

## Entscheidungsbäume (CART)



Regression Trees

...   
Peter      Marie

Verunreinigung der Information. Ein Dataset ist "sauber" (statistisch sauber), wenn der Gini-Index des Datensets NULL ist, und ist ..schmutzig", wenn der Gini-Index EINS ist.

Dataset: { Apfel, Orange, Erdbeer, Apfel, Zitrone }

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

$p_i$  := Wahrscheinlichkeit dafür, dass ..i.. auftritt.

$$\text{Gini} = 1 - \left[ \frac{2}{5} \right]^2 - \left[ \frac{1}{5} \right]^2 - \left[ \frac{1}{5} \right]^2 - \left[ \frac{1}{5} \right]^2 = 0'72$$

Apfel    Orange    Erdbeer

wir nutzen das Konzept um Entscheidungen zu treffen, auf Basis von vorherigen Entscheidungen.

SEX Ja/Nein .

Sauberkeit Emotionale Fitness Mond Sex  
Gespräche

1.	schnützig	oft	mittel	voll	Ja
2.	schnützig	selten	intensiv	abnehmend	Nein *
3.	sauber	selten	kein	neu	Ja
4.	mittel-sauber	oft	mittel	zunehmend	Ja
5.	Schnützig	selten	intensiv	voll	Nein
6.	schnützig	oft	kein	abnehmend	Ja
7.	sauber	selten	mittel	neu	Ja
8.	Schnützig	oft	kein	zunehmend	Nein
9.	sauber	selten	mittel	voll	Ja
10.	mittel-sauber	oft	intensiv	abnehmend	Nein

Wir gehen durch die Kriterien und suchen wir den mit den kleinsten Gini-Index .

Sauberkeit	Ja	Nein	#
schnützig	2	3	5
mittel-sauber	1	1	2
sauber	3	0	3
			10

$$\text{Gini}(\text{Saubерkeit Schmutzig}) = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 = 0'48$$

$$\text{Gini}(\text{Sauberkeit mittel/Sauber}) = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0'5$$

$$\text{Gini}(\text{Sauberkeit sauber}) = 1 - \left(\frac{3}{3}\right)^2 - \left(\frac{0}{3}\right)^2 = 0$$

$$\text{Gini}(\text{Sauberkeit}) = \frac{5}{10} \cdot 0'48 + \frac{2}{10} \cdot 0'5 + \frac{3}{10} \cdot 0 = 0'34$$

<u>Emotionale Gespräche (EG)</u>	Ja	No	#
oft	3	2	5
selten	3	2	5

$$\text{Gini}(\text{EG oft}) = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0'48$$

$$\text{Gini}(\text{EG selten}) = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0'48$$

$$\text{Gini}(\text{EG}) = \frac{5}{10} \cdot 0'48 + \frac{5}{10} \cdot 0'48 = 0'48$$

<u>Fitness</u>	Ja	Nein	#
kein	2	1	3
mittel	4	0	4
intensiv	0	3	3

$$\text{Gini}(F \text{ kein}) = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 0'44$$

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

$$\text{Gini}(F \text{ intensiv}) = 1 - \left(\frac{3}{3}\right)^2 = 0$$

$$\text{Gini}(\text{Fitness}) = \frac{3}{10} \cdot 0'44 + 0 + 0 = 0'133$$

<u>Mond</u>	Ja	Nein	#
voll	2	1	3
abnehm.	1	2	3
neu	2	0	2
zunehm.	1	1	2

$$\text{Gini}(M \text{ voll}) = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 0'44$$

