

Entscheidungsbäume (CART)

Konzept. Verunreinigung der Information in einer Stichprobe

Die Verunreinigung misst die Homogenität in der Datensammlung. Wenn die Daten in der Probe homogen sind, gehören die Stichproben zur gleichen Klasse und die Verunreinigung ist 0.

Wir messen die Verunreinigung mit dem Gini-Index:

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

Gini-Index ist ein Maß für die Verunreinigung einer Stichprobe. Hat einen Wert zw. [0, 1].

- Gini-Index = 0 bedeutet, dass die Stichprobe vollkommen homogen ist und alle Elemente ähnlich sind.
- Gini-Index = 1 bedeutet, maximale Verunreinigung bzw. Ungleichheit zw. den Elementen.

Ziel. Entscheiden ob ein Fußballspiel stattfindet anhand verschiedener nicht numerische (nominal) Variablen (Wetterbedingungen).

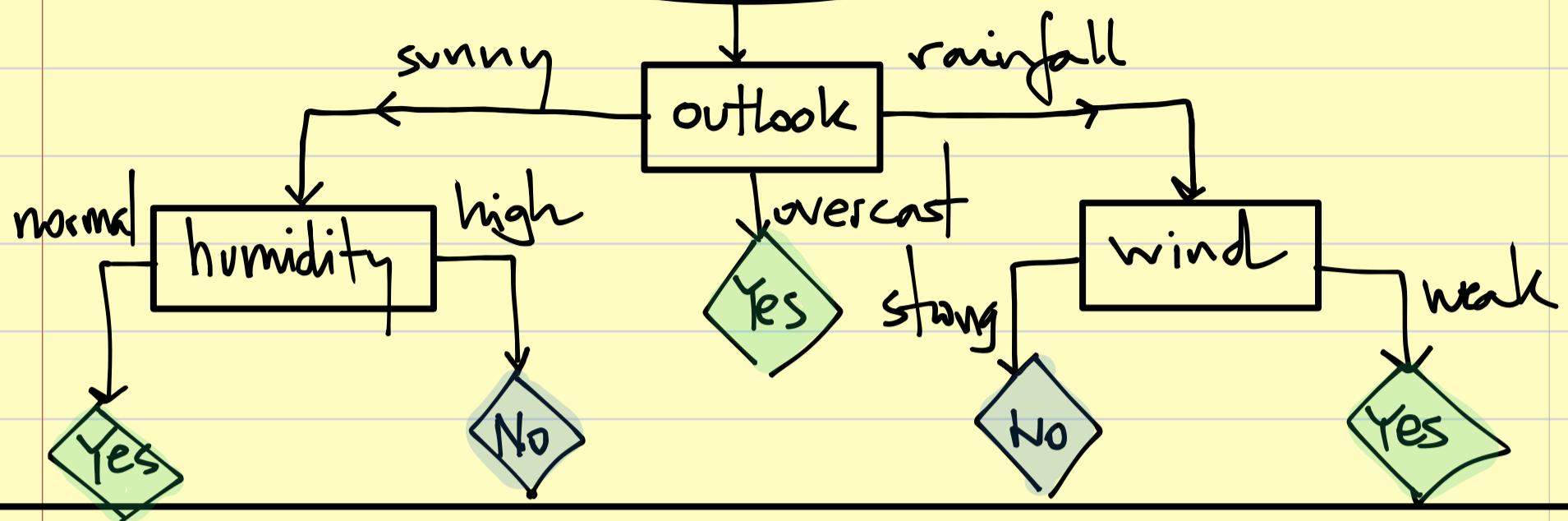
Bitte einen Entscheidungsbau erstellen.

Daten:

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

Lösung:

Entscheidung



Wir suchen den ersten Knoten mit dem geringsten Gini-Index (minimale Verunreinigung):

outlook : outlook ist ein nominals Merkmal. Es kann drei Wert annehmen: (sunny, overcast, rainfall)

<u>outlook</u>	Yes	No	#
sunny	2	3	5
overcast	4	0	4
rainfall	3	2	5

$\sum = 14$

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	No
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

$$\text{Gini}(\text{outlook}=\text{sunny}) = 1 - \left(\frac{2}{5} \right)^2 - \left(\frac{3}{5} \right)^2 = 0.48$$

$$\text{Gini}(\text{outlook}=\text{overcast}) = 1 - \left(\frac{4}{4} \right)^2 - \left(\frac{0}{4} \right)^2 = 0$$

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

Die gewichtete Summe von outlook:

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

Temperature: (hot, cold, mild)

Temperature	Yes	No	#
hot	2	2	4
cold	3	1	4
mild	4	2	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

$$\sum = 14$$

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

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

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

$$\text{Gewichtete Summe } \text{Gini}(\text{Temp.}) = \frac{4}{14} \cdot 0'5 + \frac{4}{14} \cdot 0'375 + \frac{6}{14} \cdot 0'445 = \\ = 0'43$$

humidity (hoch, normal)

humidity	Yes	No	#
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

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

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

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

Wind (strong, weak)

wind	Yes	No	#
weak	6	2	8
strong	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

$$\text{Gini}(\text{Wind}=\text{weak}) = 1 - \left(\frac{6}{8} \right)^2 - \left(\frac{2}{8} \right)^2 = 0'375$$

