

## CART (Classification and Regression Trees)

### Concept. Dirtyness of Information.

The dirtyness measures the homogeneity of the data set. If the data set is homogeneous, the data belong to the same class and dirtyness is zero.

We measure dirtyness with the Gini-Index:

$$Gini = 1 - \sum_{i=1}^n p_i^2$$

The Gini-Index is a kpi that measures the dirtyness of data and is always between ZERO and ONE:  $gini \in [0,1]$

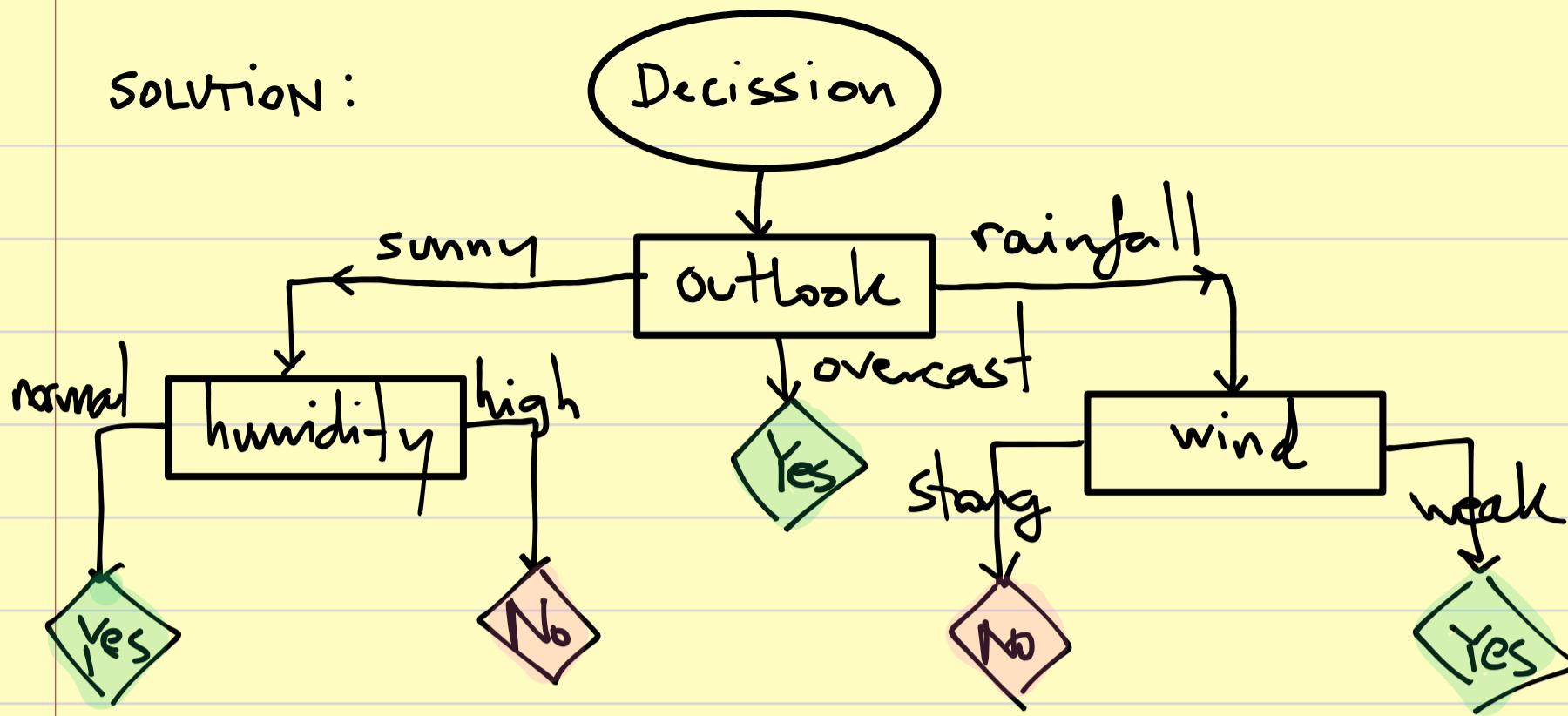
- Gini-Index = 0 means that the data is purely homogeneous and all elements belong to the same class.
- Gini-Index = 1 means that the data is maximally dirty and all elements belong to different classes.

Example. We have taken data of previous decisions made by soccer associations base on the weather about the realization of a soccer match.

Day	outlook	temperature	humidity	wind	Decision
1	sunny	hot	high	weak	No
2	sunny	hot	high	strong	No
3	overcast	hot	high	weak	Yes
4	rainfall	mild	high	weak	Yes
5	rainfall	cool	normal	weak	Yes
6	rainfall	cool	normal	strong	No
7	overcast	cool	normal	strong	Yes
8	sunny	mild	high	weak	No
9	sunny	cool	normal	weak	Yes
10	rainfall	mild	normal	weak	Yes
11	sunny	mild	normal	strong	Yes
12	overcast	mild	high	strong	Yes
13	overcast	hot	normal	weak	Yes
14	rainfall	mild	high	strong	No

Goal: create a decision tree.

