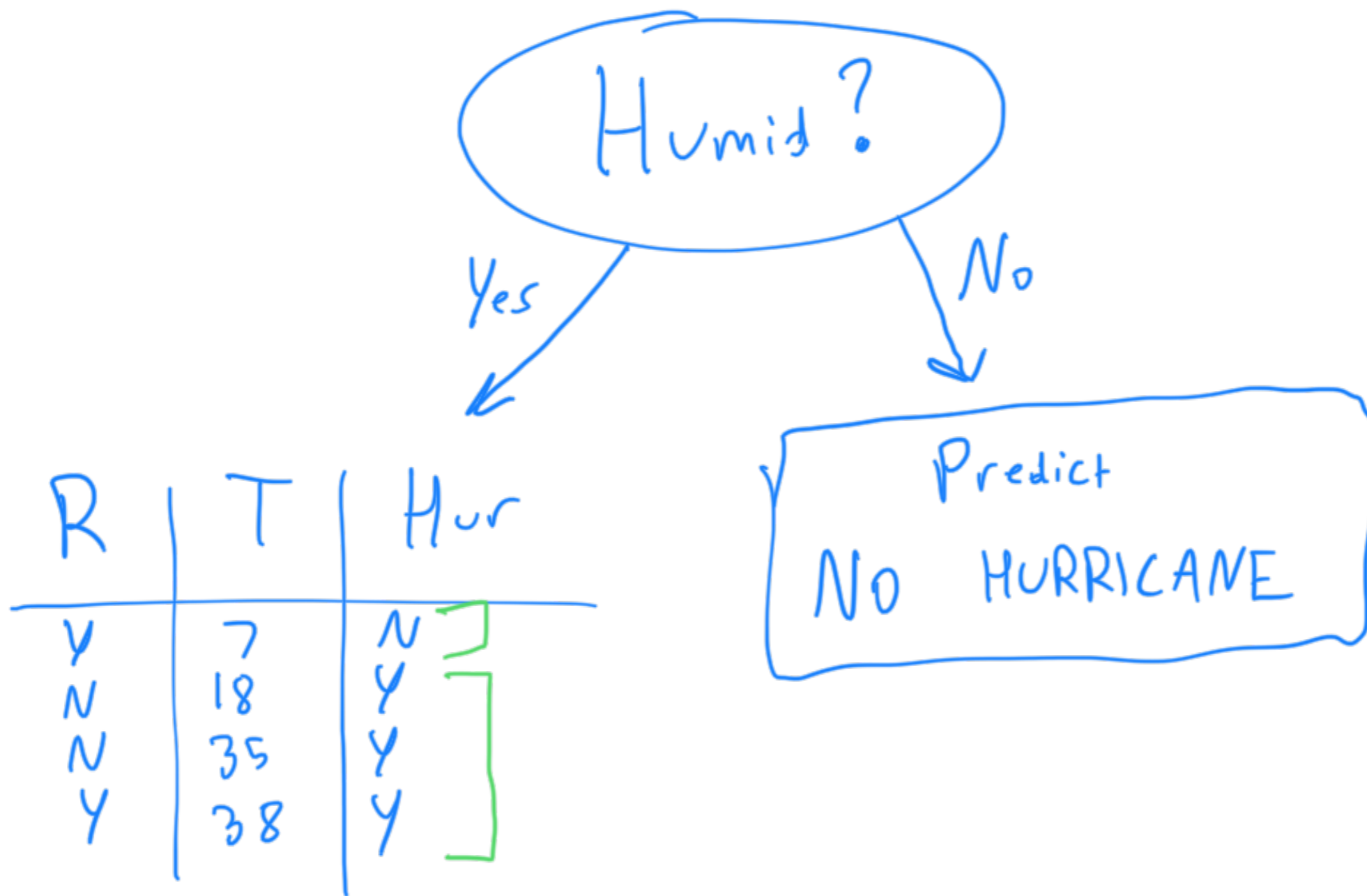
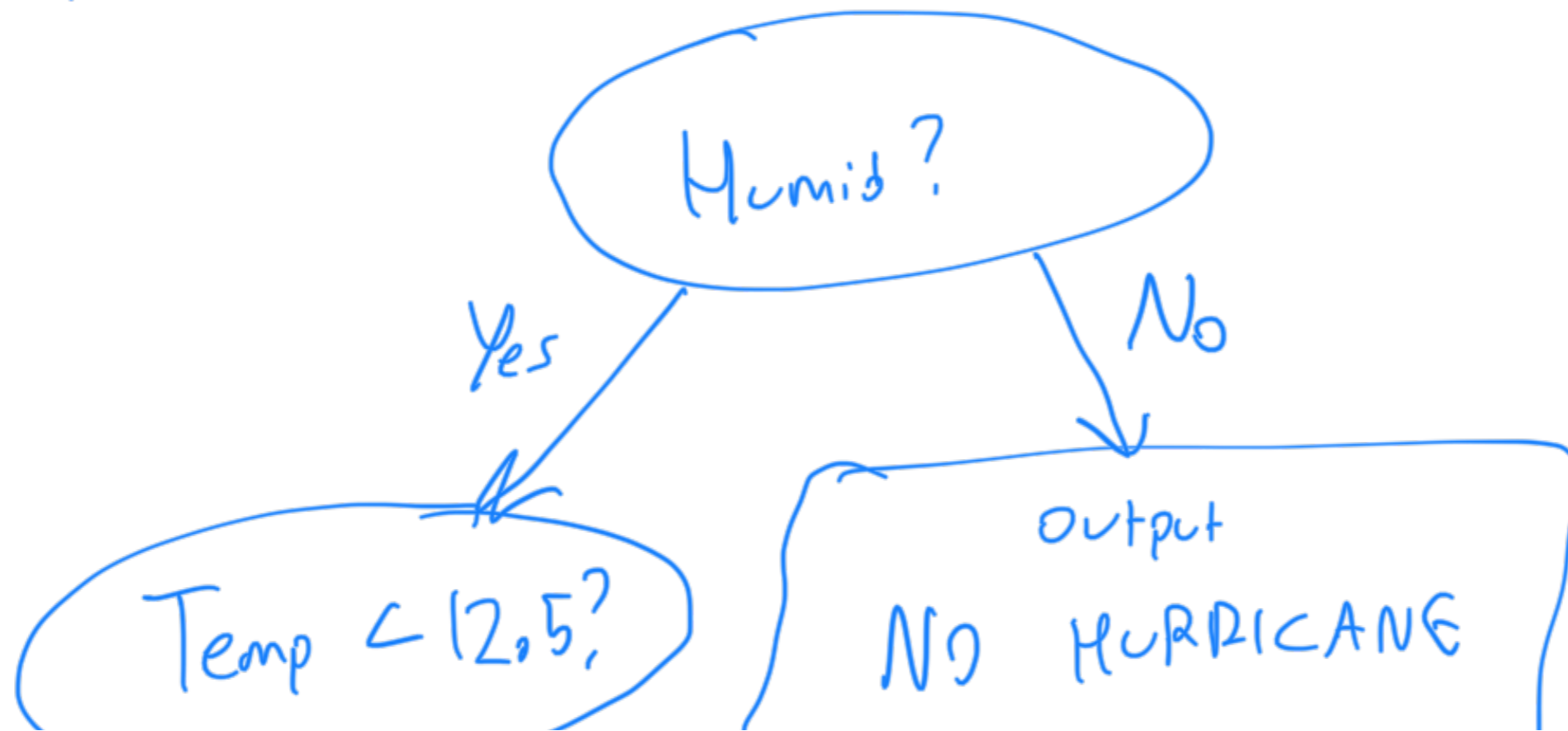