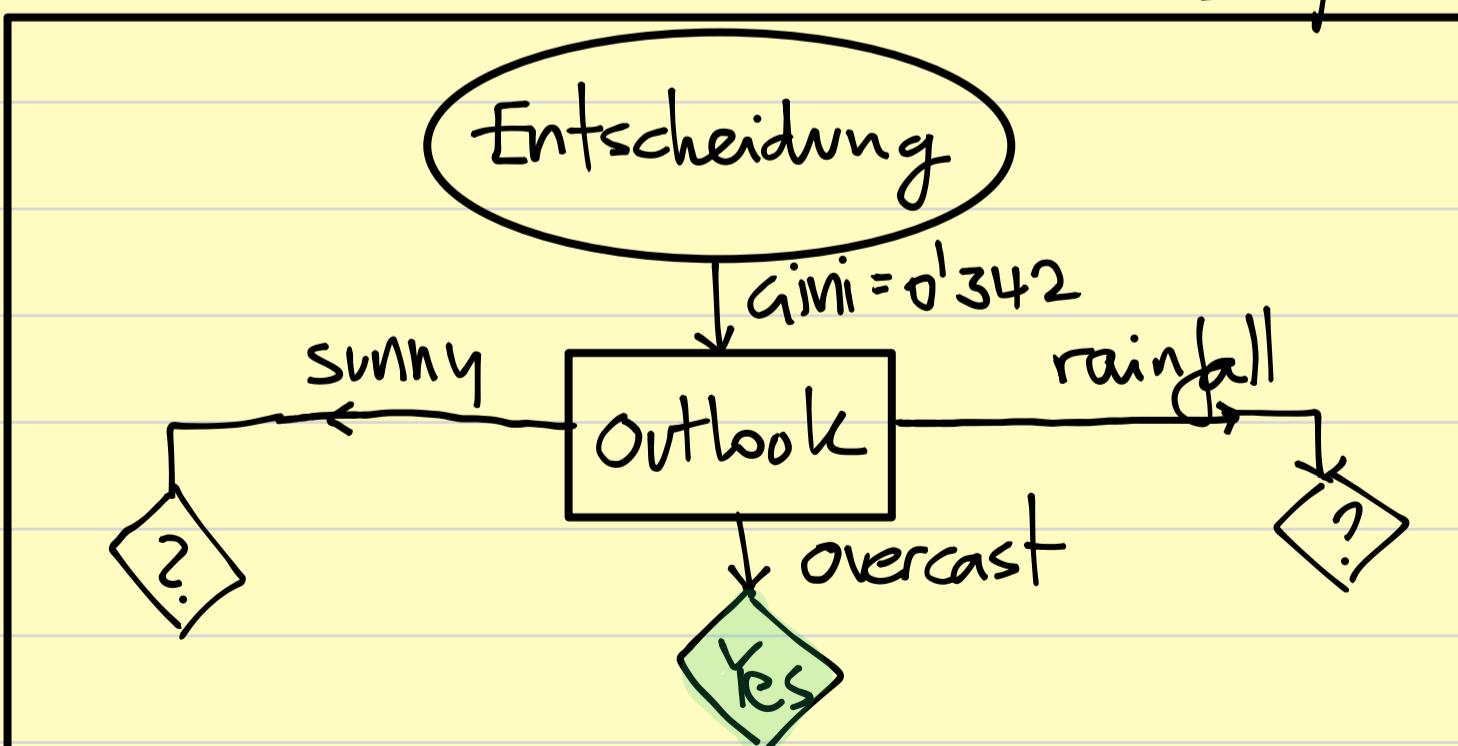
$$\text{Gini}(\text{Wind}=\text{strong}) = 1 - \left(\frac{3}{6} \right)^2 - \left(\frac{3}{6} \right)^2 = 0'5$$

$$\text{Gewichtete Gini (Wind)} = \frac{8}{14} \cdot 0'375 + \frac{6}{14} \cdot 0'5 = 0'428$$

1. Entscheidung:

feature	Gini
outlook	0'342
Temp.	0'439
Humidity	0'367
Wind	0'428

outlook hat die geringste Verunreinigung und wird als Knoten gewählt.