SOLUTION :



OUTLOOK . Outlook is a nominal variable. It can take 3 values: sunny, overcast, rainfall.

Day	outlook	temperature	humidity	wind	Decision
1	sunny	hot	high	weak	No
2	sunny	hot	high	strong	No
3	overcast	hot	high	weak	Yes
4	rainfall	mild	high	weak	Yes
5	rainfall	cool	normal	weak	Yes
6	rainfall	cool	normal	strong	No
7	overcast	cool	normal	strong	Yes
8	sunny	mild	high	weak	No
9	sunny	cool	normal	weak	Yes
10	rainfall	mild	normal	weak	Yes
11	sunny	mild	normal	strong	Yes
12	overcast	mild	high	strong	Yes
13	overcast	hot	normal	weak	Yes
14	rainfall	mild	high	strong	No

OUTLOOK	YES	NO	#
sunny	2	3	5
overcast	4	0	4
rainfall	3	2	5

$$\sum = 14$$

$$\text{Gini}(\text{outlook Sunny}) = 1 - \sum p_i^2 = 1 - \left[ \left( \frac{2}{5} \right)^2 + \left( \frac{3}{5} \right)^2 \right] = 0.48$$

$$\text{Gini}(\text{outlook Overcast}) = 1 - \sum p_i^2 = 1 - \left[ \left( \frac{4}{4} \right)^2 + \left( \frac{0}{4} \right)^2 \right] = 0$$

$$\text{Gini}(\text{outlook Rainfall}) = 1 - \sum p_i^2 = 1 - \left[ \left( \frac{3}{5} \right)^2 + \left( \frac{2}{5} \right)^2 \right] = 0.48$$

The Gini Index for outlook is the weighted sum:

$$\text{Gini}(\text{outlook}) = \frac{5}{14} \cdot 0.48 + \frac{4}{14} \cdot 0 + \frac{5}{14} \cdot 0.48 = 0.342$$

The dirtyness of Outlook is 0'342.

## TEMPERATURE .

<u>Temperature</u>	<u>Yes</u>	<u>No</u>	<u>#</u>
hot	2	2	4
cool	3	1	4
mild	4	2	6
			$\sum = 14$

Day	outlook	temperature	humidity	wind	Decision
1	sunny	hot	high	weak	No
2	sunny	hot	high	strong	No
3	overcast	hot	high	weak	Yes
4	rainfall	mild	high	weak	Yes
5	rainfall	cool	normal	weak	Yes
6	rainfall	cool	normal	strong	No
7	overcast	cool	normal	strong	Yes
8	sunny	mild	high	weak	No
9	sunny	cool	normal	weak	Yes
10	rainfall	mild	normal	weak	Yes
11	sunny	mild	normal	strong	Yes
12	overcast	mild	high	strong	Yes
13	overcast	hot	normal	weak	Yes
14	rainfall	mild	high	strong	No

$$\text{Gini}(\text{Temp. hot}) = 1 - \left[ \left( \frac{2}{4} \right)^2 + \left( \frac{2}{4} \right)^2 \right] = 0'5$$

$$\text{Gini}(\text{Temp. cool}) = 1 - \left[ \left( \frac{3}{4} \right)^2 + \left( \frac{1}{4} \right)^2 \right] = 0'375$$

$$\text{Gini}(\text{Temp. mild}) = 1 - \left[ \left( \frac{4}{6} \right)^2 + \left( \frac{2}{6} \right)^2 \right] = 0'445$$

$$\text{Gini}(\text{Temp}) = \frac{4}{14} \cdot 0'5 + \frac{4}{14} \cdot 0'375 + \frac{4}{14} \cdot 0'445 = 0'439$$

## HUMIDITY

<u>Humidity</u>	<u>Yes</u>	<u>No</u>	<u>#</u>
high	3	4	7
normal	6	1	7
			$\sum 14$

Day	outlook	temperature	humidity	wind	Decision
1	sunny	hot	high	weak	No
2	sunny	hot	high	strong	No
3	overcast	hot	high	weak	Yes
4	rainfall	mild	high	weak	Yes
5	rainfall	cool	normal	weak	Yes
6	rainfall	cool	normal	strong	No
7	overcast	cool	normal	strong	Yes
8	sunny	mild	high	weak	No
9	sunny	cool	normal	weak	Yes
10	rainfall	mild	normal	weak	Yes
11	sunny	mild	normal	strong	Yes
12	overcast	mild	high	strong	Yes
13	overcast	hot	normal	weak	Yes
14	rainfall	mild	high	strong	No

$$\text{Gini}(\text{Humidity high}) = 1 - \left[ \left( \frac{3}{7} \right)^2 + \left( \frac{4}{7} \right)^2 \right] = 0'489$$

$$\text{Gini}(\text{Humidity normal}) = 1 - \left[ \left( \frac{6}{7} \right)^2 + \left( \frac{1}{7} \right)^2 \right] = 0'244$$

$$\text{Gini}(\text{Humidity}) = \frac{7}{14} \cdot 0'489 + \frac{7}{14} \cdot 0'244 = 0'367$$