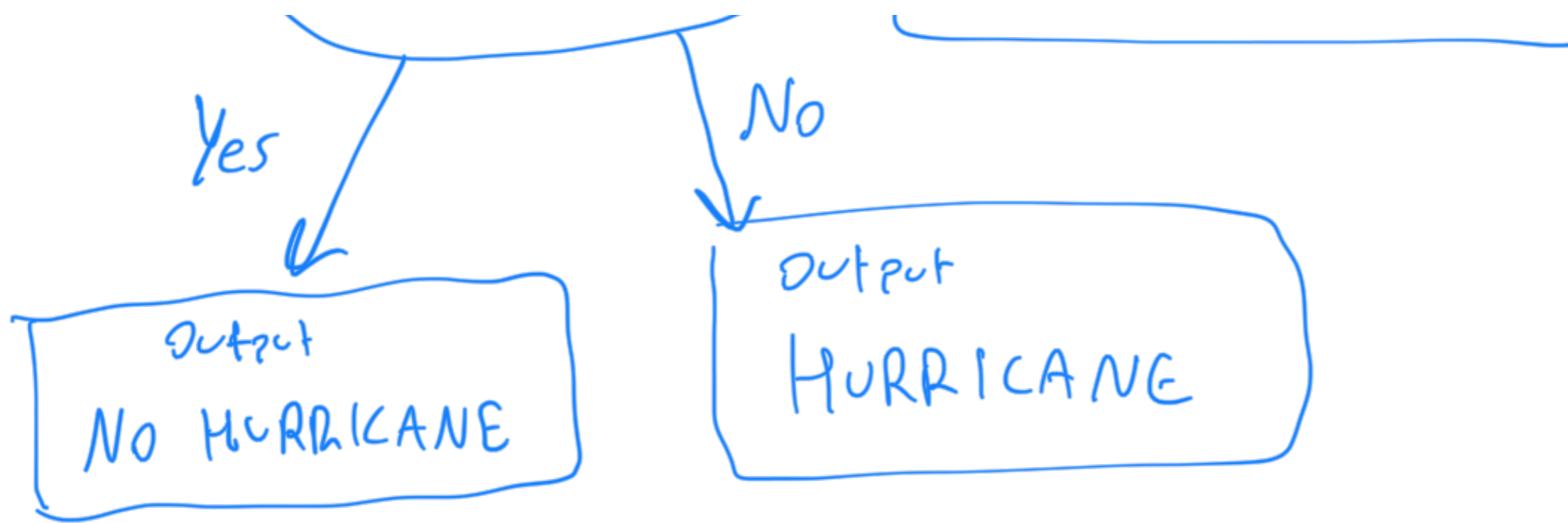


