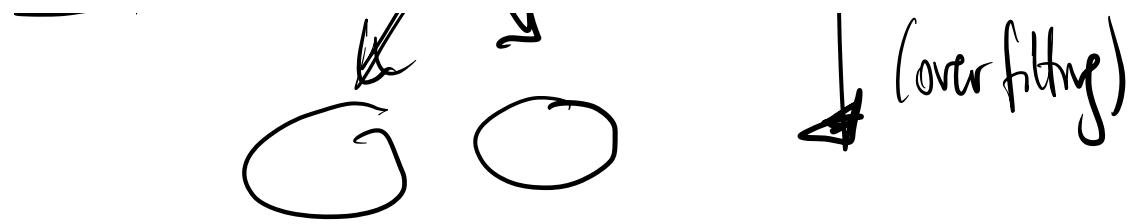
	mask	cape	tie	cars	smokes	height	class
Batman	y	y	n	y	n	180	good
Robin	y	y	n	n	n	176	u
Alfred	n	n	y	n	n	185	?
Penguin	n	n	y	n	y	140	evil
Catwoman	y	n	y	y	n	170	?
Joker	n	n	y	n	n	179	?

Butgirl	y	y	n	y	n	165	?
Riddler	y	n	n	n	n	182	?
Yordate	n	y	y	y	y	181	?

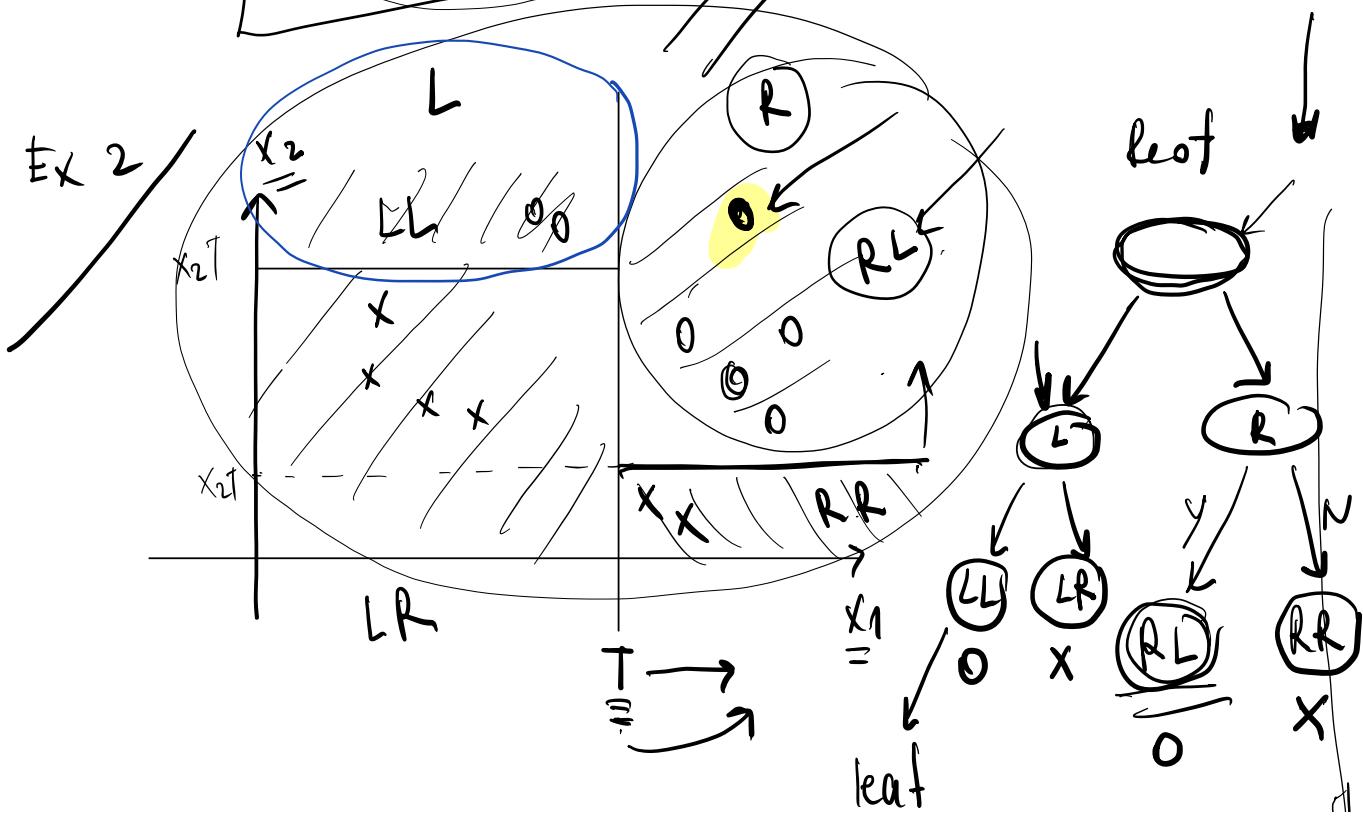
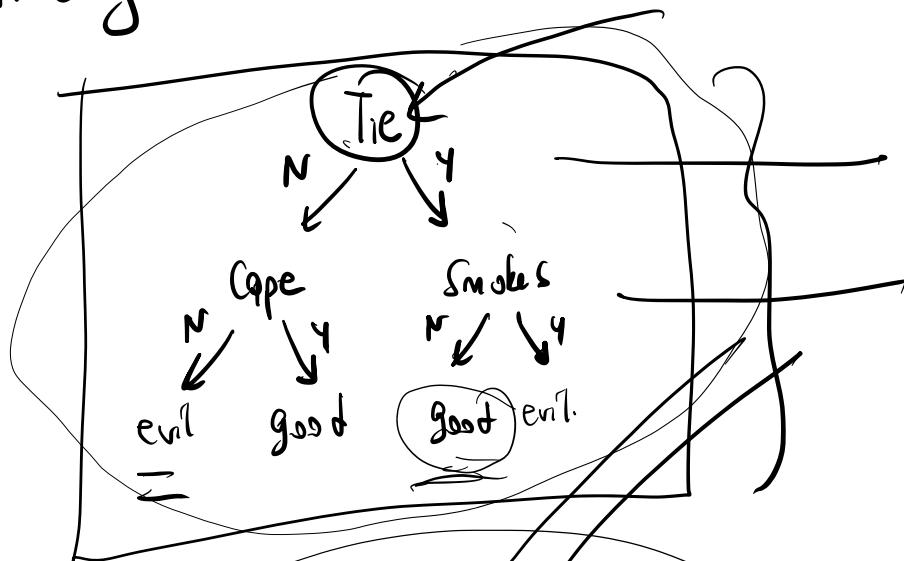