## WIND

WIND	Yes	No	#
strong	6	2	8
weak	3	3	6

Day	outlook	temperature	humidity	wind	Decision
1	sunny	hot	high	weak	No
2	sunny	hot	high	strong	No
3	overcast	hot	high	weak	Yes
4	rainfall	mild	high	weak	Yes
5	rainfall	cool	normal	weak	Yes
6	rainfall	cool	normal	strong	No
7	overcast	cool	normal	strong	Yes
8	sunny	mild	high	weak	No
9	sunny	cool	normal	weak	Yes
10	rainfall	mild	normal	weak	Yes
11	sunny	mild	normal	strong	Yes
12	overcast	mild	high	strong	Yes
13	overcast	hot	normal	weak	Yes
14	rainfall	mild	high	strong	No

$$Gini(\text{Wind strong}) = 1 - \left[ \left( \frac{6}{8} \right)^2 + \left( \frac{2}{8} \right)^2 \right] = 0.375$$

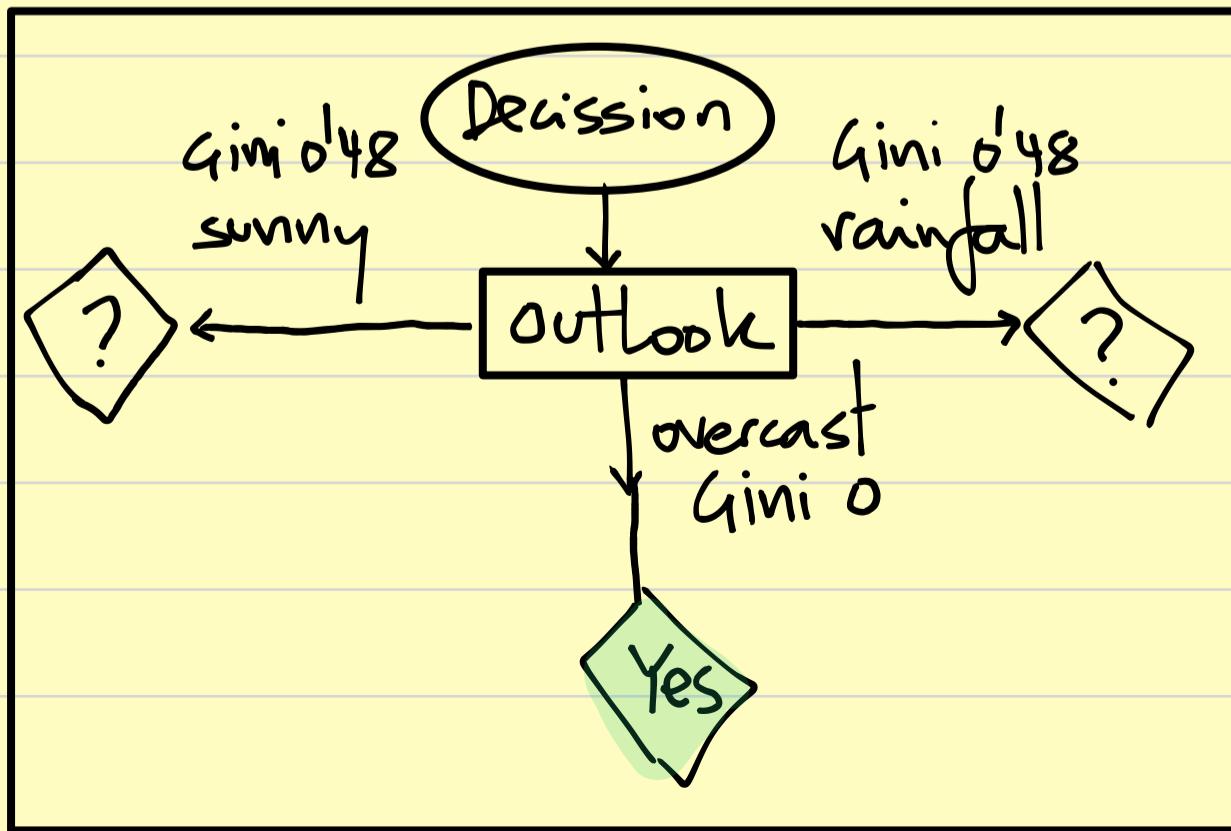
$$Gini(\text{wind mild}) = 1 - \left[ \left( \frac{3}{6} \right)^2 + \left( \frac{3}{6} \right)^2 \right] = 0.5$$

$$Gini(\text{Wind}) = \frac{8}{14} \cdot 0.375 + \frac{6}{14} \cdot 0.5 = 0.428$$

- Based on the values of the Gini for all variables, we decide which is our first node in our decision. The variable with the lowest Gini is our first node.

