

Entscheidungsäume ...

Wir messen die Verunreinigung der Daten um eine Entscheidung zu treffen.

Dies wird mit dem Gini-Index gemessen.

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

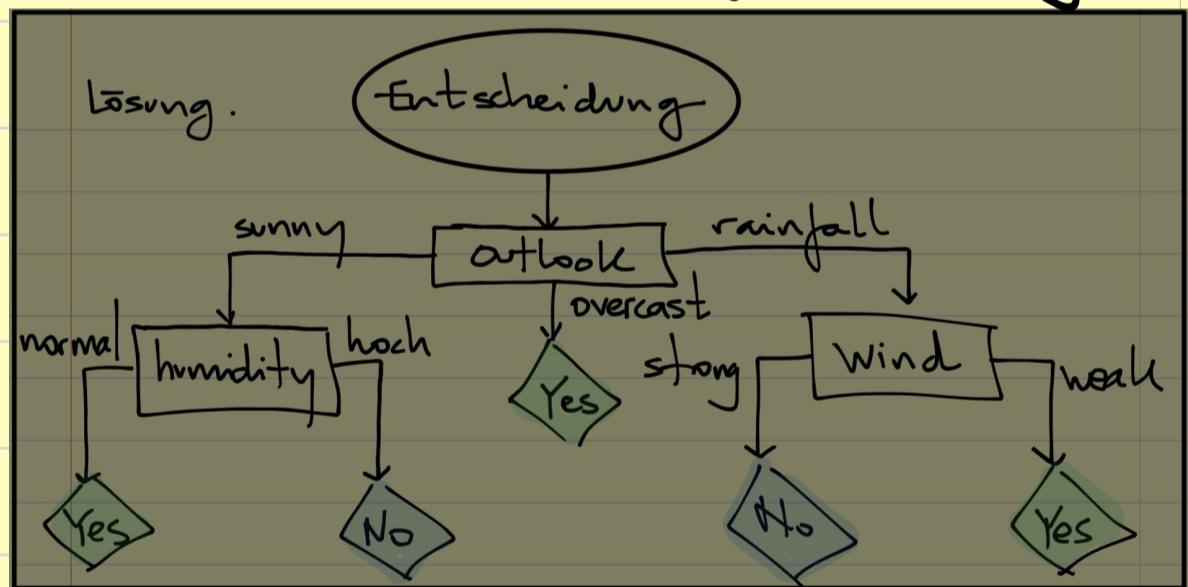
$p_i \in [0,1] \equiv w$, dafür, dass die Probe zur Klasse gehört.

$Gini = 0 \equiv$ die Stichprobe ist homogen.

$Gini = 1 \equiv$ die Stichprobe ist maximal verunreinigt.

| 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 ...



Wir suchen den ersten Knoten mit dem geringsten Gini-Index.

OUTLOOK. Ja Nein #

| | | | |
|----------|---|---|---|
| Sunny | 2 | 3 | 5 |
| Overcast | 4 | 0 | 4 |
| Rainfall | 3 | 2 | 5 |

$$Gini(\text{outlook Sunny}) = 1 - \left[\left(\frac{2}{5} \right)^2 + \left(\frac{3}{5} \right)^2 \right] = 0'48$$

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

$$Gini(\text{outlook Rain}) = 1 - \left[\left(\frac{3}{5} \right)^2 + \left(\frac{2}{5} \right)^2 \right] = 0'48$$

$$Gini(\text{outlook}) = \frac{5}{14} \cdot 0'48 + \frac{4}{14} \cdot 0 + \frac{5}{14} \cdot 0'48 = 0'342$$

| | Ja | Nein | # |
|------|----|------|---|
| hot | 2 | 2 | 4 |
| mild | 4 | 2 | 6 |
| cool | 3 | 1 | 4 |

$$\text{Gini}(\text{temp hot}) = 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right] = 0.5$$

$$\text{Gini}(\text{Temp mild}) = 1 - \left[\left(\frac{4}{6} \right)^2 + \left(\frac{2}{6} \right)^2 \right] = 0.44$$