- is this a good tree?
- is there any misclassified?
- Alfred ↑





* find smallest tree that classify everything correctly.



How to split a tree node?

Data: $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$, $y_i \in \{1, \dots, c\}$

c is the # of classes.

Impurity functions: measures how pure a set is in terms of label (same label)

Gini impurity



0 (perfect purity, all elements in the subset belong to the same class).

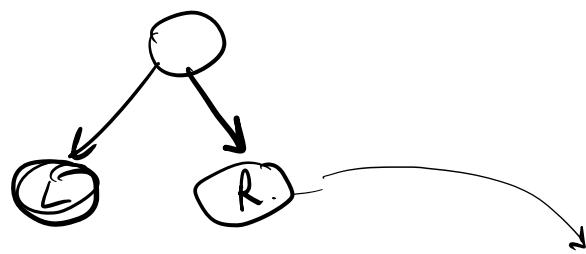
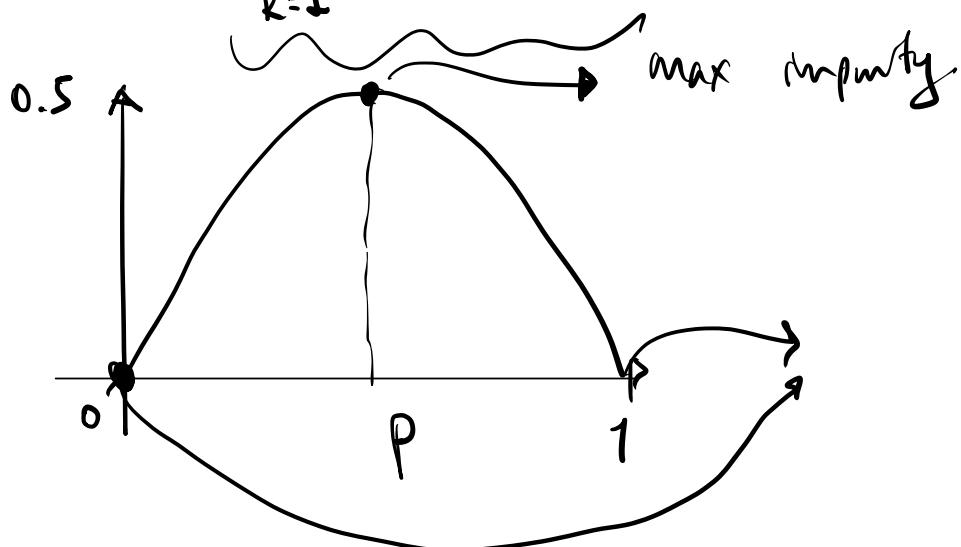
0.5 (complete impurity, where instances are evenly dist. across the classes)

Let's $S_k \subseteq S$ where $S_k = \{(x, y) \in S : y=k\}$

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

$$P_k = \frac{|S_k|}{|S|} \leftarrow \text{fraction of inputs in } S \text{ with label } k$$

$$G(S) = \sum_{k=1}^c P_k (1 - P_k)$$



$$G^T(S) = \underbrace{\frac{|S_L|}{|S|}}_{\sim} G^T(S_L) + \underbrace{\frac{|S_R|}{|S|}}_{\sim} G^T(S_R)$$

$\text{Ex/ } \textcircled{10}$ instances of class "Yes"

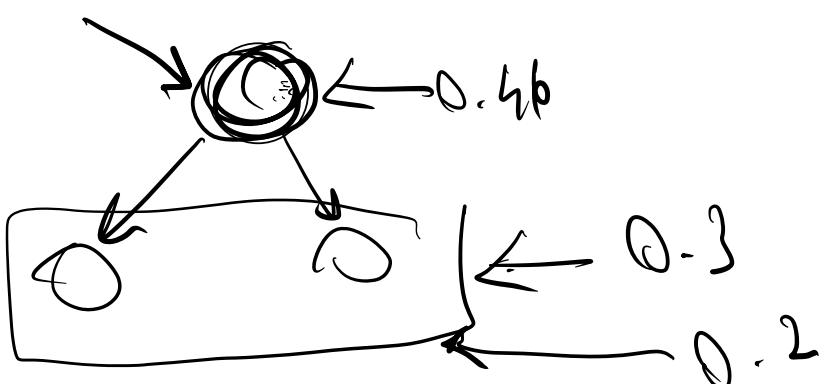
(~~6~~) u u u "No".

