A factory is producing papers. The quality control unit applies two types of testing (durability test and strength test) to assess paper quality. The data for the same is given below:

Table III

| S. No. | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|--------|------|------|------|------|------|------|------|------|
| Durability | 7 | 6 | 7 | 6 | 3 | 1 | 4 | 3 |
| Strength | 7 | 4 | 4 | 5 | 4 | 4 | 3 | 5 |
| Quality | Good | Bad | Good | Good | Bad | Bad | Bad | Bad |

*(annotation: "features" pointing to Durability & Strength rows; "y" pointing to Quality row)*

In general, the factory produces ~~720 good quality papers out of 1000~~. Use k-nearest neighbor (KNN) with $k = 1$, and 3 to predict the quality of a new paper (durability = 5, strength = 5).  **[2+1] [CO3] [L3]**

*(annotation: "test datapoint ?")*

| D | S | labels | Distance | | K=1 | k=3 |
|---|---|--------|----------|---|-----|-----|
| 7 | 7 | G | $\sqrt{(7-5)^2+(7-5)^2} = 2\sqrt{2}$ | | | |
| 6 | 4 | B | $\sqrt{(6-5)^2+(4-5)^2} = \sqrt{2}$ | | | ✔ B |
| 7 | 4 | G | $\sqrt{(7-5)^2+(4-5)^2} = \sqrt{5}$ | | | |
| 6 | 5 | G | $\sqrt{(6-5)^2+(5-5)^2} = 1$ | | ✔ G | ✔ G |
| 3 | 4 | B | $\sqrt{(3-5)^2+(4-5)^2} = \sqrt{5}$ | | | |
| 1 | 4 | B | $\sqrt{17}$ | | | |
| 4 | 3 | B | $\sqrt{5}$ | | | |
| 3 | 5 | B | $2$ | | | ✔ B |

K=1: Good

k=3: Bad

**Naive Bayes Classifier:**

Supervised

Conditional Probability:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$
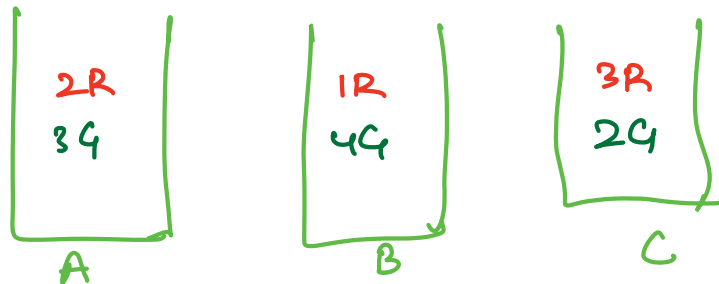
$$P(A \cap B) = P(A|B) \cdot P(B)$$

$$P(B|A) = \frac{P(B \cap A)}{P(A)}$$

$$P(B \cap A) = P(B|A) \cdot P(A)$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B \cap A)}{P(B)} = \frac{P(B|A) \cdot P(A)}{P(B)}$$

$$\boxed{P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}}$$

Eg:



| 2R | 1R | 3R |
| 3G | 4G | 2G |
| A | B | C |

Q: Prob. of getting a red ball given that A box is chosen

$$P(R|A) = \frac{2}{5}$$

Q: Prob. of getting a red ball?

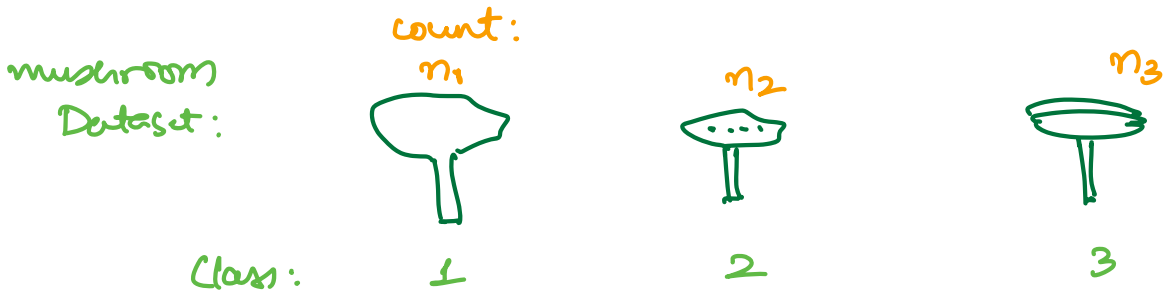$$P(R) = P(R \cap A) + P(R \cap B) + P(R \cap C)$$