$$\text{Gini}(\text{Temp cool}) = 1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] = 0.37$$

$$\text{Gini}(\text{Temp}) = \frac{4}{14} \cdot 0.5 + \frac{6}{14} \cdot 0.44 + \frac{4}{14} \cdot 0.37 = 0.43$$

| | Ja | Nein | # |
|--------|----|------|---|
| high | 3 | 4 | 7 |
| normal | 6 | 1 | 7 |

$$\text{Gini}(\text{Hum. high}) = 1 - \left[\left(\frac{3}{7} \right)^2 + \left(\frac{4}{7} \right)^2 \right] = 0.49$$

$$\text{Gini}(\text{Hum. high}) = 1 - \left[\left(\frac{6}{7} \right)^2 + \left(\frac{1}{7} \right)^2 \right] = 0.24$$

$$\text{Gini}(\text{humidity}) = \frac{7}{14} \cdot 0.49 + \frac{7}{14} \cdot 0.24 = 0.367$$

| | Ja | Nein | # |
|--------|----|------|---|
| weak | 6 | 2 | 8 |
| strong | 3 | 3 | 6 |

$$\text{Gini}(\text{Wind weak}) = 1 - \left[\left(\frac{6}{8} \right)^2 + \left(\frac{2}{8} \right)^2 \right] = 0.37$$

$$\text{Gini}(\text{Wind strong}) = 0.5$$

$$\text{Gini}(\text{Wind}) = \frac{8}{14} \cdot 0.37 + \frac{6}{14} \cdot 0.5 = 0.428$$

Entscheidung: Variablen

outlook

Gini

0.342 → outlook hat die

Temp

0.439

geringste
Verunreinigung

Humidity

0.367

Wind

0.428

Entscheidung

sunny Gini=0.48

Gini=0.342

rainfall Gini=0.48

OUTLOOK

Gini=0

overcast
Yes

?

OUTLOOK SUNNY { Temp
Hum.
Wind

Q. SUNNY + Temp Ja Nein #

| | | | |
|------|---|---|---|
| hot | 0 | 2 | 2 |
| mild | 1 | 1 | 2 |
| cool | 1 | 0 | 1 |