$$P_{\text{Yes}} = \frac{10}{10+6} = \frac{10}{16}, \quad P_{\text{No}} = \frac{6}{16}$$

$$G_{\text{mi}}(S) = 1 - (P_{\text{Yes}}^2 + P_{\text{No}}^2)$$

$$= 1 - \left(\left(\frac{10}{16} \right)^2 + \left(\cancel{\frac{6}{16}} \right)^2 \right)$$

$$\approx 0.46875$$



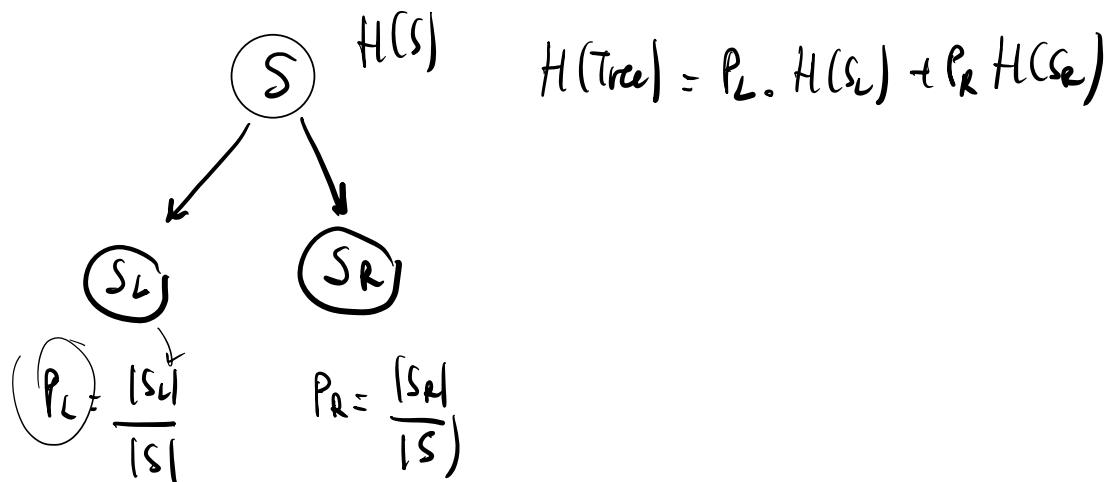
Entropy

$$\text{Entropy } (S) = - \sum_{i=1}^C p_i \log_2 (p_i)$$

Ex / 9 samples of class \rightarrow "YES"
 5 samples of a \rightarrow "NO"

$$p_{\text{Yes}} = 9/14 \quad p_{\text{No}} = 5/14$$

$$\begin{aligned}\text{Entropy } (S) &= - (p_{\text{Yes}} \log_2 (p_{\text{Yes}}) + p_{\text{No}} \log_2 (p_{\text{No}})) \\ &= - \left(\frac{9}{14} \log_2 \frac{9}{14} + \frac{5}{14} \log_2 \frac{5}{14} \right) \\ &\approx 0.939\end{aligned}$$

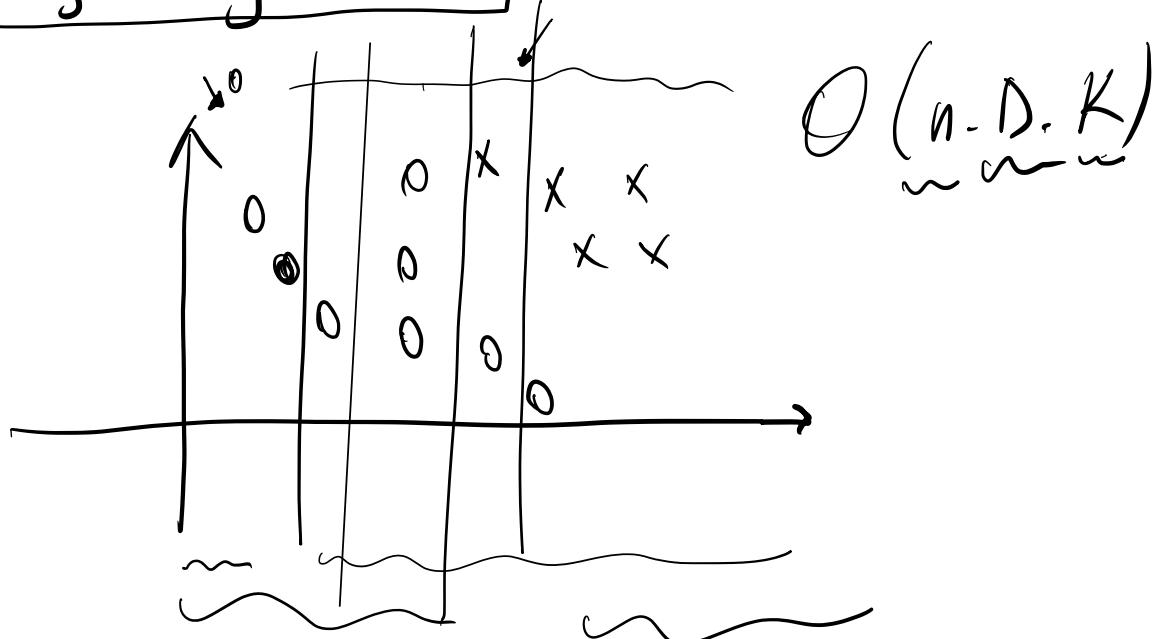


- How do you find the optimal split?