Gini : Sunny + Temp.

outlook	Temp	Yes	No	#
Sunny	hot	0	2	2
"	cold	1	0	1
"	mild	1	1	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	wtrong	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

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

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

$$\text{Gini}(\text{sunny+mild}) = 1 - \left(\frac{1}{2} \right)^2 - \left(\frac{1}{2} \right)^2 = 0.5$$

$$\text{Gewichtete Gini}(\text{sunny+Temp}) = \frac{2}{5} \cdot 0 + \frac{1}{5} \cdot 0 + \frac{2}{5} \cdot 0.5 = 0.25$$

Gini: Sunny + Humidity

outlook	Humidity	Yes	No	#
sunny	high	0	3	3
"	normal	2	0	2

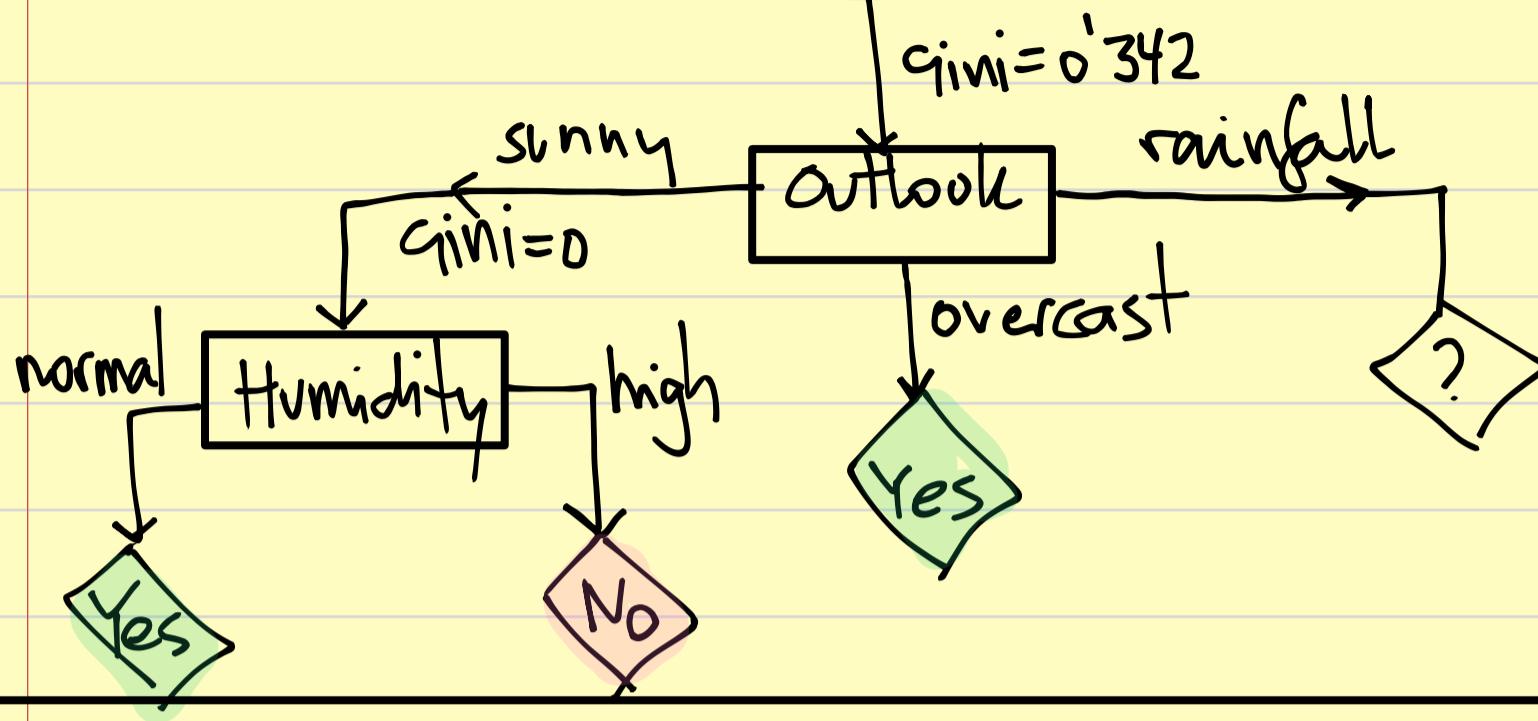
Day	outlook	temperature	humidity	wind	Decision
1	sunny	hot	high	weak	No
2	sunny	hot	high	strong	No
3	overcast	hot	high	weak	Yes
4	rainfall	mild	high	weak	Yes
5	rainfall	cool	normal	weak	Yes
6	rainfall	cool	normal	strong	No
7	overcast	cool	normal	wtrong	Yes
8	sunny	mild	high	weak	No
9	sunny	cool	normal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

$$\text{Gini}(\text{sunny+Hum. high}) = 1 - \left(\frac{0}{3} \right)^2 - \left(\frac{3}{3} \right)^2 = 0$$

$$\text{Gini}(\text{sunny+Hum. norm.}) = 1 - \left(\frac{2}{2} \right)^2 - \left(\frac{0}{2} \right)^2 = 0$$

$$\text{Gewichtete Gini}(\text{sunny+Humidity}) = \frac{3}{5} \cdot 0 + \frac{2}{5} \cdot 0 = 0$$

Entscheidung



Gini (Rainfall + Temperature)

outlook	Temp.	Yes	No	#
Rainfall	cold	1	1	2
	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	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

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

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

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

Gini (Rainfall + Wind)

outlook	wind	Yes	No	#
Rainfall	weak	3	0	3
	strong	0	2	2

Gini (Rainfall + weak wind) = 0

Gini (Rainfall + strong wind) = 0

Gewichtete Gini (Rainfall, Wind) = 0