# Day 12: Ensemble Models



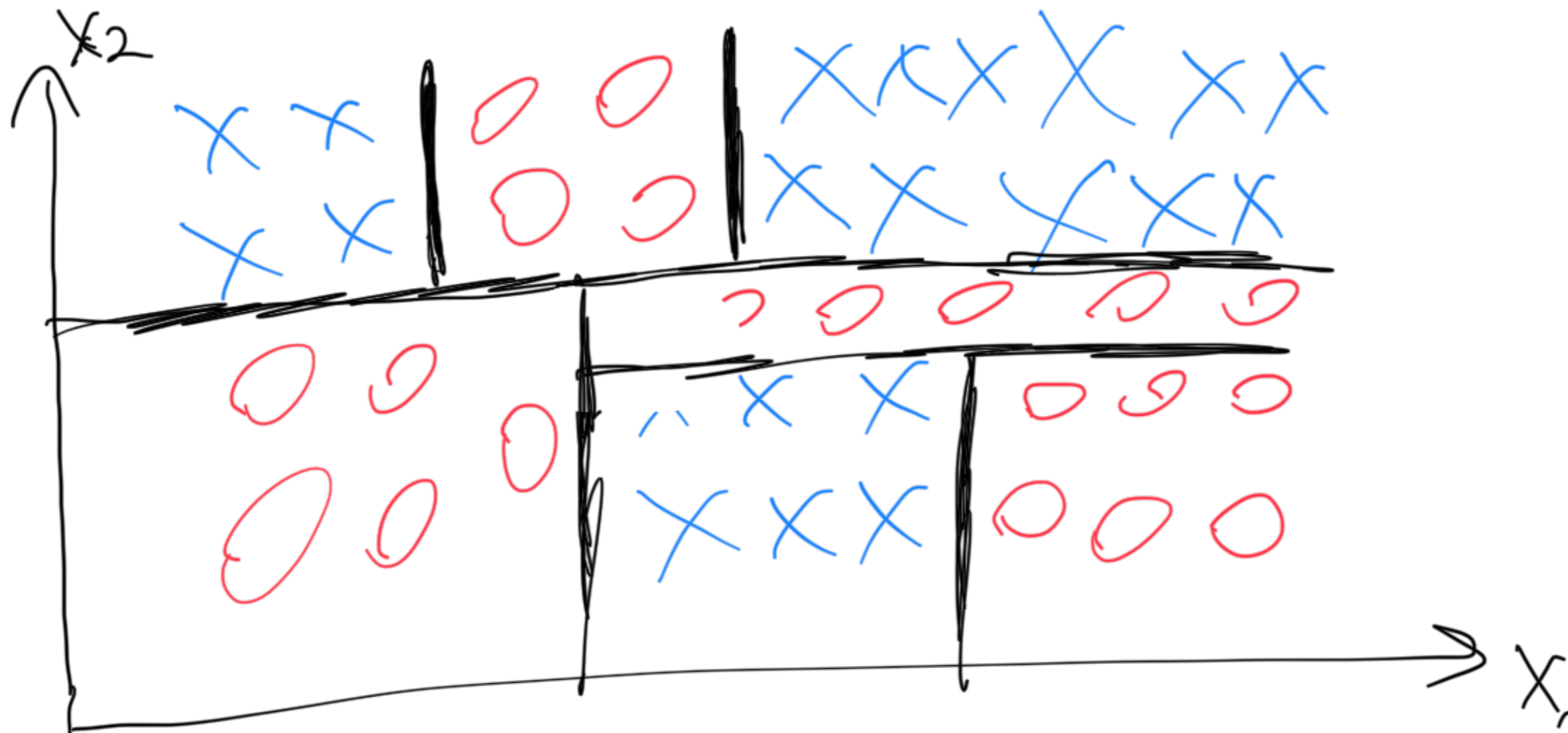
Our final decision tree is:





Decision Trees Can Learn Complex  
Non-Linear Separating Boundaries

---



# Many possible Splitting metrics

## Entropy

$$\text{Entropy}(X) = H(X) = - \sum_{i=1}^n p(x_i) \log_2 [p(x_i)]$$

$n = \# \text{ classes}$

"Disorder"

"Uncertainty"

"The number of bits required, on average, to encode information"

High entropy: high level of uncertainty

( $\underset{a}{25\%}$ ,  $\underset{b}{25\%}$ ,  $\underset{c}{25\%}$ ,  $\underset{d}{25\%}$ )

Low entropy: low level of uncertainty  
(97%, 1%, 1%, 1%)  
(split on low entropy)

## Information Gain

$$\text{Information Gain}(y, a) = \underbrace{H(y)}_{\text{Uncertainty/entropy before you split}} - \underbrace{H(y|a)}_{\text{Uncertainty/entropy after you split}}$$

"Expected reduction in entropy from splitting on an example"

Example

$$IG(\text{"Crash"}, \text{"Excess Gasoline"}) \\ = H(\text{"Crash"}) - H(\text{"Crash"} \mid \text{"Excess Gasoline"})$$

$$\underline{H(\text{"Crash"})} \\ = - \left[ \left(\frac{4}{6}\right) \log_2\left(\frac{4}{6}\right) + \left(\frac{2}{6}\right) \log_2\left(\frac{2}{6}\right) \right] \\ = 0.918$$

$$\underline{H(\text{"Crash"} \mid \text{"Excess Gasoline"})}$$

"Excess Gasoline"

$$\text{Yes: } \left(\frac{4}{6}\right) \cdot - \left[ \left(\frac{2}{4}\right) \log_2\left(\frac{2}{4}\right) + \left(\frac{2}{4}\right) \log_2\left(\frac{2}{4}\right) \right]$$

$$\text{No: } \left(\frac{2}{6}\right) \cdot - \left[ \left(\frac{2}{2}\right) \log_2\left(\frac{2}{2}\right) + \left(\frac{0}{2}\right) \log_2\left(\frac{0}{2}\right) \right]$$



→ add up the two terms:

$$H(\text{"Crash"} | \text{"Excess Gasoline"}) = 0.667$$

## Decision Trees for Regression

Similar process, except:

- \* use MSE or sum of square error instead of IG, G, or H

- \* prediction is the average of values

in a split

\* Stop when # of points in a node  
is below a threshold, or if you  
exceed some # of iterations

# Random Forests

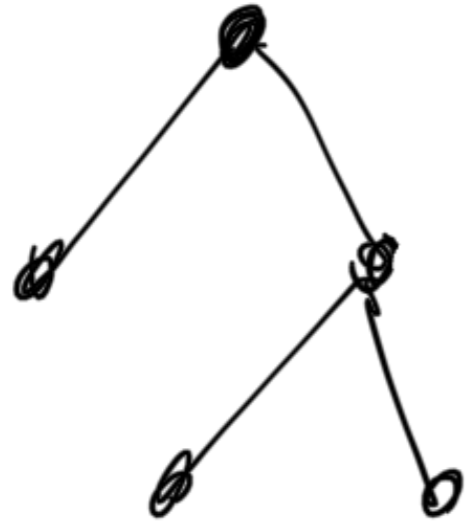
(and more broadly, Ensemble Learning)

General idea

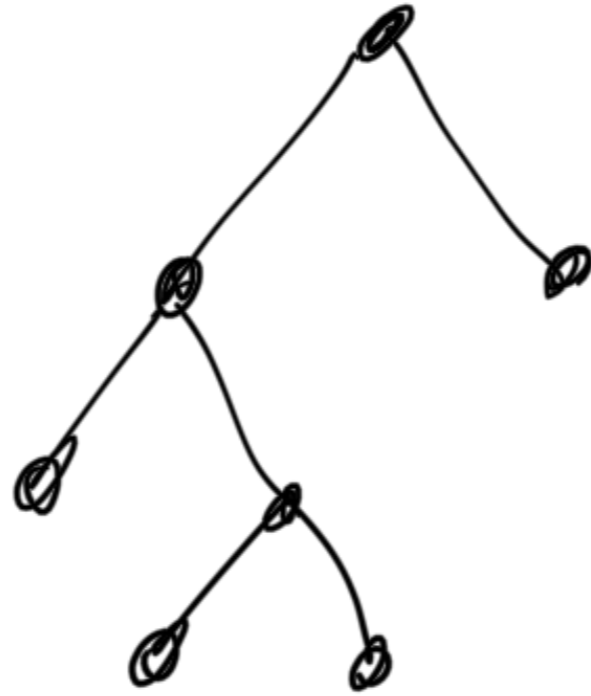
- ① Builds many decision trees\*
- ② Take the majority vote