Feature	Gini
outlook	0.342
Temp.	0.439
Humidity	0.367
Wind	0.428

→ The outlook is the cleanest variable and delivers therefore a sharper decision.



# OUTLOOK SUNNY . TEMPERATURE

- HUMIDITY
- WIND

Day	outlook	temperature	humidity	wind	Decision
1	sunny	hot	high	weak	No
2	sunny	hot	high	strong	No
3	overcast	hot	high	weak	Yes
4	rainfall	mild	high	weak	Yes
5	rainfall	cool	normal	weak	Yes
6	rainfall	cool	normal	strong	No
7	overcast	cool	normal	strong	Yes
8	sunny	mild	high	weak	No
9	sunny	cool	normal	weak	Yes
10	rainfall	mild	normal	weak	Yes
11	sunny	mild	normal	strong	Yes
12	overcast	mild	high	strong	Yes
13	overcast	hot	normal	weak	Yes
14	rainfall	mild	high	strong	No

outlook	SUNNY	TEMPERATURE	Yes	No	#
sunny		hot	0	2	2
sunny		cold	1	0	1
sunny		mild	1	1	2

$$\sum 5$$

$$Gini(\text{Sunny+hot}) = 1 - \left[ \left( \frac{0}{2} \right)^2 + \left( \frac{2}{2} \right)^2 \right] = 0$$

$$Gini(\text{Sunny+cold}) = 1 - \left[ \left( \frac{1}{1} \right)^2 + \left( \frac{0}{1} \right)^2 \right] = 0$$

$$Gini(\text{Sunny+mild}) = 0.5$$

$$Gini(\text{sunny+Temperature}) = \frac{2}{5} \cdot 0 + \frac{1}{5} \cdot 0 + \frac{2}{5} \cdot 0.5 = 0.2$$

outlook	SUNNY	HUMIDITY	Yes	No	#
sunny		high	0	3	3
sunny		normal	2	0	2

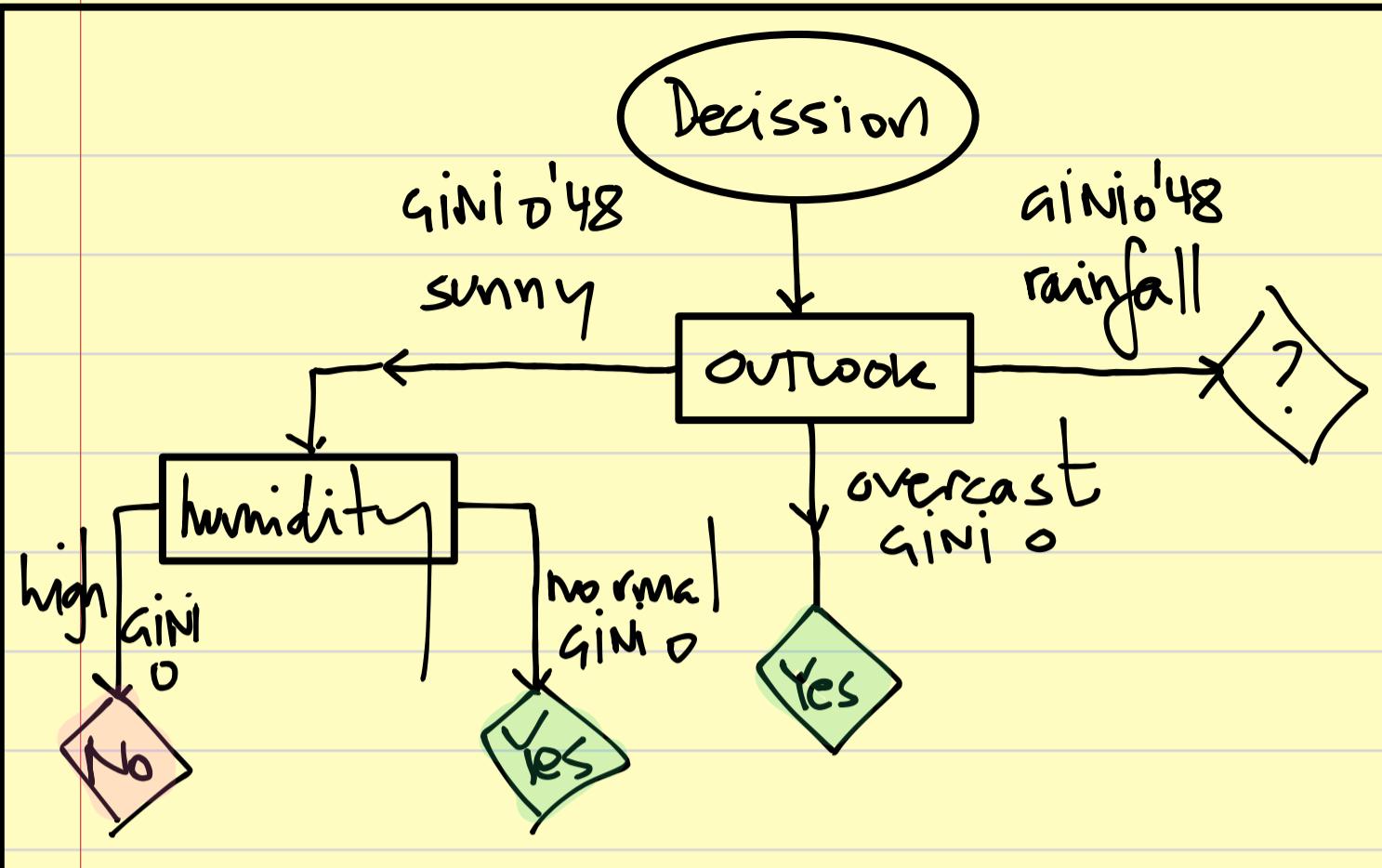
Day	outlook	temperature	humidity	wind	Decision
1	sunny	hot	high	weak	No
2	sunny	hot	high	strong	No
3	overcast	hot	high	weak	Yes
4	rainfall	mild	high	weak	Yes
5	rainfall	cool	normal	weak	Yes
6	rainfall	cool	normal	strong	No
7	overcast	cool	normal	strong	Yes
8	sunny	mild	high	weak	No
9	sunny	cool	normal	weak	Yes
10	rainfall	mild	normal	weak	Yes
11	sunny	mild	normal	strong	Yes
12	overcast	mild	high	strong	Yes
13	overcast	hot	normal	weak	Yes
14	rainfall	mild	high	strong	No