$$\text{Gini}(M \text{ abn.}) = 0'44$$

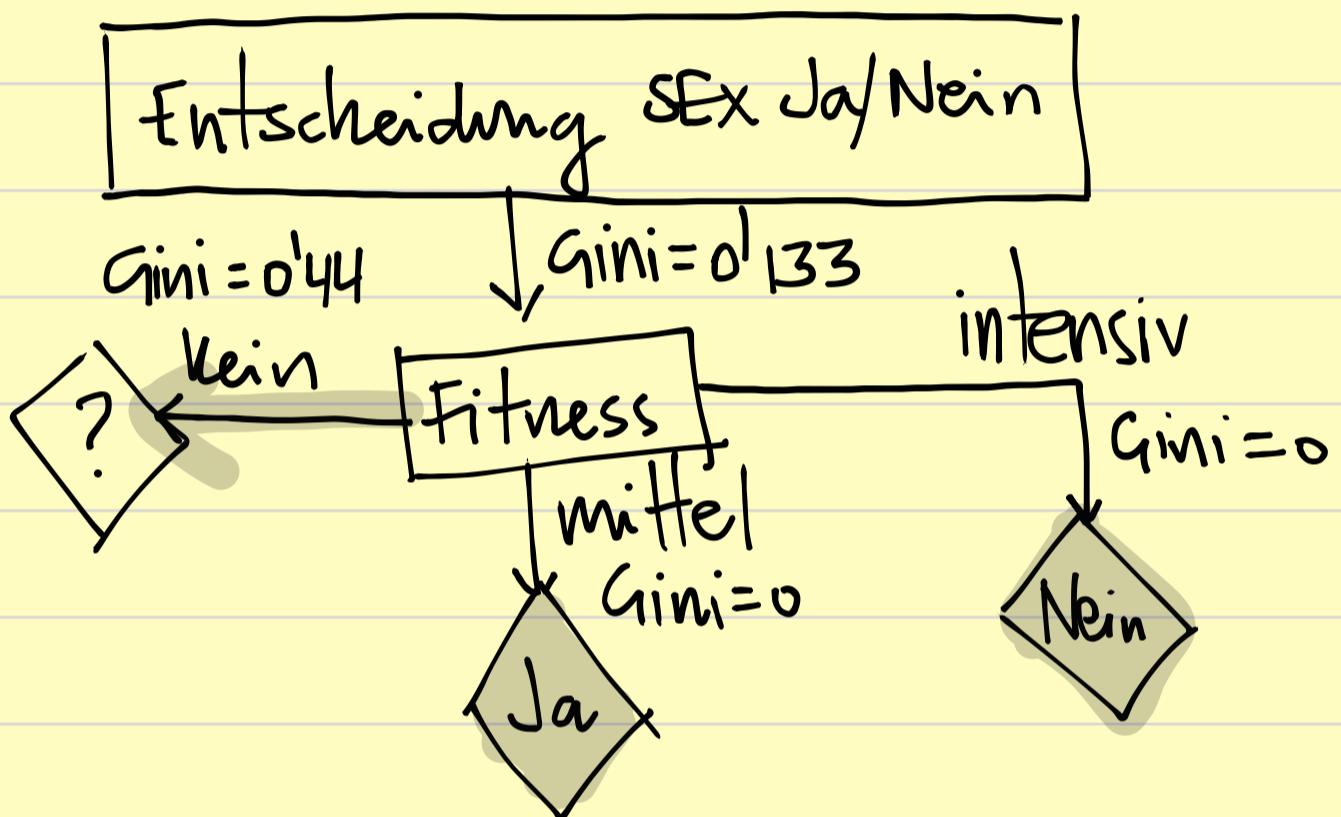
$$\text{Gini}(M \text{ neu}) = 0$$

$$\text{Gini}(M \text{ zun.}) = 0'5$$

$$\text{Gini}(\text{Mond}) = \frac{3}{10} \cdot 0'44 + \frac{3}{10} \cdot 0'44 + 0 + \frac{2}{10} \cdot 0'5 = 0'364$$

	Gini
Sauberkeit	0'34
EG	0'48
F	0'133
M	0'364

Der erste Knoten in unserem Baum ist FITNESS mit dem Kleinsten Gini.