↳ NP-hard

// greedy approach ↵

→ try every feature.



Information Gain.

$$IG(S, A) = H(S) - H(S, A)$$

Adv of D.T.

$$Gini(E) = 1 - \sum_{j=1}^C p_j^2 \quad \text{or} \quad \sum_j p_j (1-p_j)$$

f_1	f_2	f_3	f_4	Labels
outlook	Temperature	Humidity	Wind	Played football
Sunny	Hot	High	Weak	No
S	H	H	Strong	N
Overcast	H	H	W	Yes
Rain	Mild	Normal	W	Y
R	Cool	C	S	Y
R	C	C	S	Z
O	M	N	W	Y
O	C	C	WW	Y
R	M	M	S	Y
S	M	H	S	Y
R	M	H	W	Y
S	H	H	S	N
R	M	H	S	Y
O	H	H	W	Y
O	M	H	S	N
R	H	H	W	Y
O	M	H	S	Y
R	H	H	W	Y
O	M	H	S	N

Σ No, η Yes

Σ Sunny, η overcast, Σ rain.

1/1

$$H(S) = - \left[\underbrace{\eta/14 \cdot \log(\eta/14)} + \underbrace{(S/14) \cdot \log(S/14)} \right]$$

$$\approx 0.94$$

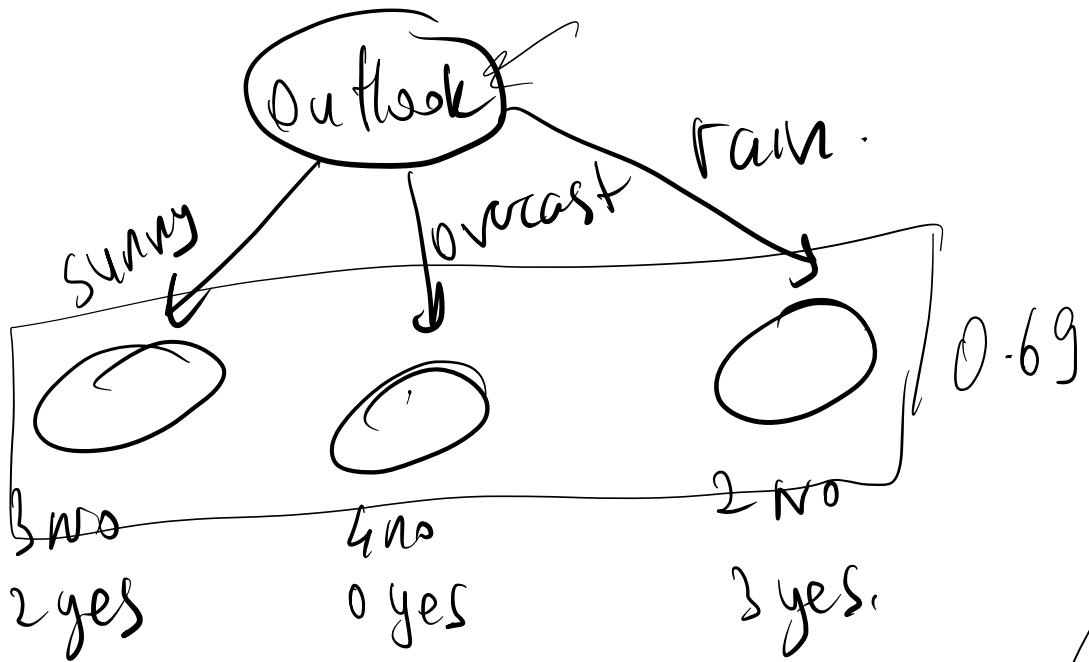
$$f_1(S, \text{outlook}) = \underbrace{S/14}_{\text{No}} \cdot \underbrace{E(3, 2)}_{\text{Yes}} + \underbrace{4/14}_{\text{No}} \cdot \underbrace{f(4/14)}_{\text{Yes}}$$