Q: Prob. that bag A is chosen given that Red ball is drawn

$$P(A|R) = \frac{P(A \cap R)}{P(R)}$$

$$= \frac{P(R|A) \cdot P(A)}{P(A)}$$

$$= \frac{\overset{2/5}{P(R|A)} \cdot \overset{1/3}{P(A)}}{P(R \cap A) + P(R \cap B) + P(R \cap C)}$$

## Naive Bayes Classifier:

(likelihood)
Conditional Probability

Prior Prob.

$$\underbrace{P(A|B)}_{\text{Posterior Probability}} = \frac{P(A \cap B)}{P(B)} = \frac{\overbrace{P(B|A)} \; \overbrace{P(A)}}{\underbrace{P(B)}_{\to \text{ Prob. of } B}}$$

mushroom Dataset:

count:
$n_1$

$n_2$

$n_3$



Class:       1            2            3

features: shape, color, radius, weight ....

$$P(y=1) = \frac{n_1}{n_1 + n_2 + n_3}$$

$$P(y=2) = \frac{n_2}{n_1 + n_2 + n_3}$$

$$P(y=3) = \frac{n_3}{n_1 + n_2 + n_3}$$

Test Mushroom        Class $\to$ 1
$\to$ 2
$\to$ 3  ?

$P(y=1|x) \to 0.25$

$P(y=2|x) \to 0.15$

$P(y=3|x) \to 0.6$

max is 0.6
Test Mushroom belongs to Class 3.

$$P(\underset{A}{y=1}|\underset{B}{x}) = \frac{P(x|y=1) ** P(y=1)}{P(x)}$$

$$P(A|B) = \frac{P(B|A) ** P(A)}{P(B)}$$

$$= \frac{P(x|y=1) ** P(y=1)}{P(x \cap y=1) + P(x \cap y=2) + P(x \cap y=3)}$$

$$P(y=1|x) = \frac{P(x|y=1) ** P(y=1)}{P(x|y=1) \cdot P(y=1) + P(x|y=2)P(y=2) + P(x|y=3) \cdot P(y=3)}$$

$$P(y=2|x) = \frac{P(x|y=2) ** P(y=2)}{P(x|y=1) \cdot P(y=1) + P(x|y=2)P(y=2) + P(x|y=3) \cdot P(y=3)}$$

$$P(y=3|x) = \frac{P(x|y=3) ** P(y=3)}{P(x|y=1) \cdot P(y=1) + P(x|y=2)P(y=2) + P(x|y=3) \cdot P(y=3)}$$

$$P(y=1|x) = \frac{P(x|y=1) \cdot P(y=1)}{Denom}$$

$$P(y=2|x) = \frac{P(x|y=2) \cdot P(y=2)}{Denom}$$

$$P(y=3|x) = \frac{P(x|y=3) \cdot P(y=3)}{}$$

$$P(y=1|x) \quad \alpha \quad P(x|y=1) \cdot P(y=1)$$

$$P(y=2|x) \quad \alpha \quad P(x|y=2) \cdot P(y=2)$$

$$P(y=3|x) \quad \alpha \quad P(x|y=3) \cdot P(y=3)$$

$$P(x|y=1) = P(x_1|y=1) \cdot P(x_2|y=1) \cdot P(x_3|y=1) \quad \ldots \ldots \quad P(x_n|y=1)$$