Fitness Keint Sauberkeit    Ja    Nein #

	Ja	Nein	#
schnürtig	1	1	2

	Ja	Nein	#
sauker	1	0	1

$$\text{Gini}(F \text{ Keint Saub. Schn.}) = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0'5$$

$$\text{Gini}(F \text{ kein + Saub. saub}) = 0$$

$$\text{Gini}(F \text{ kein + Sauberkeit}) = \frac{2}{3} \cdot 0'5 + \frac{1}{3} \cdot 0 = 0'33$$

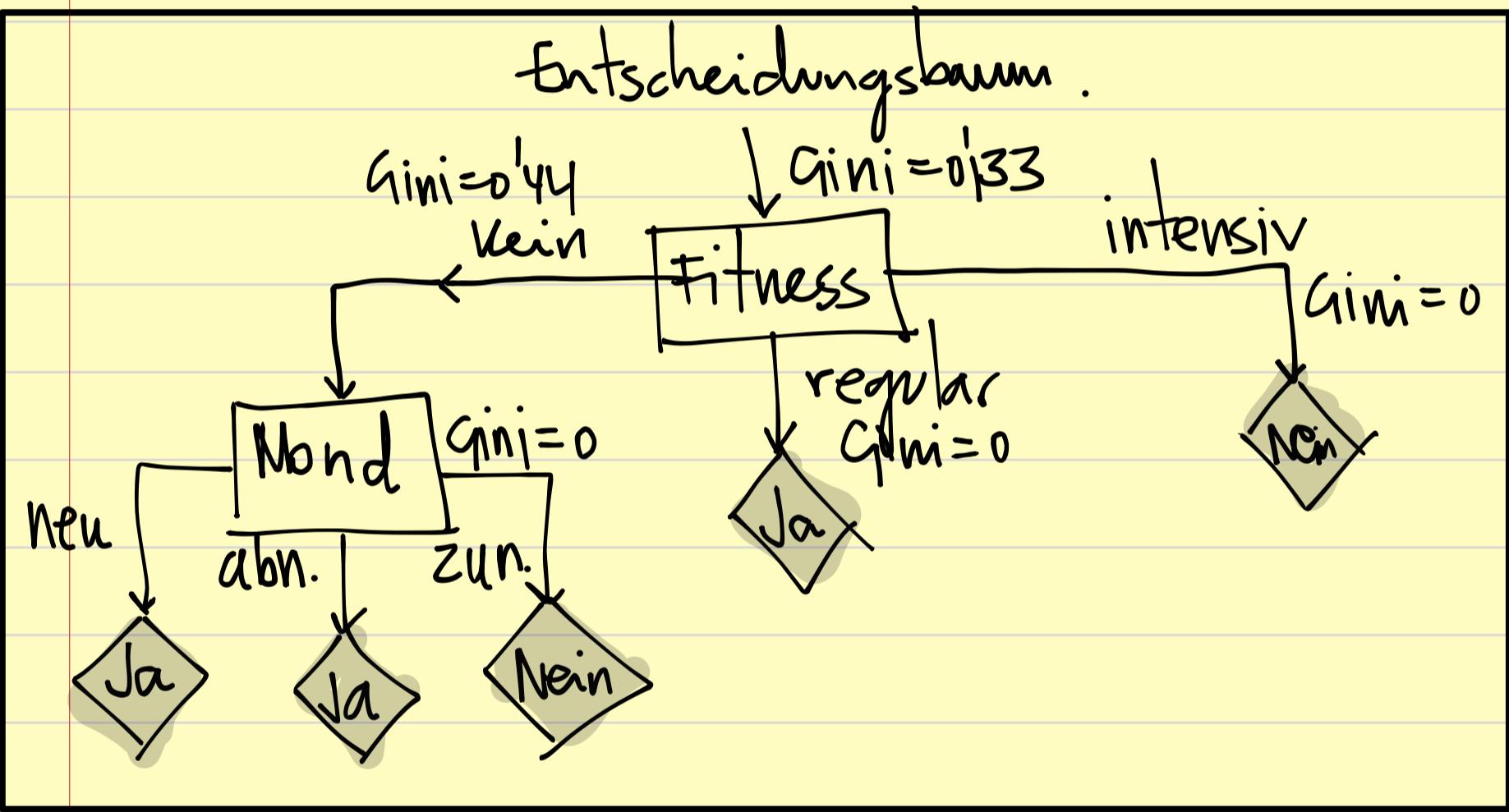
Fitness Kein + EG.    Ja    Nein #

	Ja	Nein	#
oft	1	1	2
selten	1	0	1

$$\text{Gini}(\text{Fitness Kein + EG}) = 0|33$$

<u>Fitness Kein + Mond</u>	Ja	Nein	#
neu	1	0	1 Gini=0
abnehmend	1	0	1 Gini=0
voll	0	0	0
zunehmend	0	1	1 Gini=0

$$\text{Gini}(F \text{ Kein + Mond}) = 0$$



# Verein Fußball spielen Ja/Nein

<u>outlook</u>	<u>Ja</u>	<u>Nein</u>	<u>#</u>
sunny	2	3	5
overcast	4	0	4
rainfall	3	2	5
			<u>14</u>

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

$$\text{Gini}(\text{outlook overcast}) = 0$$