$$Gini(\text{sunny+high}) = 0$$

$$Gini(\text{sunny+normal}) = 0$$

$$Gini(\text{outlook sunny + humidity}) = 0$$

I have found a criteria (HUMIDITY) which combined with OUTLOOK SUNNY cuts perfectly the data  $Gini = 0$ .



outlook RAINFALL . TEMPERATURE

. WIND

outlook RAINFALL TEMP. Yes No #

Rainfall	cool	1	1	2
Rainfall	mild	2	1	3

Day	outlook	temperature	humidity	wind	Decision
1	sunny	hot	high	weak	No
2	sunny	hot	high	strong	No
3	overcast	hot	high	weak	Yes
4	rainfall	mild	high	weak	Yes
5	rainfall	cool	normal	weak	Yes
6	rainfall	cool	normal	strong	No
7	overcast	cool	normal	wrong	Yes
8	sunny	mild	high	weak	No
9	sunny	cool	normal	weak	Yes
10	rainfall	mild	normal	weak	Yes
11	sunny	mild	normal	strong	Yes
12	overcast	mild	high	strong	Yes
13	overcast	hot	normal	weak	Yes
14	rainfall	mild	high	strong	No

$$\text{Gini}(\text{Rainfall} + \text{Temp cool}) = 1 - \left[ \left( \frac{1}{2} \right)^2 + \left( \frac{1}{2} \right)^2 \right] = 0.5$$

$$\text{Gini}(\text{Rainfall} + \text{Temp mild}) = 1 - \left[ \left( \frac{2}{3} \right)^2 + \left( \frac{1}{3} \right)^2 \right] = 0.44$$

