## * Entropy

$$P = [P_1, P_2 \ldots P_k]$$
$$P_i \geqslant 0$$
$$\Sigma P_i = 1$$

$$\text{entropy}(P) = -\sum_{i=1}^{k} P_i \cdot \log_k P_i$$

↳ non negative number

Single entry with P=1
$$P_j = 1$$
$$0$$

$$[\tfrac{1}{k} \ldots \tfrac{1}{k}]$$
1 evenly distribute

14 Points

$$P = \left[\frac{9}{14}, \frac{5}{14}\right]$$

Yes's     No's

entropy of a feature
Outlook

Sunny → how many have Yes and No

5      2    3

$$\left(\frac{2}{5}, \frac{3}{5}\right)$$

→ Entropy($D$ | Outlook = Sunny) = $E(p_1)$
→ Entropy($D$ | Outlook = Rain) = $E(p_2)$
→ Entropy($D$ | Outlook = Overcast) = $E(p_3)$

avg entropy

$$E(\text{outlook}) = \frac{5}{14} E(p_1) + \frac{5}{14} E(p_2) + \frac{4}{14} E(p_3)$$

---

all features.

E [outlook]      0-7
E [temperature]   0·2 → smallest → becomes first question.
E [humidity]     0·3
E [wind.]      0·6

Hyperparameter

Entropy (P)    $P = [P_1 \ldots P_u]$

$\hookrightarrow$ gini index(p)

① gini (p) = $1 - \sum_{i=1}^{n} P_i^2$ $\Big]_0^{1 - \frac{k}{k^2}}$ } max value } default

whenever we have high gini index)

It's a good feature

② when to stop asking question?
Good to ask not too many questions.

10k days Large
$\downarrow$
10 days $\rightarrow$ Overcast (very small subset)
$\downarrow$
doesn't have much value
$\downarrow$
if we ask questions about it
doesn't make sense

$\rightarrow$ Small SubDataset –Size of min dataset

$\rightarrow$ Depth of the tree (we can restrict it)
if reach dept x $\rightarrow$ stop

# *Regression

Calculate mean of numerical →

Variance/MSE → $\sum_{j=1}^{n} (y_j - \bar{y})^2 / n$

MSE (D) = 0 | $y_j = \bar{y}$

→ avg

↙ ↓

entire Dataset    labels

(how good?)

labels in output - same

↓

Stop asking QS


* Trees on numerical features

$-100^{0} \ldots \ldots 100^{10}$ | human intuition

bucketisation | automatic

⤵ → hyperparameter
              or
              decided by algorithm

num feature → k binary feature

i: $X_i$ — [0, .... ]

⤵ ⤵

0·33    0·66

divide into sub interval

thresholds $\begin{cases} (X_i > 0·33) \\ (X_i > 0·66) \end{cases}$ if 0·5 which is numerical

example ↓

convert to 1

if $X_i > 0·33$ ⤵ assign

Example → fully numerical.

3.0 —
2.5 —
2.0 —
1.5 —
1 —
0.2   0.4   0.6   0.8   1.0

y → numerical

x → numerical

X > 0.25
X > 0.5  ⟶  X > 0.5 → first question
0 > 0.75

0.25 ———————————— High Variance
0.5 —/———————— → Great → first question
                                    (not high variance)
0.75 ————/———— High Variance

          X > 0.5
      no  /      \  yes
    X > 0.25      X > 0.75
  no /    \ yes   no /   \ yes

┌─────────────┐
│ avg of y's  │
│ in 0.25     │
└─────────────┘

## Random Forests    forest (multiple decision trees)

## TRAINING

Subdataset

$$X \xrightarrow{\hspace{1cm}} \begin{array}{c} X_L \\ n' \times d' \\ \vdots \\ X_k \end{array}$$

$n \times d$

→ reselect row (repetition)

(random selected n' with replacement (able to reselect)
↓
select d' attributes (randomly not replacement)

$$\text{Train} \begin{array}{c} T_L \\ \vdots \\ T_n \end{array} \Big\} \text{over datasets} \begin{array}{c} X_1 \\ X_2 \\ \vdots \\ X_k \end{array} \Big\}$$

## EVALUATION / PREDICTION

$$X \rightarrow \begin{array}{c} T_1(x) \\ \vdots \\ \vdots \\ T_n(x) \end{array} \qquad \begin{array}{c} y_1 \; [0,1,2] \\ \vdots \\ \vdots \\ y_k \end{array} \Big\}$$

majority
which label
comes more
often

k diff labels

↘ output

Early stopping questions $\Big]$ 10 decisio trees → 10 protocols
↓
overfitting will
cancel out.