

①

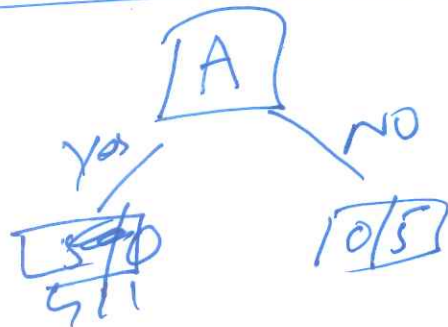
Decision Tree

$$L_{ini} = 1 - \sum p_i^2$$

Target  $\frac{Y}{N} = \frac{6}{10} = 0.6$   
 $\frac{N}{N} = \frac{4}{10} = 0.4$

$$L_{ini_{root}} = 1 - ((0.6)^2 + (0.4)^2)$$

$$= 1 - (0.16 + 0.36) = 1 - 0.52 = 0.48$$



$$L_{in_A (yes)} = 1 - \left( \left( \frac{4}{5} \right)^2 + \left( \frac{0}{5} \right)^2 \right)$$

$$= 1 - (0.64 + 0.32)$$

$$= 0.36$$

$$L_{in_A (no)} = 1 - \left( \left( \frac{0}{5} \right)^2 + \left( \frac{5}{5} \right)^2 \right)$$

$$= 1 - (1) = 0$$

$$\text{Weighted } L_{in_A} = \frac{15}{10} \times 0.32 + \frac{5}{10} \times 0$$

$$= 0.16$$

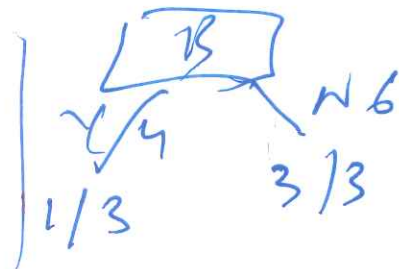
2

$$\text{Crim}_B(\text{Yes}) = 1 - \left( \left( \frac{1}{4} \right)^2 + \left( \frac{3}{4} \right)^2 \right)$$

$$= 1 - (0.0625 + 0.5625)$$

$$= 1 - 0.625$$

$$= 0.375$$



$$\text{Crim}_B(\text{No}) = 1 - \left( \left( \frac{3}{6} \right)^2 + \left( \frac{3}{6} \right)^2 \right)$$

$$= 1 - (0.25 + 0.25)$$

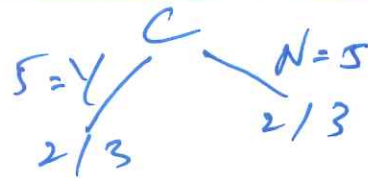
$$= 0.5$$

$$\text{Weighted Crim}_B = \frac{4}{10} \times 0.375 + \frac{6}{10} \times 0.5$$

$$= 0.45$$

$$\text{Crim}_C(Y) = 1 - \left( \left( \frac{2}{5} \right)^2 + \left( \frac{3}{5} \right)^2 \right)$$

$$= 0.48$$



$$\text{Crim}_C(N) = 1 - \left( \left( \frac{2}{5} \right)^2 + \left( \frac{3}{5} \right)^2 \right)$$

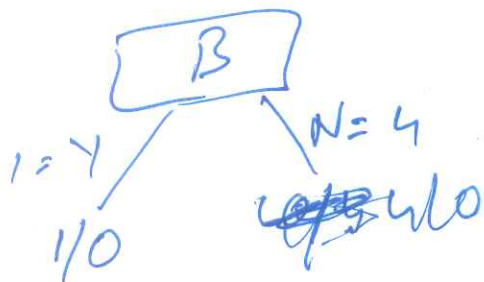
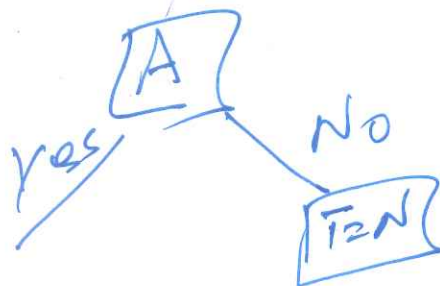
$$= 0.48$$

$$\text{Weighted Crim}_C = \frac{5}{10} \times 0.48 + \frac{5}{10} \times 0.48$$

$$= 0.48$$

Crim A has best split — lowest error

4



$$L_{inB}(1) = 1 - \left(\frac{1}{1}\right)^2 + \left(\frac{0}{1}\right)^2$$

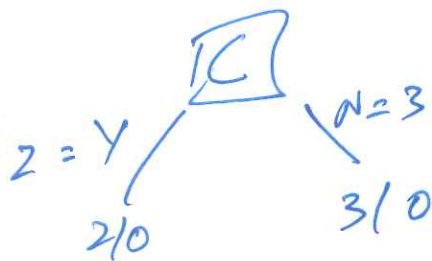
$$= 0$$

$$L_{inB}(N) = 1 - \left(\frac{0}{4}\right)^2 + \left(\frac{4}{4}\right)^2$$

$$= 0$$

$$Weight_{linB} = \frac{1}{5} \times 0 + \frac{4}{5} \times 0$$

$$= 0$$

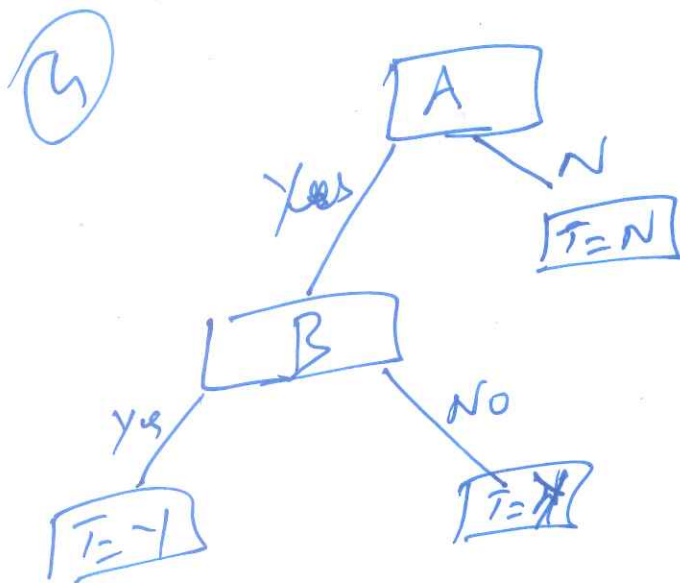


$$L_{inC}(1) = 0$$

$$L_{inC}(N) = 0$$

$$Weight_{linC} = 0$$

We can select B or C.



- 
- easy to build, easy to use, easy to interpret.
  - good to create with data
    - ↳ not good to predict

Overfit — memorize the training data

Sensitive → small change = diff. tree

↳ instability.

Imbalanced data = bias towards dominant class

# Random Forest

Bootstrapping = Bagging.

More than one trees

↳ Forest.

↳ same tree as the same data?

---

Bootstrapped data set

↳ ~~Allow~~ Select sample from the <sup>original</sup> data  
+ Allow to repeat. | Bag

↳ some data may not be part of  
the Bag — out of bag data

default for classification

↳ Select a few features (randomly)

↳  $\sqrt{n}$  or 2

↳ different trees with different features.



# How to classify

→ for any new data

↳ run through all DTs.

↳ Voting: the class ~~with most~~ in the result of most DTs is the winner

→ how to measure the model

↳ Out of Bag samples = data not selected

↳ run these samples on DTs.

↳ we have label and predicted labels.