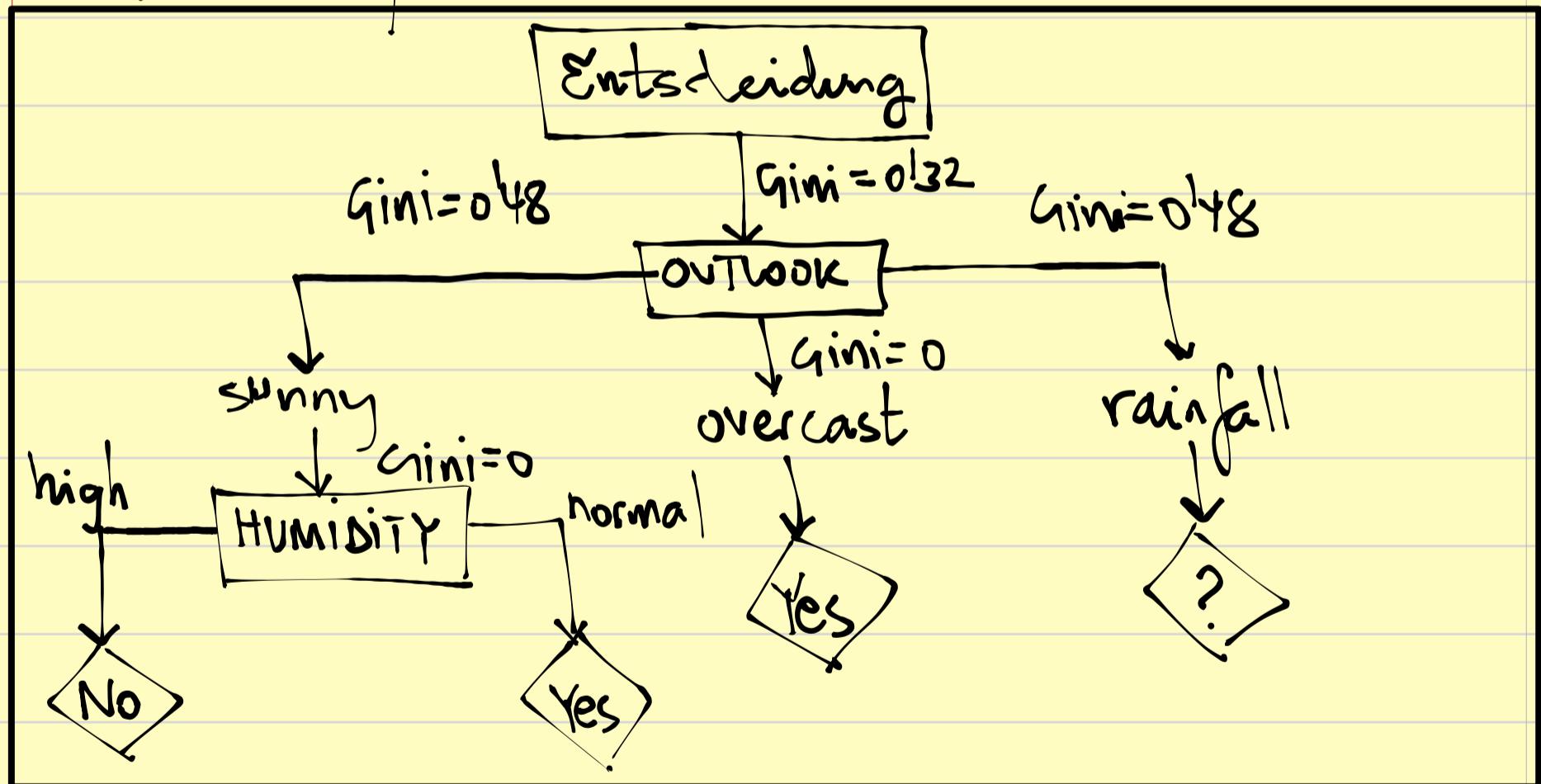
$$\text{Gini}(\text{sunny} + \text{Temp}) = \frac{2}{5} \cdot \left[1 - \left(\frac{0^2}{2} + \frac{2^2}{2} \right) \right] + \frac{2}{5} \cdot \left[1 - \left(\frac{1^2}{2} + \frac{1^2}{2} \right) \right] + \frac{1}{5} \cdot \left[1 - \left(\frac{1^2}{1} \right) \right]$$