$$+ \frac{5}{14} \cdot 1 \in (2, 3)$$

~~~~~

$$\Rightarrow \frac{5}{14} \left( -\left(\frac{3}{5}\right) \log \left(\frac{3}{5}\right) - \left(\frac{2}{5}\right) \log \left(\frac{2}{5}\right) \right) + \frac{4}{14} (0) +$$

$$\frac{5}{14} \left( -\left(\frac{2}{5}\right) \log \left(\frac{2}{5}\right) - \left(\frac{3}{5}\right) \log \left(\frac{3}{5}\right) \right) = 0.693$$

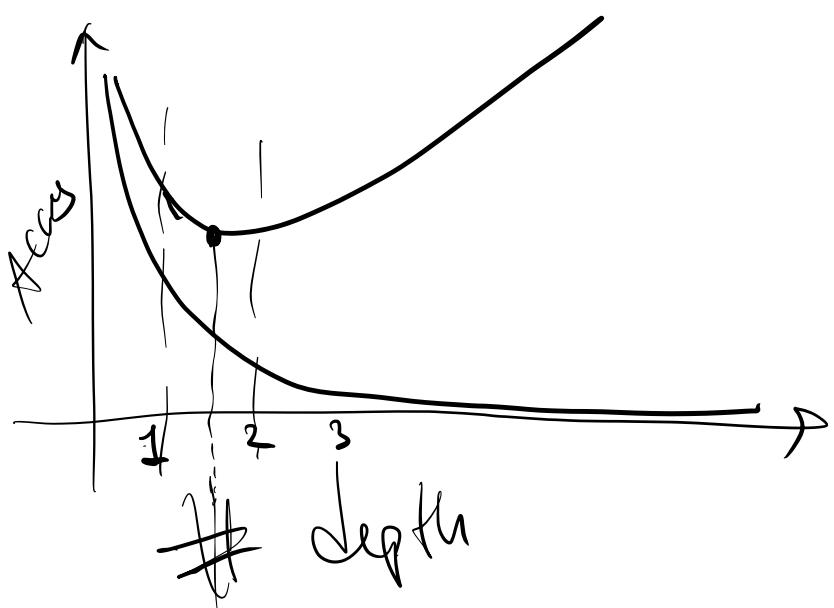
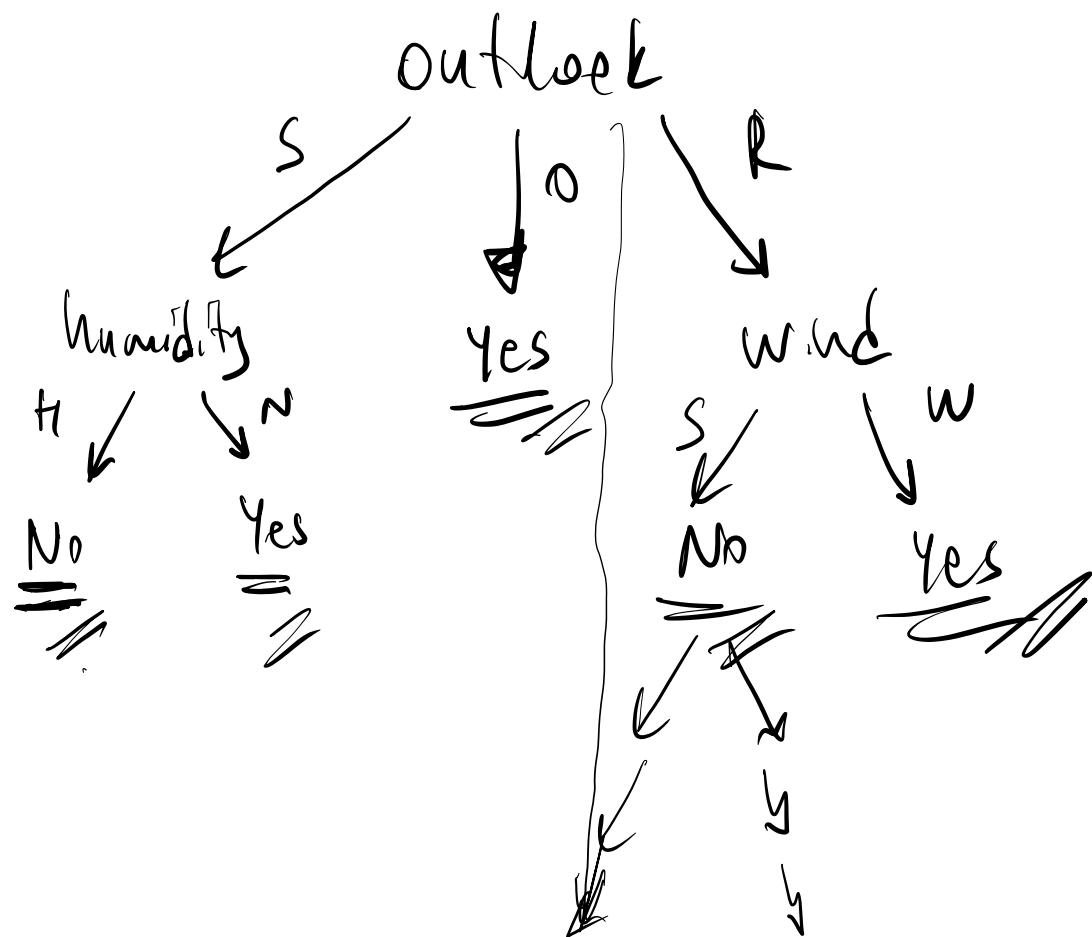


$$I_6(S, \boxed{\text{Outlook}}) = 0.94 - \boxed{0.693} = \boxed{0.247}$$

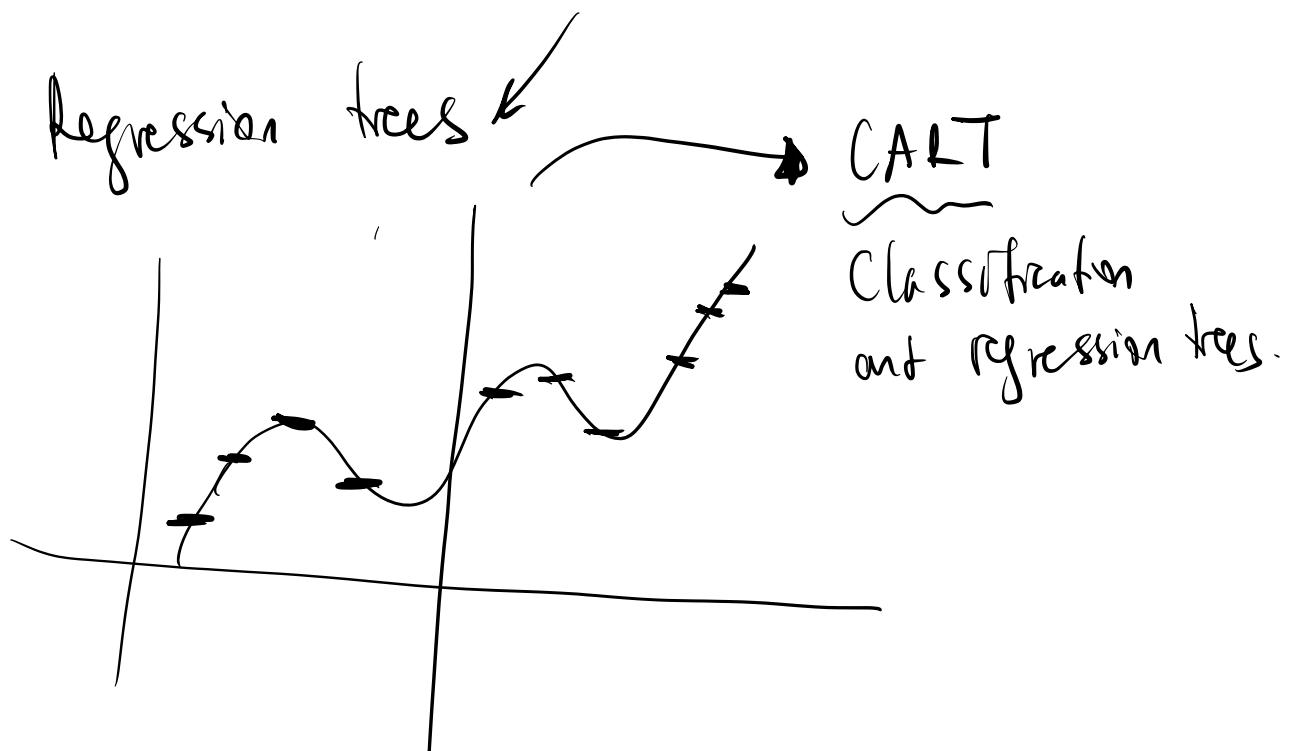
$$I_6(S, \text{temperature}) = 0.94 - \boxed{0.911} = 0.029$$