$$\text{Gini}(\text{outlook rainfall}) = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0'48 \quad \left\{ \begin{array}{l} \text{Gini}(\text{outlook}) = 0'343 \\ \left( \frac{5}{14} \cdot 0'48 + \frac{5}{14} \cdot 0'48 \right) \end{array} \right.$$

<u>TEMP.</u>	<u>Ja</u>	<u>Nein</u>	<u>#</u>
hot	2	2	4
mild	4	2	6
cool	3	1	4
			<u>14</u>

$$\text{Gini}(T \text{ hot}) = 0'5$$

$$\text{Gini}(T \text{ mild}) = 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = 0'44$$

$$\text{Gini}(T \text{ cool}) = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = 0'375$$

$$\text{Gini}(\text{Temp}) = \frac{4}{14} \cdot 0'5 + \frac{6}{14} \cdot 0'44 + \frac{4}{14} \cdot 0'375 = 0'438$$

<u>LUFTFEUCHTIGKEIT</u>	<u>Ja</u>	<u>Nein</u>	<u>#</u>
high	3	4	7
normal	6	1	7

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

$$\text{Gini}(LF \text{ n}) = 1 - \left(\frac{6}{7}\right)^2 - \left(\frac{1}{7}\right)^2 = 0'245 \quad \left\{ \begin{array}{l} \text{Gini}(LF) = 0'367 \\ \frac{14}{7} \end{array} \right.$$