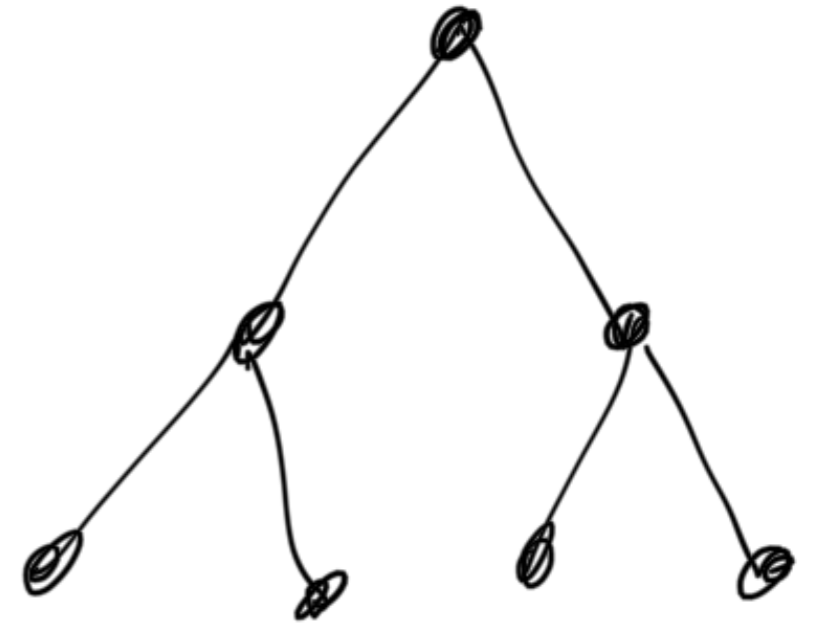
\* there are many ways to build  
Multiple decision trees



"Cancer"



"No Cancer"



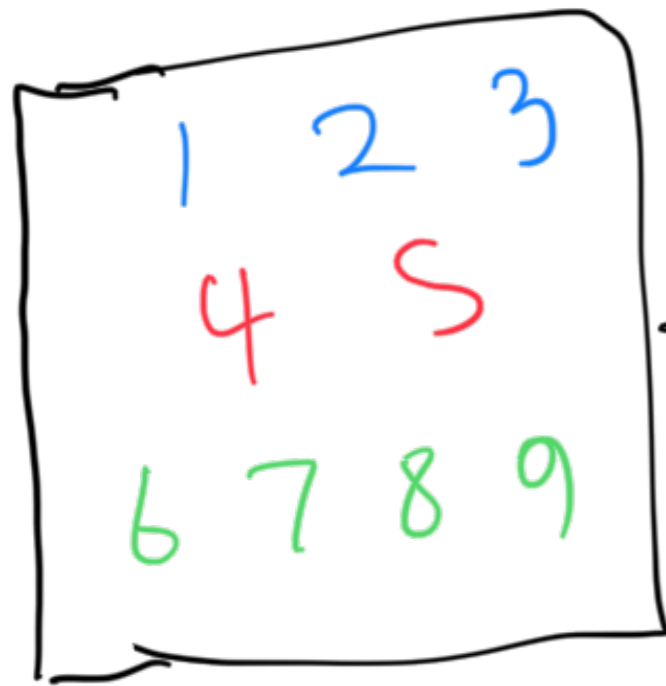
"Cancer"

predict "Cancer"

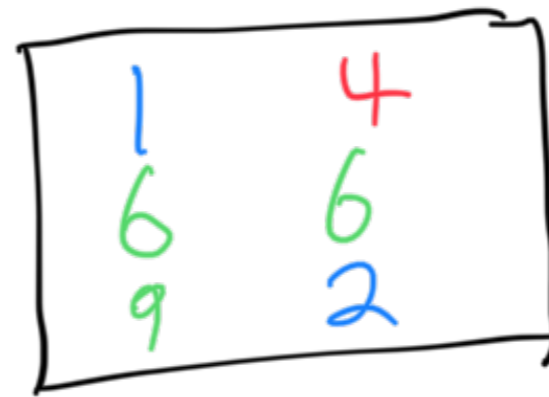
Bagging: bootstrapped aggregating

① bootstrap (sample with replacement)

dataset into  $M$  "bootstrapped datasets"