$x$ is a test data point

Shape, Color, Radius, weight

$$P(x|y=1) = \prod_{i=1}^{n} P(x_i|y=1)$$

$$P(y=1|x) \alpha \prod_{i=1}^{n} P(x_i|y=1) \cdot P(y=1)$$

$$c \in \{1, 2, 3\}$$

$$P(y=c|x) \alpha \prod_{i=1}^{n} P(x_i|y=c) \cdot P(y=c)$$

Posterior Probability

Conditional Probability / Likelihood

Prior

eg :-

| Day | Outlook | Temp | Humidity | Wind | PlayTennis |
|-----|---------|------|----------|------|------------|
| | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D8 | Sunny | Mild | High | Weak | No |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |
| D14 | Rain | Mild | High | Strong | No |

$$c \in \{Yes, No\}$$

$$P(y = Yes) = \frac{9}{14}$$

$$P(y = No) = \frac{5}{14}$$

**Outlook**

⑨

| | Yes | No |
|--------|-----|-----|
| Sunny | 2/9 | 3/5 |
| Overcast | 4/9 | 0/5 |
| Rain | 3/9 | 2/5 |

→ $P(\text{Outlook} = \text{Sunny} \mid y = Yes)$

$P(\text{temp} = \text{cool} \mid y = Yes)$

**Temp**

| | Yes | No |
|------|-----|-----|
| Hot | 2/9 | 2/5 |
| Mild | 4/9 | 2/5 |
| Cool | 3/9 | 1/5 |

**Humidity**

| | Yes | No |
|--------|-----|-----|
| High | 3/9 | 4/5 |
| Normal | 6/9 | 1/5 |

**Windy**

| | Yes | No |
|--------|-----|-----|
| Strong | 3/9 | 3/5 |
| Weak | 6/9 | 2/5 |

Test Datapoint

$x$ [ Outlook = Sunny , Temp = Cool , Humidity = high , wind = Strong

$P(y=Yes|x) = $ P (outlook = Sunny | y = Yes) . ✔

$\quad\quad\quad\quad\quad\quad\quad$ P (temp = cool | y = Yes) .

$\quad\quad\quad\quad\quad\quad\quad$ P (humidity = high | y = Yes) .

$\quad\quad\quad\quad\quad\quad\quad$ P (wind = strong | y = Yes) .

$\quad\quad\quad\quad\quad\quad\quad$ P (y = Yes)

$$= \frac{2}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot \frac{9}{14} = \frac{1}{9 \cdot 3 \cdot 7} = \frac{1}{189}$$

$$= 0.0053$$

$P(y=No|x) = $ P (outlook = Sunny | y = No) .

$\quad\quad\quad\quad\quad\quad\quad$ P (temp = cool | y = No) .

$\quad\quad\quad\quad\quad\quad\quad$ P (humidity = high | y = No) .

$\quad\quad\quad\quad\quad\quad\quad$ P (wind = strong | y = No) .

$\quad\quad\quad\quad\quad\quad\quad$ P (y = No)

$$= \frac{3}{5} \cdot \frac{1}{5} \cdot \frac{4}{5} \cdot \frac{3}{5} \cdot \frac{5}{14} = \frac{18}{875} = 0.0206$$