<u>WIND</u>	<u>Ja</u>	<u>Nein</u>	<u>#</u>
weak	6	2	8
strong	3	3	6

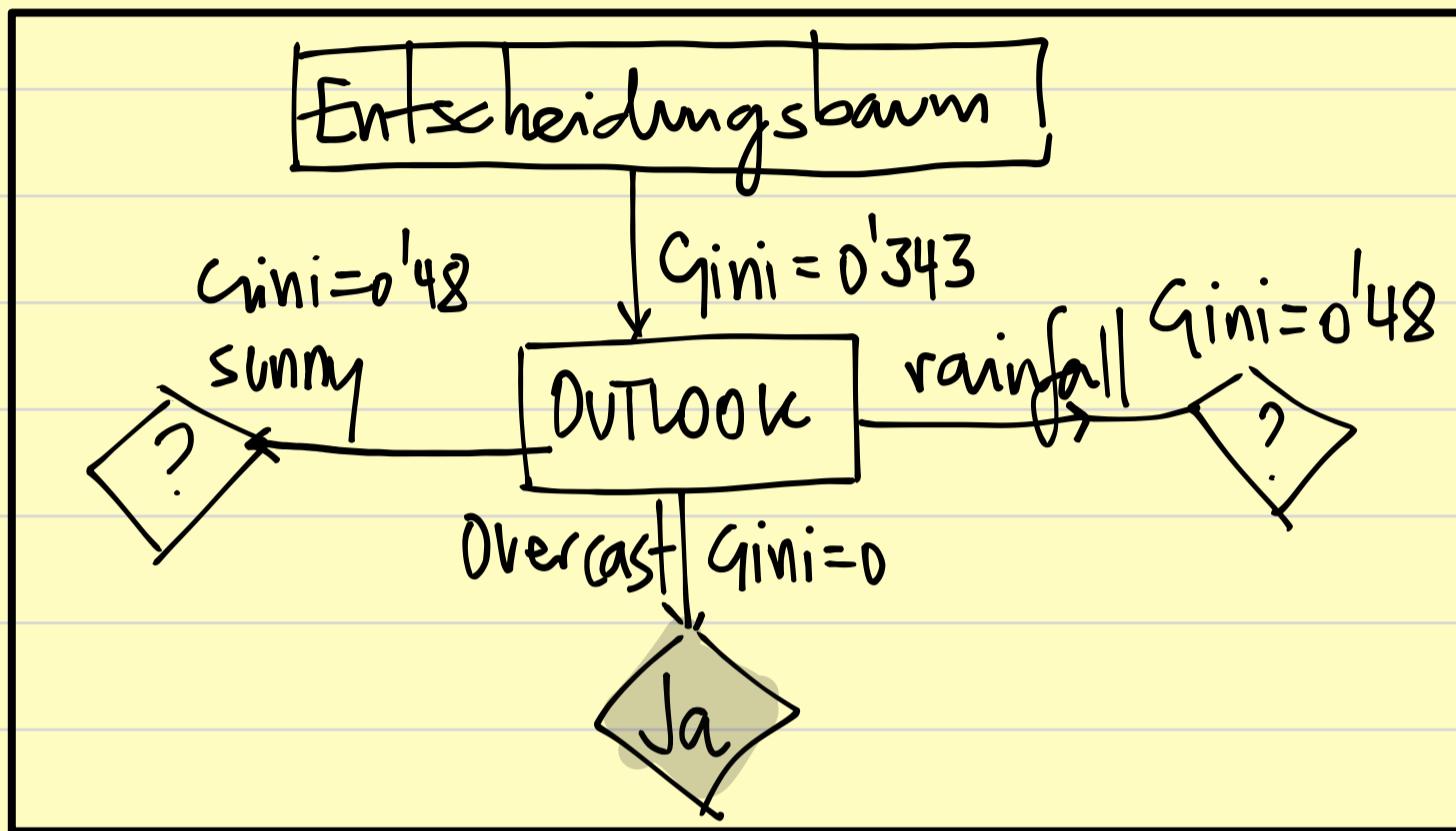
$$\text{Gini}(W \text{ weak}) = 1 - \left(\frac{6}{8}\right)^2 - \left(\frac{2}{8}\right)^2 = 0'375$$

$$\text{Gini}(W \text{ strong}) = 0'5$$

$$\text{Gini}(\text{wind}) = 0'428$$

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
Outlook	0'343 → höchste Sauberkeit.
Temp.	0'438
LF	0'367
Wind	0'428



<u>outlook S + Temp</u>	Ja	Nein	#
hot	0	2	2 → Gini = 0
mild	1	1	2 → Gini = 0'5
cool	1	0	1/5 → Gini = 0