$$I_6(S, \text{humidity}) = 0.94 - \boxed{0.788} = 0.152$$

$$IG(s, \text{windy}) = [0.94] - [0.8932] = 0.048$$



$\uparrow \downarrow \rightarrow$  depth.



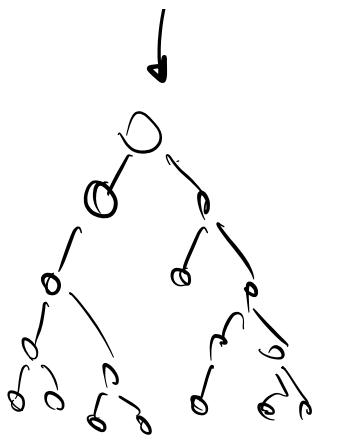
## Ensemble learning

### bagging

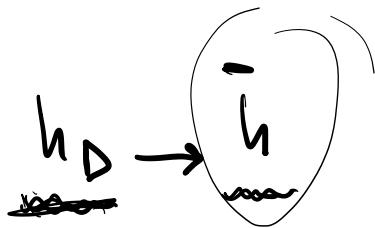
- D.T face with variance  $\sigma^2$

$$E[(h_D(x) - h(x))^2]$$

variance      outlined & expected.



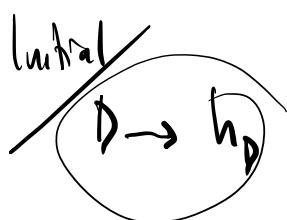
we want



## weak law of large numbers

$$\frac{1}{m} \sum_{i=1}^m x_i \rightarrow \bar{x} \text{ as } m \rightarrow \infty$$

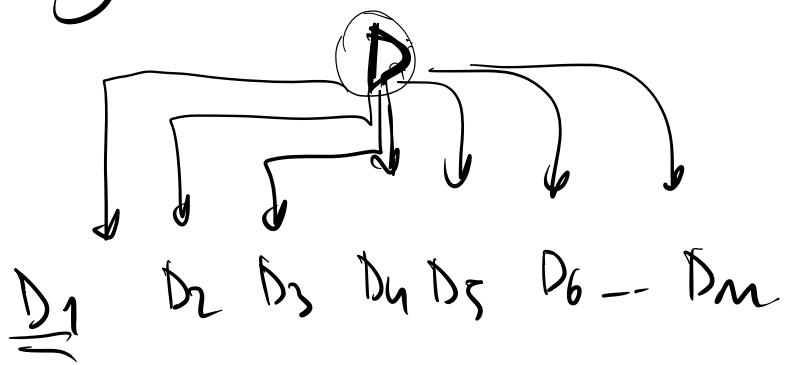
Assume we have "m" training sets  $D_1, D_2, \dots, D_m$   