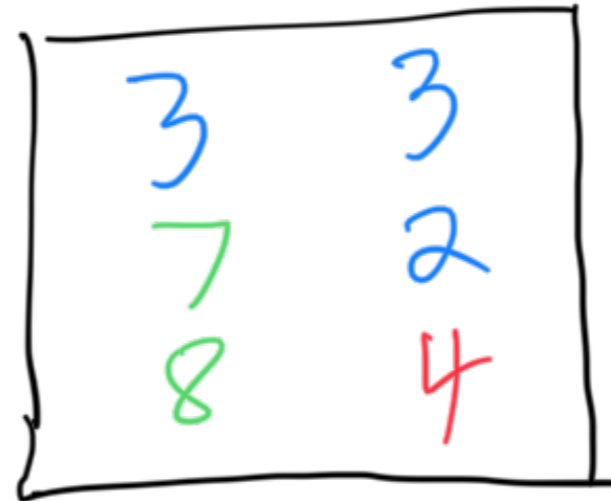


original dataset



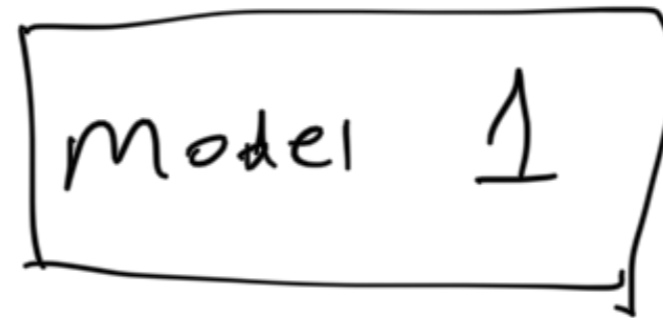
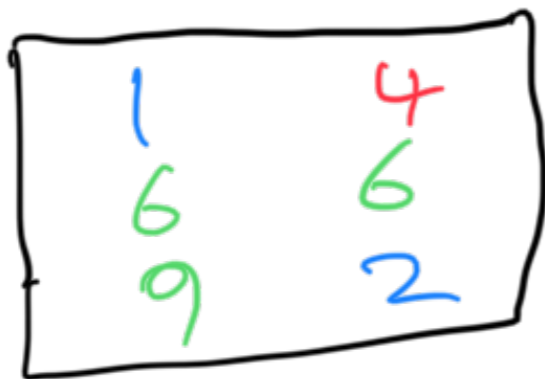
bootstrap 1

...



bootstrap  $m$

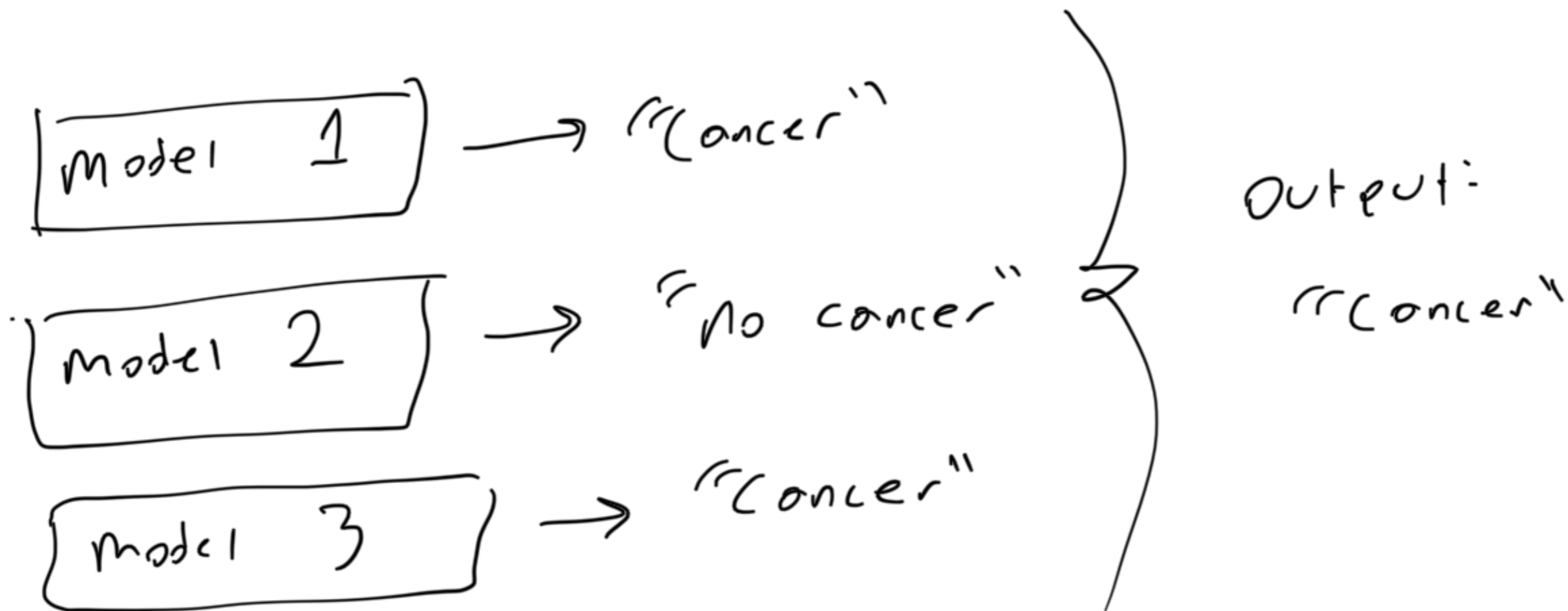
② train a model (e.g., a decision tree) on each bootstrapped sample