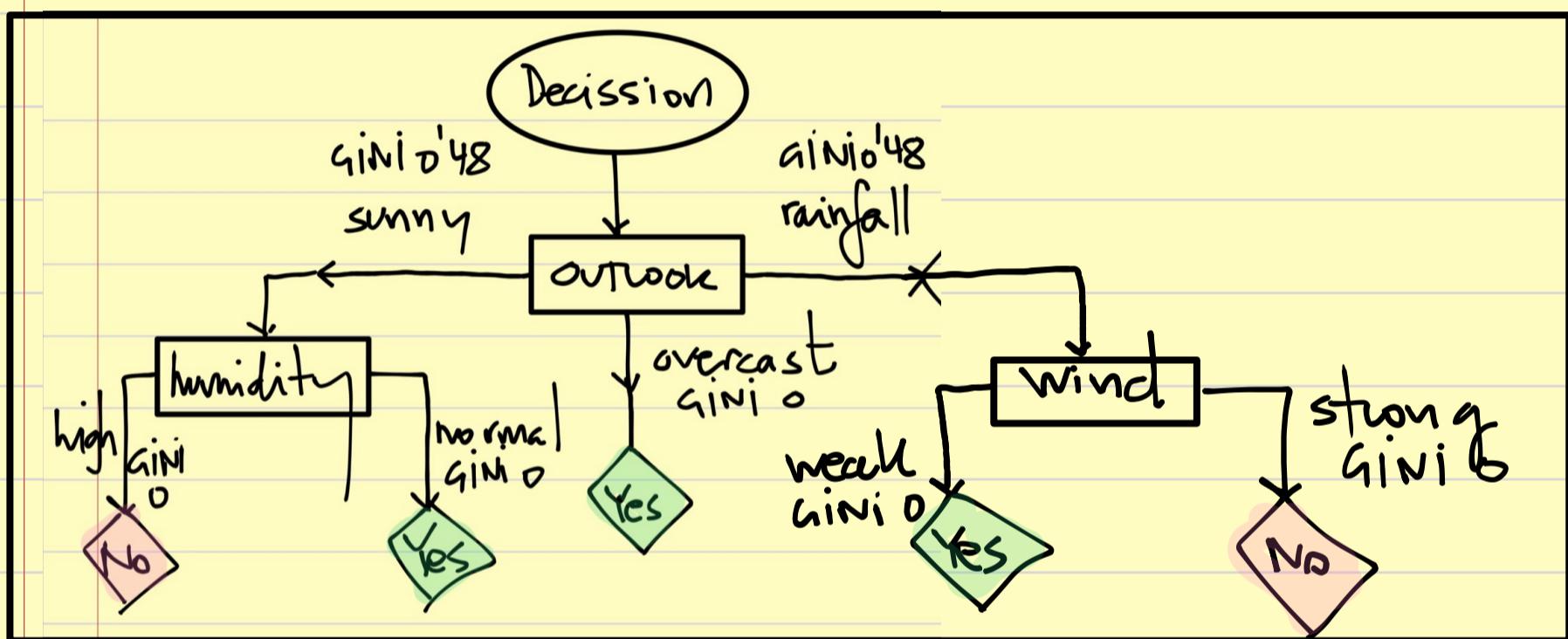
$$\text{Gini}(\text{Rainfall} + \text{Temp}) = \frac{2}{5} \cdot 0.5 + \frac{3}{5} \cdot 0.44 = 0.467$$