$$= \frac{2}{5} \cdot 0 + \frac{2}{5} \cdot 0.5 + \frac{1}{5} \cdot 0 = 0.2$$

Q. SUNNY + Humidity Ja Nein #

| | | | |
|--------|---|---|---|
| high | 0 | 3 | 3 |
| normal | 2 | 0 | 2 |

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

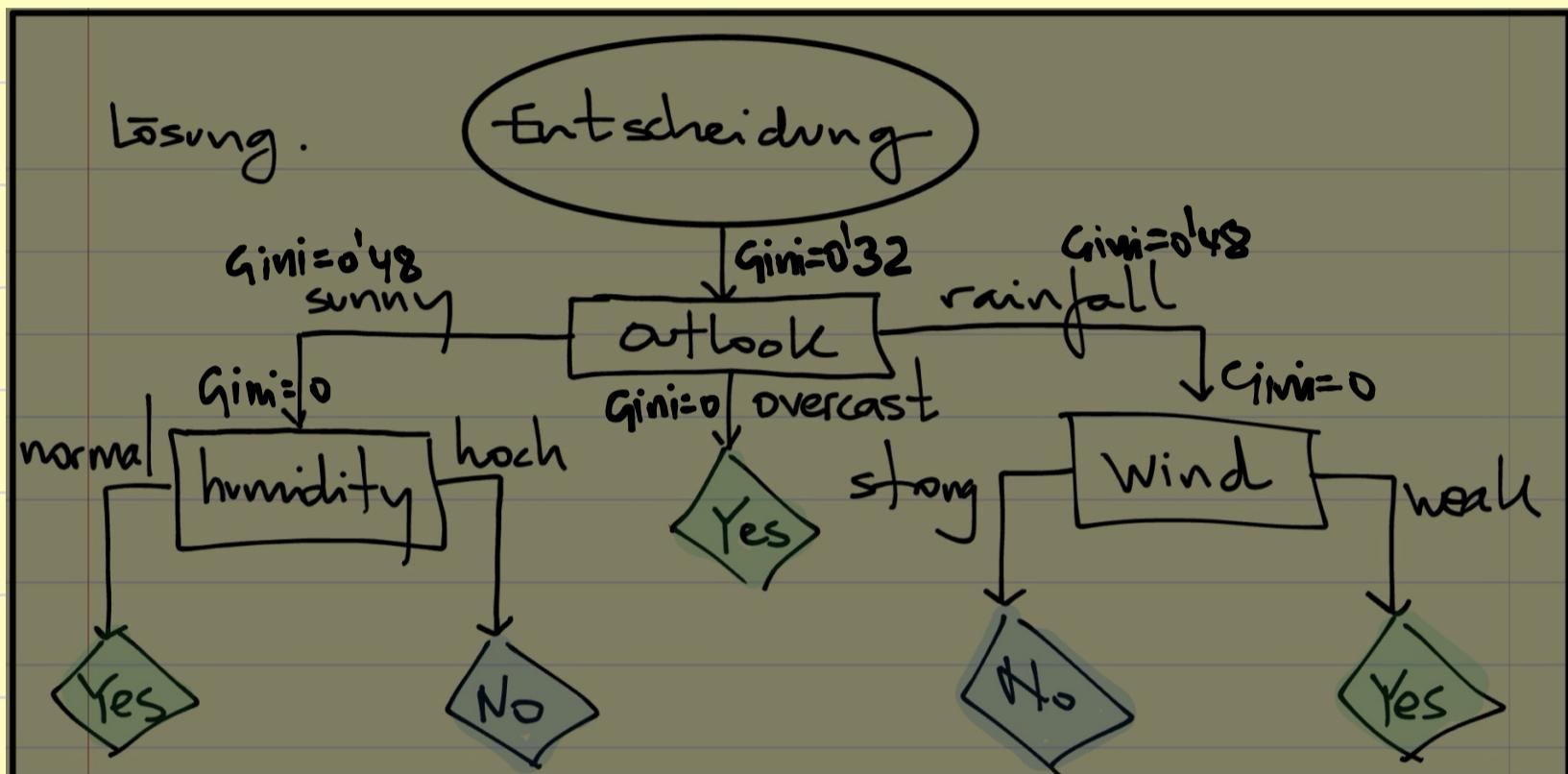


OUTLOOK RAINFALL + { Temp
Wind

0. RAINFALL + WIND Ja Nein #

| | | | |
|--------|---|---|---|
| Weak | 3 | 0 | 3 |
| Strong | 0 | 2 | 2 |

$$\text{Gini}(0. \text{Rainfall} + \text{Wind}) = 0$$



Entscheidungsbaum \equiv CART

Übung.

Beispiel. SEX Ja/Nein.

| | Wohnungsverschmutzung | Sinnvolle Gespräche | Fitness Niveau | Mond | Sex |
|-----|-----------------------|---------------------|----------------|-----------|------|
| 1. | stark | oft | hoch | voll | Ja |
| 2. | schwach | oft | gering | wachsend | Nein |
| 3. | sauber | selten | hoch | voll | Ja |
| 4. | stark | oft | mittel | abnehmend | Ja |
| 5. | stark | selten | hoch | voll | Nein |
| 6. | sauber | oft | hoch | wachsend | Ja |
| 7. | schwach | oft | mittel | voll | Nein |
| 8. | stark | oft | gering | voll | Ja |
| 9. | schwach | selten | gering | neu | Ja |
| 10. | sauber | oft | hoch | neu | Nein |