...



③ take majority vote of each of the models



When we apply bootstrapping to decision trees, we call that a "random forest".

But, can also bootstrap with other types of models;

logistic regression → "Cancer"  
KNN → "No cancer"  
decision tree → "Cancer"  
Naïve Bayes → "Cancer"

In addition to looking at a subset of data points, each model can also look at a subset of features.

Note: bootstrapped analyses are common in data science and regular science

Example:

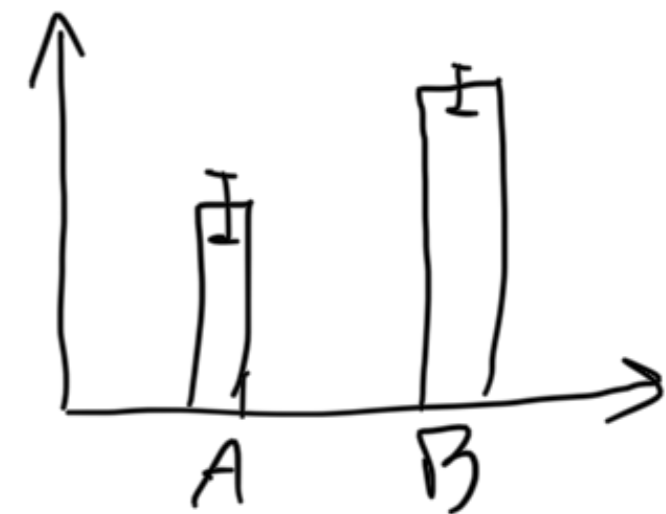
Model A	F1 Score:	0.853
Model B	F1 Score:	0.892

How do we know that this difference is not due to random chance?!!!

Instead, we can bootstrap the test set and report mean  $\pm$  standard deviation of the bootstrapped datasets

Case 1

Model A:	0.853	$\pm$ 0.0001
Model B:	0.892	$\pm$ 0.0001



Case 2

Model A:	0.853	$\pm$ 0.06	$\uparrow$	T	$\pm$
----------	-------	------------	------------	---	-------



Model B: 0.892 +/- 0.01



Boosting: iteratively fit models which prioritize misclassified data points in the previous iteration

Ada Boost (one boosting algorithm)