OUTLOOK	RAINFALL	WIND	Yes	No	#
Rainfall		weak	3	0	3
Rainfall		strong	0	2	2

$$\text{Gini}(\text{Rainfall} + \text{Wind weak}) =$$

$$= \text{Gini}(\text{Rainfall} + \text{Wind strong}) =$$

$$= \text{Gini}(\text{Rainfall}) = 0$$



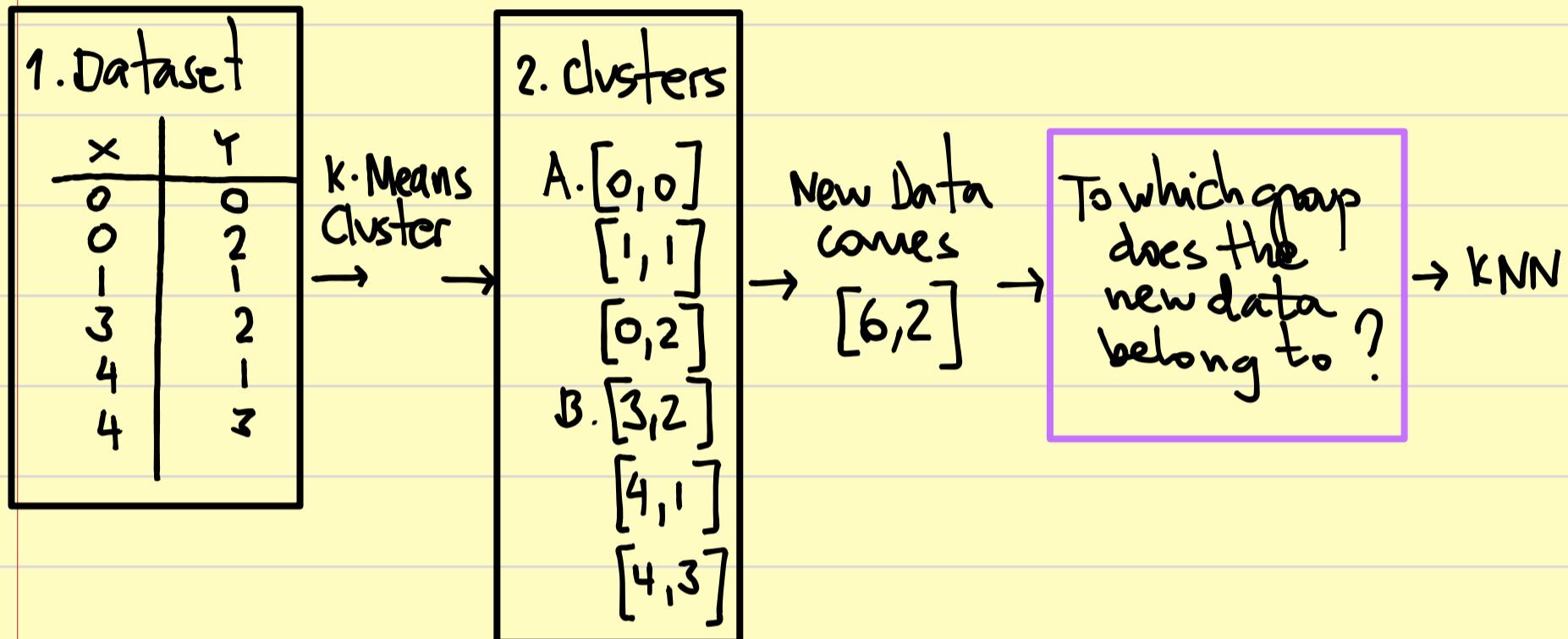
## CART for SOCCER Game Decision

Exercise . SEX . Yes or No .

	KITCHEN CLEAN	GROCERIES	STRESS @ WORK	TIME TOGETHER	D
1.	dirty	big	strong	long	Y
2.	very dirty	small	mild	short	N
3.	dirty	small	mild	short	Y
4.	clean	big	weak	long	Y
5.	very clean	small	strong	long	N

## K-NEAREST NEIGHBOUR (KNN)

Given an already clustered dataset (for instance with K-Means Clustering), we are confronted with a new datapoint and are asked to which cluster the new point belongs to.



KNN is a supervised learning algorithm used for regression and classification. KNN tries to predict the correct class for the new data, then select the k-number of points which are closest to the test data. The KNN algorithm calculates the probability of the test data belonging to the classes and assigns the class to the highest probability.

Example from data above :

A. [0, 0]

[0, 2]  
[1, 1]

B. [3, 2]

[4, 1]  
[4, 3]

New data point  $\equiv [6, 2] \equiv \alpha$

Step 1. calculate the distance from the new point to all the rest of the data.

$$d_{\alpha, A_1} = \sqrt{(6-0)^2 + (2-0)^2} = 6'32$$

$$d_{\alpha, A_2} = \sqrt{(6-0)^2 + (2-2)^2} = 6$$

$$d_{\alpha, A_3} = \sqrt{(6-1)^2 + (2-1)^2} = 5'1$$

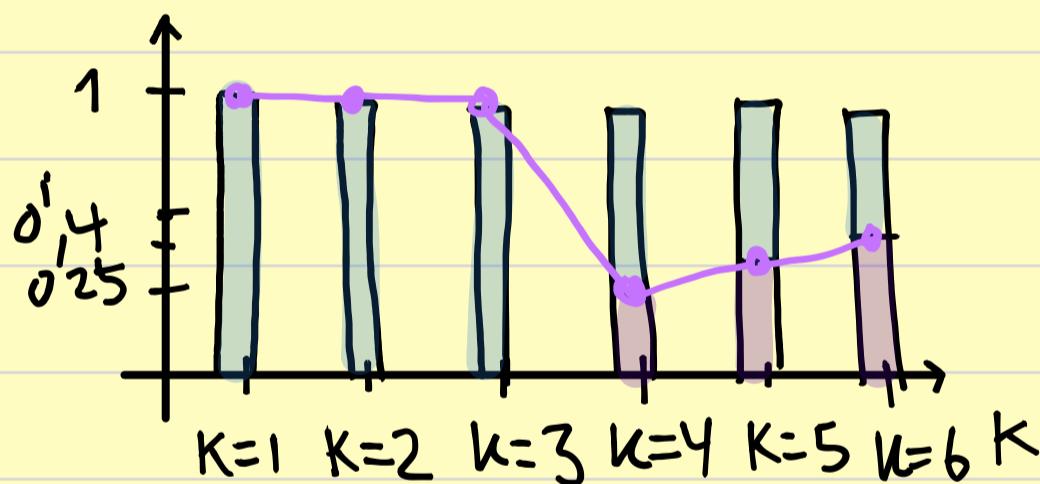
$$d_{\alpha, B_1} = \sqrt{(6-3)^2 + (2-2)^2} = 3$$

