

Entropy:

02/14/2023

$$-\sum_i \text{prob}_i \log(\text{prob}_i)$$

$$\text{prob}_i = \frac{N(i, D)}{N(D)}$$

$H(D)$

Gini:

$$\sum_i \text{prob}_i (1 - \text{prob}_i)$$



information gain

$$H(P) - \left[\frac{N(P_L)}{N(P)} H(P_L) + \frac{N(P_R)}{N(P)} H(P_R) \right]$$

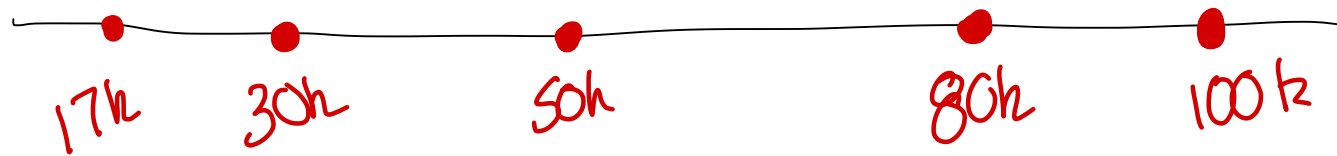
impurity of split:

$$\frac{N(P_L)}{N(P)} H(P_L) + \frac{N(P_R)}{N(P)} H(P_R)$$

Choose threshold:

1) consider threshold
data point
to consider)

value. between each
($N-1$ thresholds



2) large # of continuous values

↳ randomly choose subset of data
↳ consider thresholds as midway points
of data in subset

employed

0

unemployed

1

self-employed

2

red T

green T

blue T

yellow H

purple T

orange H

color ?

red, yellow
purple

green, blue
orange

weight

length

crocodile or snake?

2200

18

crocodile

300

17

snake

500

2200

18

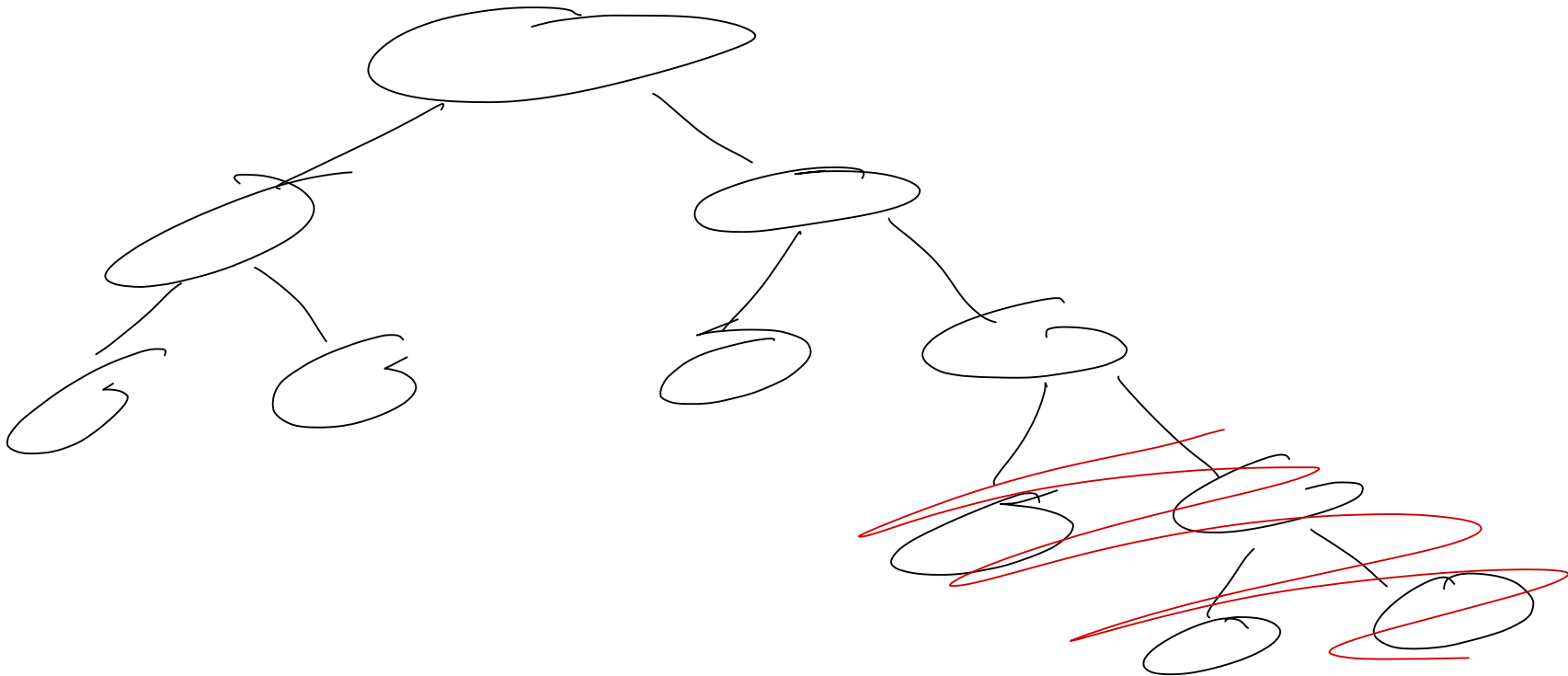
snake

When to stop:

1) if pool is too small

2) if all items in a pool
have same class

3) if depth has reached limit



Random Forest

1) at leaf, record a vote

↳ choose class with most votes

2) at leaf, record N_e votes for each label that occurs

N_e : # of times label appears in data at leaf

↳ choose label with votes