 drawn from  $P^m$ .



$$\hat{h} = \frac{1}{m} \sum_{i=1}^m h_{D_i} \rightarrow \bar{h} \text{ as } m \rightarrow \infty$$

$D_i$

## Bagging (Bootstrap Aggregation).

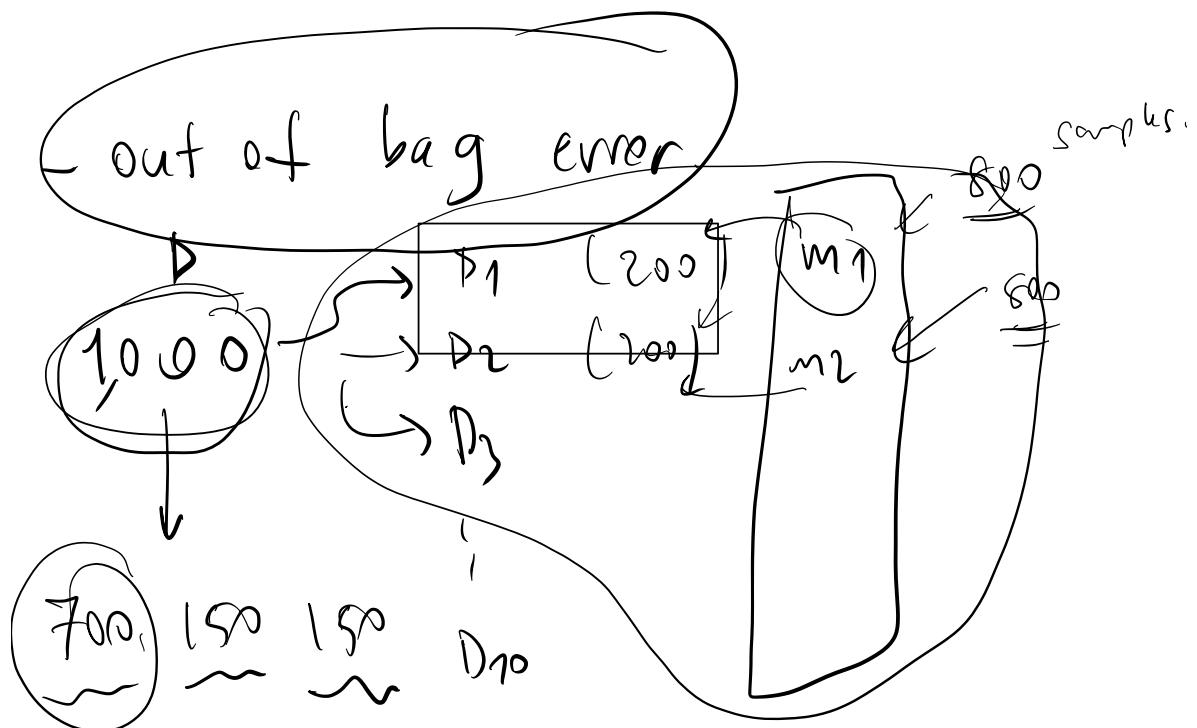


- drawing uniformly with replacement

let  $\theta(x, y | D)$  be a prob. dist. that picks  
a training sample  $(x_i, y_i)$  from  $D$  uniformly  
at random.

1. Sample  $m$  data sets  $\overbrace{D_1 \dots D_m}^{\sim}$  from  $\underline{D}$   
with replacement.
2. for each  $D_j$  train a classifier  $h_j$ .

$$3. \text{ The final classifier is } h(x) = \frac{1}{m} \sum_{j=1}^m h_j(x)$$



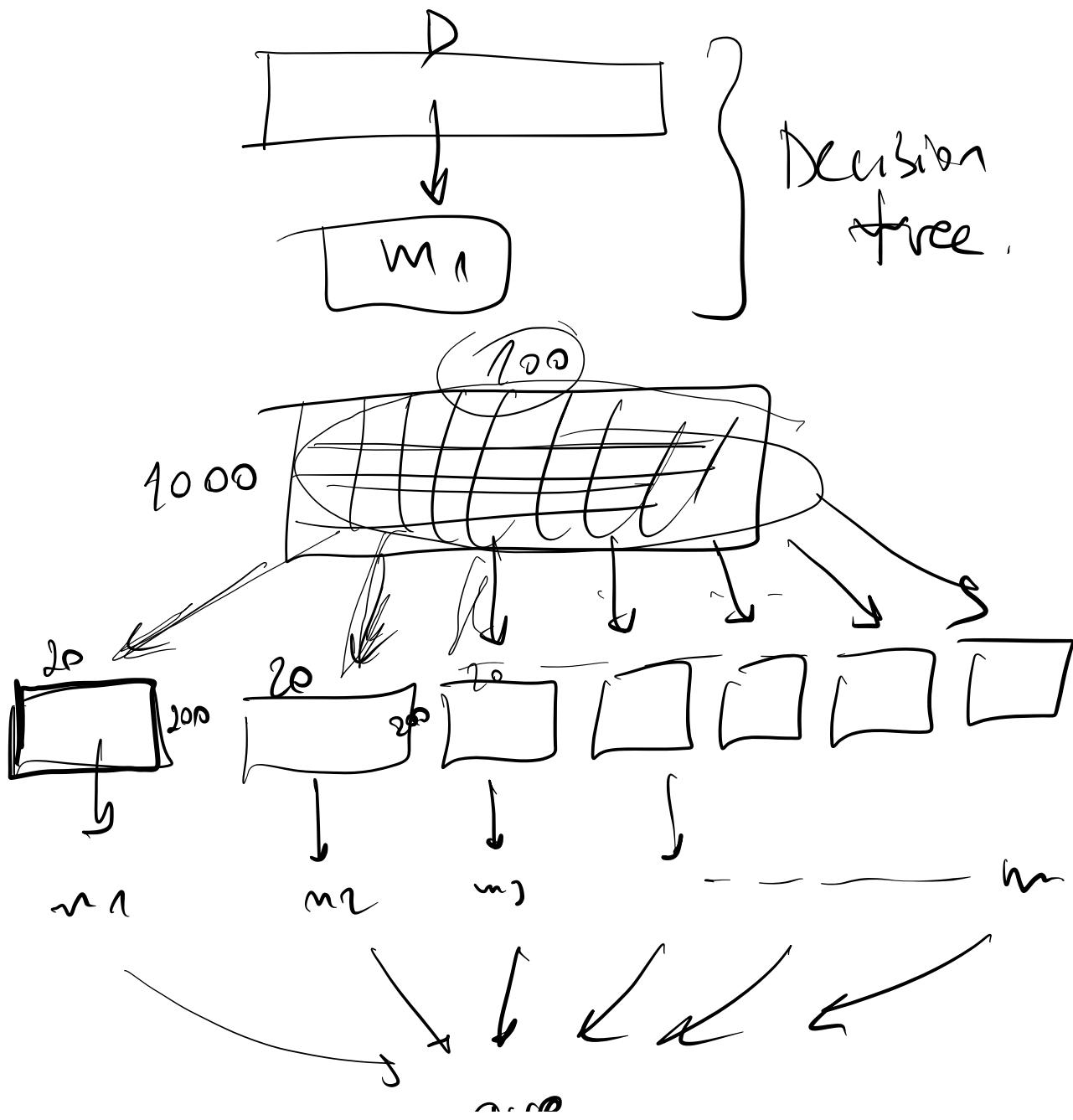