$$d_{\alpha, B_2} = \sqrt{(6-4)^2 + (2-1)^2} = 2'23$$

$$d_{\alpha, B_3} = \sqrt{(6-4)^2 + (2-3)^2} = 2'23$$

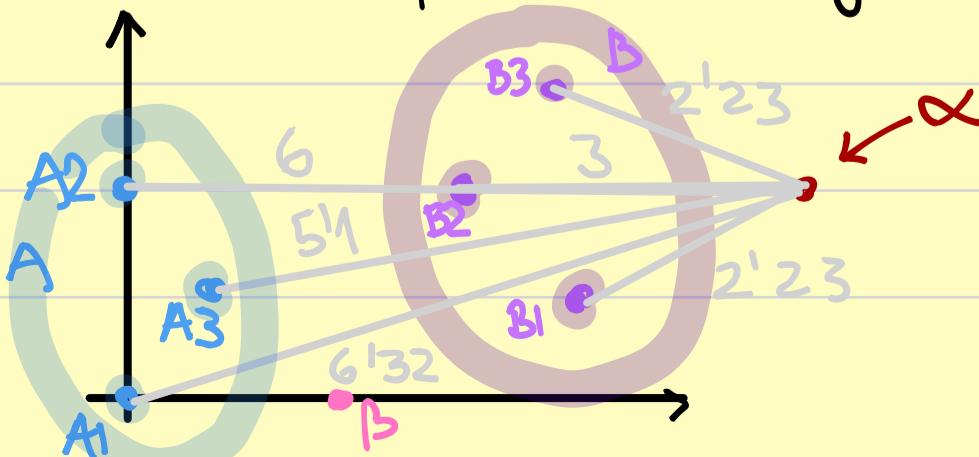
Step 2. Take different values of  $k$  (points in the group) and calculate the probability it belongs to one or other group.

$B_3, B_2, B_1, A_3, A_2, A_1$  . Hierarchy of distances



- For  $k=4$  (the point of minimum error) the probability of belonging to group B is 75%.
- Hence, the view point belongs to group B.

GRAPHICALLY



Now, with  $\alpha$  in group 3, a new point arises  $\beta = [2, 0]$

$$d_{\beta, A_1} = \sqrt{(2-0)^2 + (0-0)^2} = 2$$

$$d_{\beta, A_2} = \sqrt{(2-0)^2 + (0-2)^2} = 2^{1}83$$

$$d_{\beta, A_3} = \sqrt{(2-1)^2 + (0-1)^2} = 1^{1}414$$

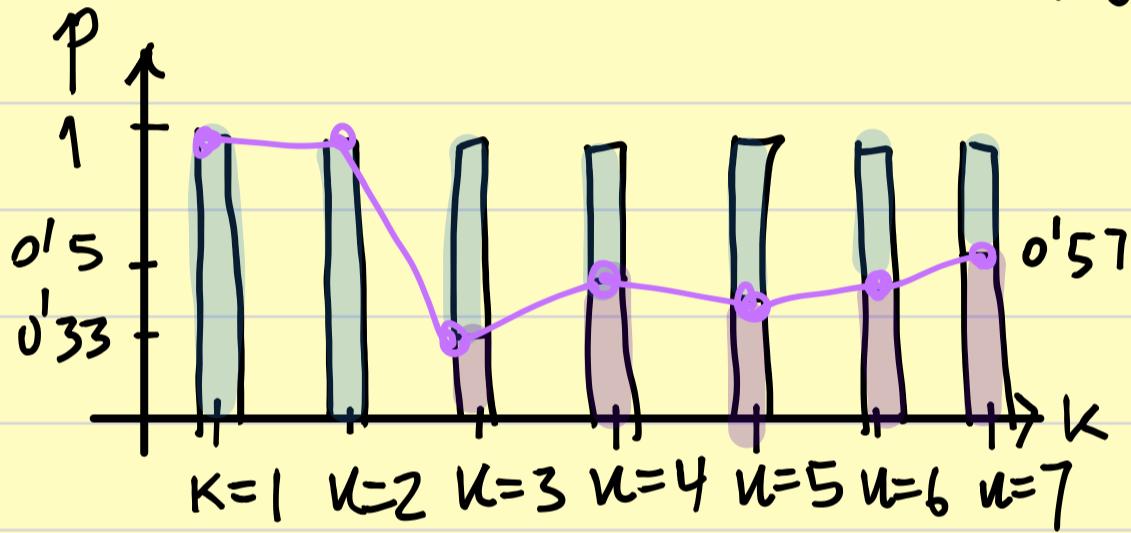
$$d_{\beta, B_1} = \sqrt{(2-3)^2 + (2-0)^2} = 2^{1}23$$

$$d_{\beta, B_2} = \sqrt{(2-4)^2 + (2-1)^2} = 2^{1}23$$

$$d_{\beta, B_3} = \sqrt{(2-4)^2 + (2-3)^2} = 3^{1}6$$

$$d_{\beta, \alpha} = \sqrt{(2-6)^2 + (0-2)^2} = 4^{1}47$$

$A_3, A_1, B_1, B_2, A_2, B_3, \alpha$ . Hierarchy of distances.



- $K=3$  delivers best probability
- $\beta$  belongs to group A.