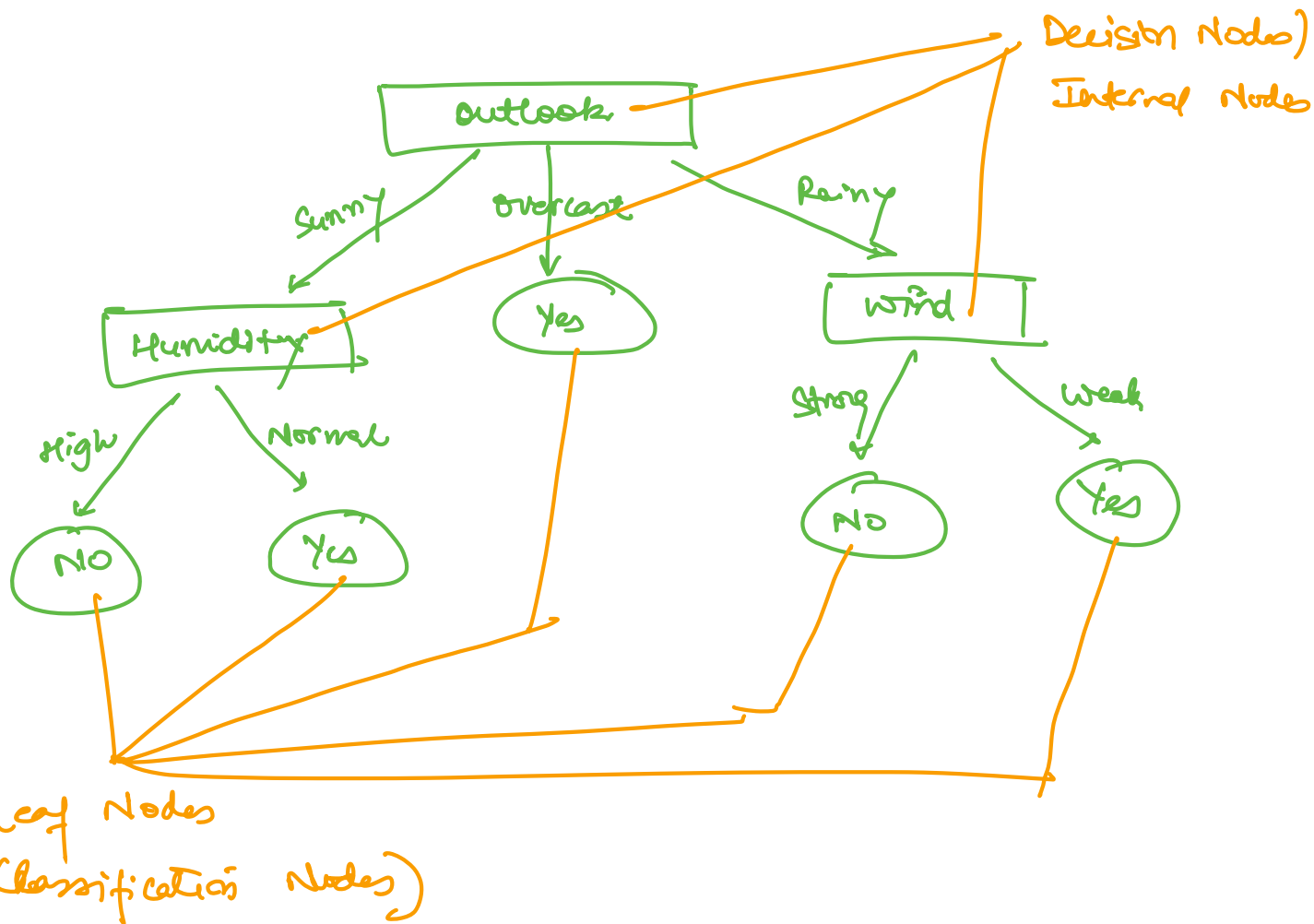Test Datapoint belongs to class No.

# Decision Trees

↳ Supervised Algo

      ↳ Classification & Regression
         $y \in$ discrete       $y \in \mathbb{R}$
             set

