# Decision Tree

Versinca

o  o  o    o

o o o  o ——— hyperplane

o  o    Versica

Ⅴ

$PL < 0$  — Root node

Yes          No ——— splitting

Setosa

$SL < 1.5$  ——— Branch /subtree

Yes        No

Versicolor    Virginica  → Pecision node

→ leaf node

$\sum P_i \, \log$

## Entropy

measure of disorder

$$E(S) = \sum_{i=1}^{c} P_i \log_2 P_i$$

$$E(D) = -P_{yes} \log_2 (P_{yes}) - P_{no} \log_2 (P_{no})$$

| | Salary | Age | Purchase |
|---|---|---|---|
| 1 | 20000 | 21 | Yes |
| 2 | 10000 | 45 | No |
| 3 | 60000 | 27 | Yes |
| 4 | 150000 | 31 | No |
| 5 | 120000 | 18 | No |

$$H(d) = -P_y \log_2 (P_y) - P_n \log_2 (P_n)$$
$$= -\tfrac{2}{5} \log_2 (2/5) - 3/5 \log_2 (3/5)$$

$$H(d) = 0.97$$

3 terms

$$-P_x \log_2 P_x - P_y \log_2 P_y - P_z \log_2 P_z$$

- More the uncertanty more is entropy
- for a 2 class problem the min entropy is 0 and the max is 1

- for more than 2

    min → 0

    max can be greater than 1

- Both $\log_2$ or $\log_e$ can be used to calculate entropy

# Information Gain

### Decrease in Entropy

$$IG = E(Parent) - \{Weighted\ Avg\}^o_{e\ \{children\}}$$
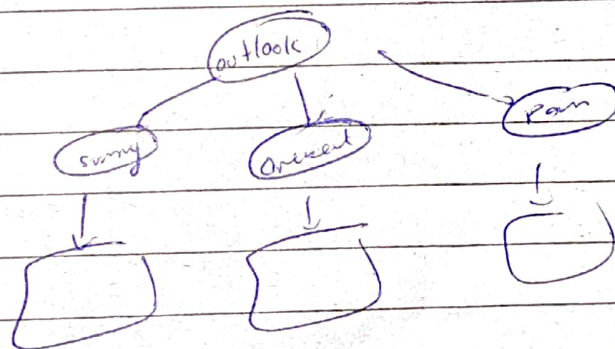
## Play Tennis dataset

### Entropy of Parent

$$E(P) = -P_y \log_2 P_y - P_n \log_2 P_n$$

$$= -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

$$= 0.94$$



$$E(S) = -\frac{2}{5} \log \frac{2}{5} - \frac{3}{5} \log \frac{3}{5}$$

$$= 0.97$$

$$E(0) = -\frac{4}{4} \log \left(\frac{4}{4}\right)$$

$$E(0) = 0$$

$$E(R) = 0.97$$

### Calculate weighted Entropy

$$W.E = \frac{5}{14} * 0.97 + \frac{4}{14} * 0 + \frac{5}{14} * 0.97$$

weight

$$W.E = 0.67$$

When Entropy is 0
than its a leaf node

$I.G = 0.97 - 0.68$
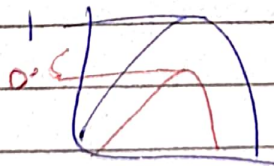
$= 0.28$

→ Gini Impurity

$$G_S = 1 - (P_y^2 + P_N^2)$$

Salay Dataset

$$G = 1 - (4/25 + 9/25)$$

$G = 0.48$

$I.G = P_G - \text{No. Tx gini (child)}$

$P_x = 0.5 \quad \leftarrow \text{Entropy} = 1$
$P_y = 0.5 \qquad \text{Gini} = 0.5$

# Handling Numerical value
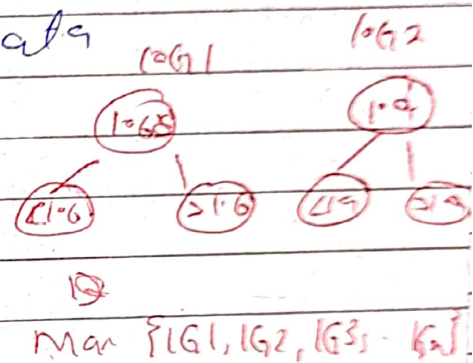
| | | |
|---|---|---|
| 1 | 3.05 | Yes |
| 2 | 4.6 | Yes |
| 3 | 2.2 | No |
| 4 | 1.6 | Yes |
| 5 | .1.05 | |

first $\downarrow$3 sort the data

| S No | User Rating | Downloaded |
|---|---|---|
| 1 | 1.6 | Yes |
| 2 | 1.9 | Yes |
| 3 | 2.2 | No |
| 4 | 2.5 | Yes |
| 5 | 2.9 | Yes |
| 6 | 3.2 | No |
| 7 | 3.3 | No |
| 8 | 3.5 | Yes |
| 9 | 3.9 | No |
| 10 | 4.1 | No |
| 11 | 4.6 | Yes |
| 12 | | |
| 13 | 4.8 | Yes |

$log_1$

$log_2$

Max $\{lG_1, lG_2, lG_3 \cdots G_n\}$

# Regression Trees

$(m, b) x$  $y = mx + b$



mark

No of hrs

Ser



85
70

1 2 3  4 5 6 7  8  No of hrs

first cut

if $n < 6$

yes        NO
|            |
if $n < 3$   mean score

yes          means

y < mean
mean

# ML

$x \leq 12$

True       False

$\boxed{A \; avg}$    $x \leq 4.5$

$\boxed{B \; avg}$    $\boxed{C \; avg}$

wages
 ↓

A

B
(·,·)

C

12     18

$$\dfrac{3 + 5}{2} = 4$$

(9, 18)

(7, 14)

(5, 12)

(3, 8)

18
16
14
12
8
4

1   2   3   4   5   6   7   8   9

$$\dfrac{12 + 14 + 28}{3} = 14.6$$

$= 18.68$

$$(14.6 - 12)^2 + (14.6 - 14)^2 +$$

$$(14.6 - 18)^2$$

$$\frac{(10-8)^2 + (10-12)^2}{8}$$

$$= \text{Error } 8/2$$

$$\frac{8+12}{2} = 10$$

$$\frac{(16-14)^2 + (16-18)^2}{8}$$

error $8/2$    $\frac{14+18}{2}$

$$= 16$$

6

$8+8 \quad > \quad 16$

$$\frac{8+12+14+18}{4} = 13$$

| y |
|---|
| 8 |
| 12 |
| 14 |
| 18 |

$$(13-8)^2 + (13-12)^2 + (13-14)^2 +$$

$$(13-18)2$$

$$SSE = \frac{}{4}$$

$$= 13$$

| | |
|---|---|
| 9 | 18 |
| 9.5 | 18.5 |
| 1.0 | 18.6 |

$$\frac{18+18.5+18.6}{3}$$

$$= 18.36$$

$$= \frac{8+12+14}{3}$$

8

$$= 11.33$$

$$= \frac{(11.33 - 8)^2 + (11.33 \cdot 12)^2 + (11.33 - 14)^2}{3}$$

$$= \frac{(18.36-18)^2 + (18.36-18.5)^2 + (18.36 - 18.6)}{4}$$

$$= \frac{}{3}$$

$$= 18.67$$

$$= 0.068$$

$$= 6.22$$