## Random forest (bagging)

original dataset

1. Sample "m" data sets  $D_1 - D_m$  from  $D$  with replacement.
2. For each  $D_j$  train a full decision tree  $h_j(\cdot)$   $\xrightarrow{\max\text{-depth} \infty}$   
with small modification: before each split randomly

Sub sample  $k \leq d$  features (without replacement) and  
only consider those for your split:

3. final classifier is  $h(x) = \sum_{f=1}^m h_f(x)$



UVJ -

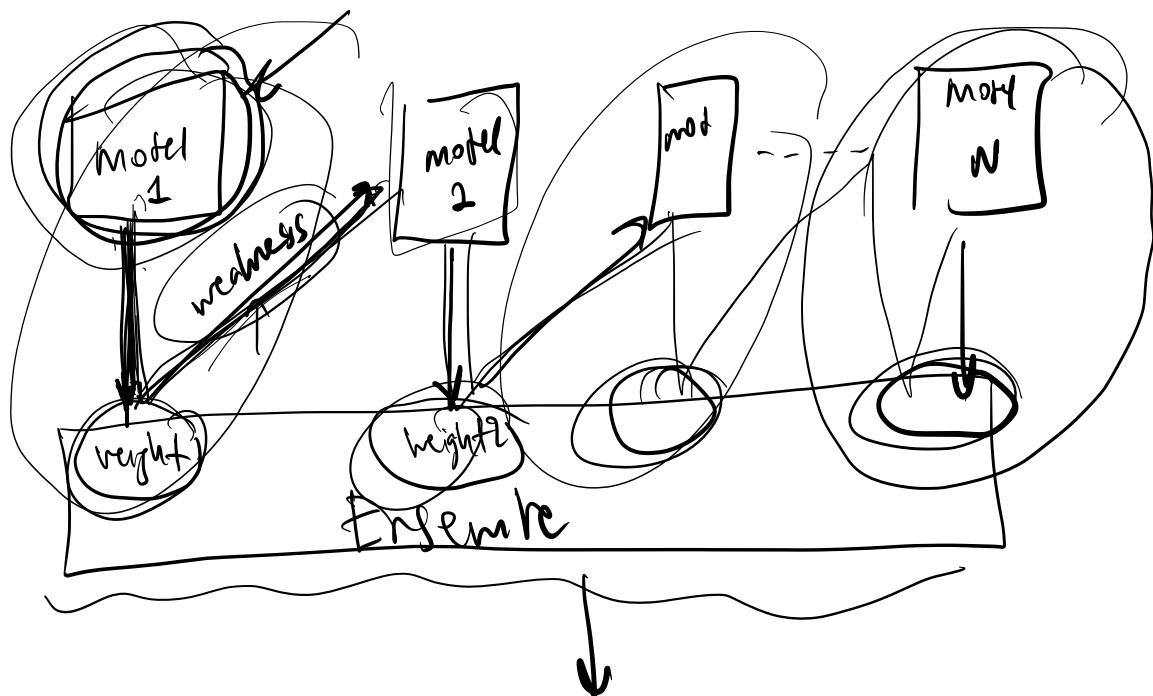
---

## Boosting.

- high bias

$$h(x) = \sum_{j=1}^m \alpha_j h_j(x).$$

→ Can we combine weak learners? (Strong learner)



Adaptive boosting: AdaBoost.

Gradient boosting:  
XGBoost