① initialize "weights" of each data point to be equal

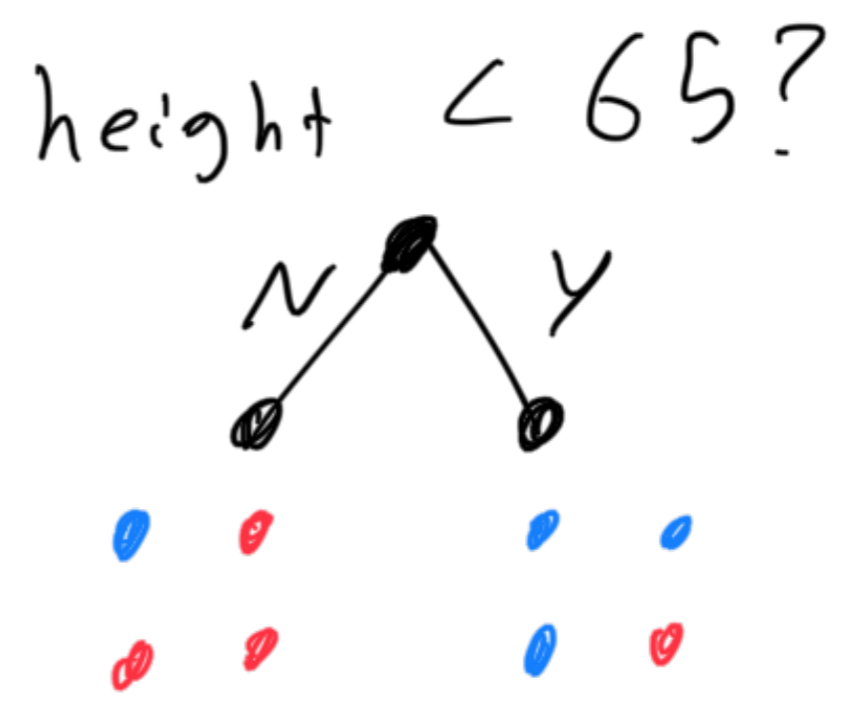
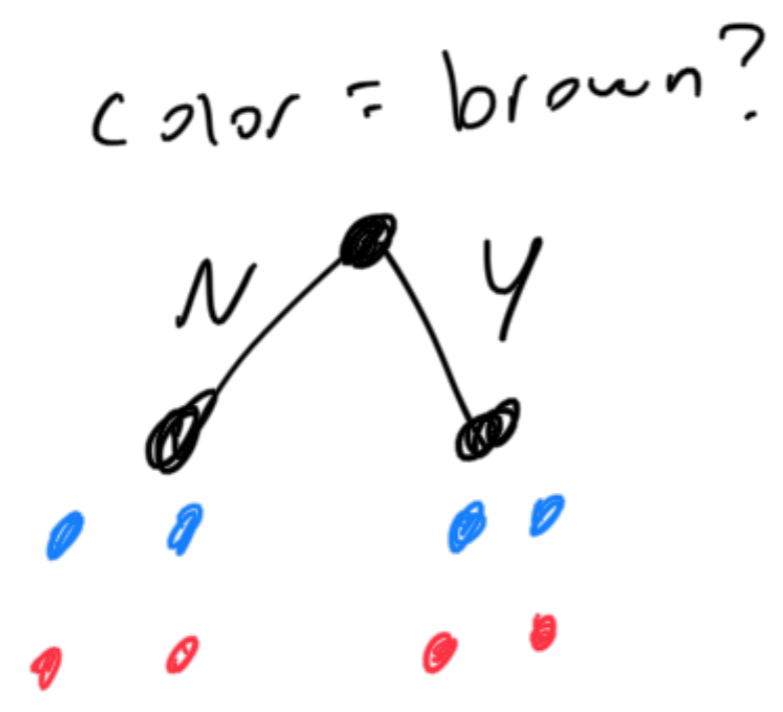
F F F F M M M M



weight:  $\frac{1}{8}$   $\frac{1}{8}$   $\frac{1}{8}$   $\frac{1}{8}$   $\frac{1}{8}$   $\frac{1}{8}$   $\frac{1}{8}$   $\frac{1}{8}$

(2) For  $N$  rounds;

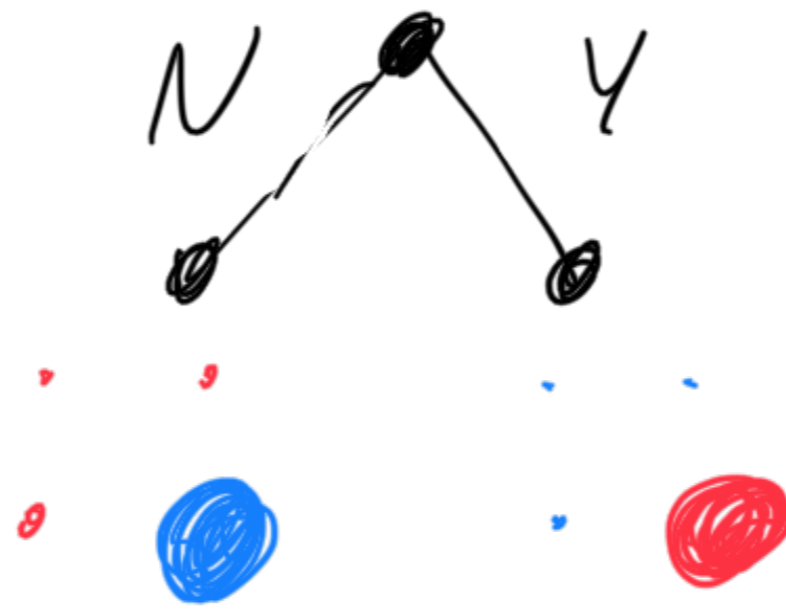
(a) train a model per feature which only uses that feature



(b) Choose model with lowest weighted error, weighted by the weights of the data points

Select 'height'  $< 65$

② Update weights of each data point:  
increase weight if incorrectly predicted,  
decrease weight otherwise



③ Final model is combination of  
models from rounds  $1$  through  $N$ ,  
weighted based on the error  
of that round

weight: 0.2      0.3      0.6      0.1      0.8



So, the final prediction is F