$x_1$: Outlook      $\in$ {Sunny, Overcast, Rainy}

$x_2$: Humidity      $\in$ {High, Normal}

$x_3$: Wind        $\in$ {Strong, Weak}

$x_4$: temperature   $\in$ {Hot, Moderate, Cold}



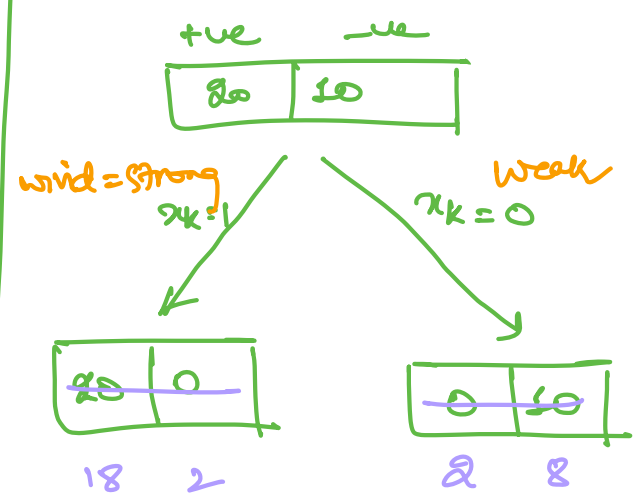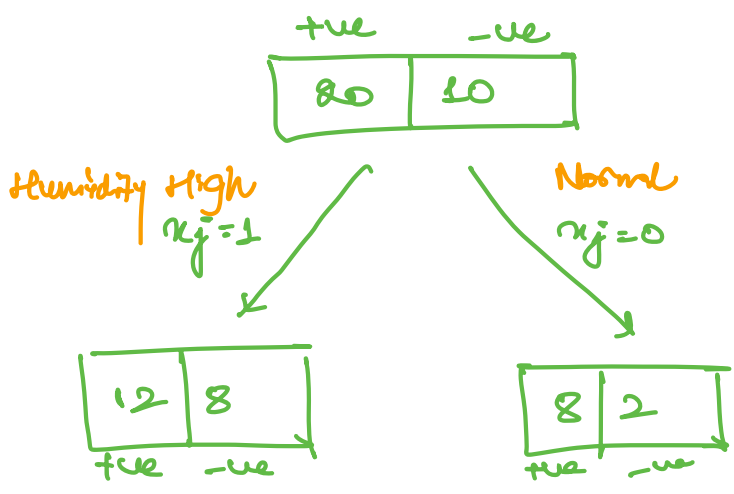Decision Nodes)
Internal Node

Leaf Nodes
(Classification Nodes)

# BUILD A DECISION TREE?

Come at a node
① ↳ data is "pure" make it a leaf Node

(2)

data  +ve, -ve

Decide
the
attribute
Split

Split
the data

Recursively
on each branch.

| +ve | -ve |
|-----|-----|
| 80  | 10  |

Humidity High
$x_j = 1$

Normal
$x_j = 0$

| 12 | 8 |
|----|---|
| +ve | -ve |

| 8 | 2 |
|---|---|
| +ve | -ve |

| +ve | -ve |
|-----|-----|
| 80  | 10  |

wind = Strong
$x_k = 1$

Weak
$x_k = 0$

| 80 | 0 |
|----|---|

18    2

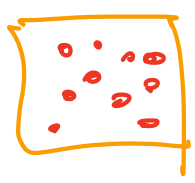| 0 | 10 |
|---|----|

2    8

Entropy          Randomness

(label)
$y \in \{1 \ldots \ldots r\}$  classes

$P(y = k) = p_k$

$$H(y) = -\sum_{i=1}^{r} p_k \log p_k$$

: max

min

$$= - \left( p_1 \log p_1 + p_2 \log p_2 + \ldots \ldots p_r \log p_r \right)$$

Boolean case   $y \in \{0, 1\}$

$$p_0 > p_1$$

$$p_0 + p_1 = 1$$

$$H(y) = -(p_0 \log p_0 + p_1 \log p_1)$$

$$= -(p_1 \log p_1 + (1-p_1) \log (1-p_1))$$

$p_1 = 0$     $\underbrace{0}$     $\underset{1}{\underbrace{1}} \; \underset{0}{\underbrace{\log 1}}$

$\underbrace{0}$

$p_1 = 1$     $\underset{0}{\underbrace{\log 1}} \quad + \quad 0$

Clain:    Cat        Dog
                      $\downarrow$

$P(Cat = 1) \longrightarrow$ Centropy $= 0$

$P(Dog = 1) \nearrow$



$H(y)$

1

concave
fx^n

0     0.5     1     $p_1$

$$\frac{d}{dp_1} \left( -(p_1 \log p_1 + (1-p_1) \log (1-p_1)) \right)$$

$$-\left( \log p_1 + \frac{p_1}{p_1} - \frac{(1-p_1)}{(1-p_1)} - \log (1-p_1) \right)$$

$$\log p_1 + \cancel{1} - \cancel{1} - \log (1-p_1) = 0$$

$$\log p_1 = \log (1-p_1)$$

$$p_1 = 1 - p_1$$

$$2 p_1 = 1 \qquad \underline{\underline{p_1 = 0.5}} = \frac{1}{2}$$

$y \in \{1, 2, \cdots r\}$

$H(Y) = -\sum_{k=1}^{r} p_k \log p_k$ is maximum at $p_k = 1/2 \; \forall k$