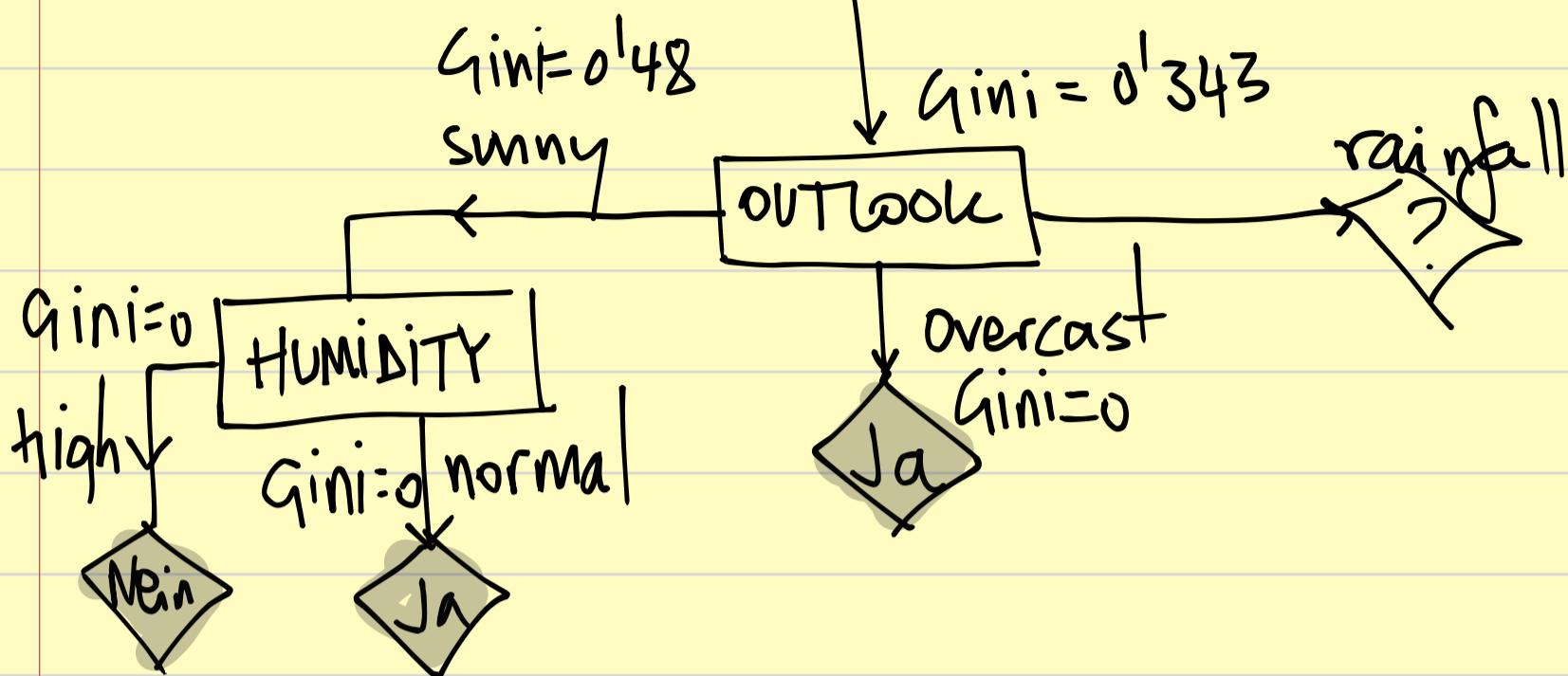
$$\text{Gini} = \frac{2}{5} \cdot 0'5 = 0'2$$

(outlook S + Temp)

<u>outlook S + LF</u>	Ja	Nein	#
high	0	3	3 → Gini = 0
normal	2	0	2 → Gini = 0

$$\text{Gini}(\text{outlook S + LF}) = 0$$

# Entscheidungsbaum



Outlook R + T	Ja	Nein	#
mild	2	1	3
cool	1	1	2

$\rightarrow \text{Gini} = 1 - \left(\frac{2}{3}\right)^2 \left(\frac{1}{3}\right)^2 = 0.44$

$\rightarrow \text{Gini} = 0.5$

$$\text{Gini}(\text{Outlook R + T}) = \frac{3}{5} \cdot 0.44 + \frac{2}{5} \cdot 0.5 = 0.364$$

Outlook R + W	Ja	Nein	#
weak	3	0	3
strong	0	2	2

$\rightarrow \text{Gini}(\text{Out. R + W}) = 0$

# Entscheidungsbaum

