

DATA SCIENCE 1

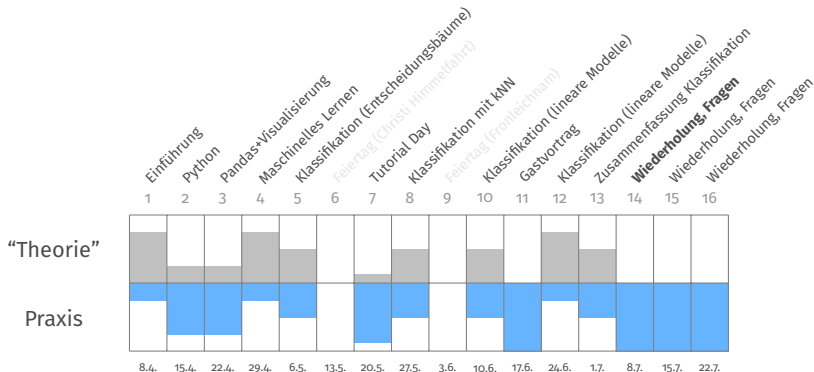
VORLESUNG 10 – OFFENE FRAGERUNDE

PROF. DR. CHRISTIAN BOCKERMANN

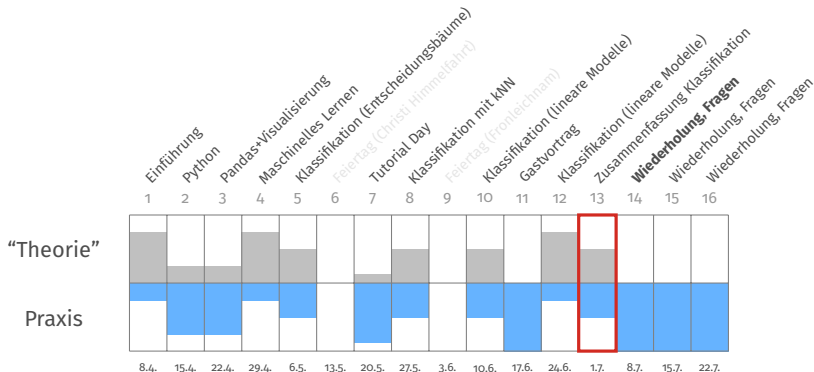
HOCHSCHULE BOCHUM

SOMMERSEMESTER 2021

Wo sind wir heute (Vorlesung 10) ?



Wo sind wir heute (Vorlesung 10) ?



Hausarbeit

- Themenausgabe: 30.7. **?Uhrzeit?**
- Bearbeitungszeit: 30.7. bis 6.8. 23:59 Uhr
- Abgabe als PDF via Moodle

Themenausgabe

Freitag, den 30.7.2021 um 9:00 Uhr
im BBB-Vorlesungsraum

<https://moodle.hs-bochum.de/mod/bigbluebuttonbn/view.php?id=126534>

Wiederholung, Fragen

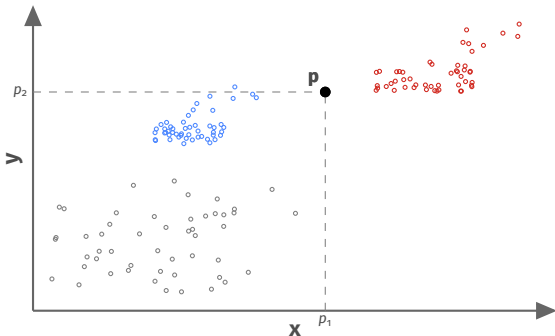
Warum brauchen wir die Normalisierung?

Beispiel: **Klassifikation von Bällen**

Wir wollen Bälle ihrer Sportart zuordnen (**Klassifikationsaufgabe**)

Umfang (cm)	Gewicht (g)	Sportart
70.29	444.30	Fussball
77.73	647.53	Basketball
53.34	427.07	Handball
57.09	406.12	Handball
68.28	440.96	Fussball
80.38	648.94	Basketball

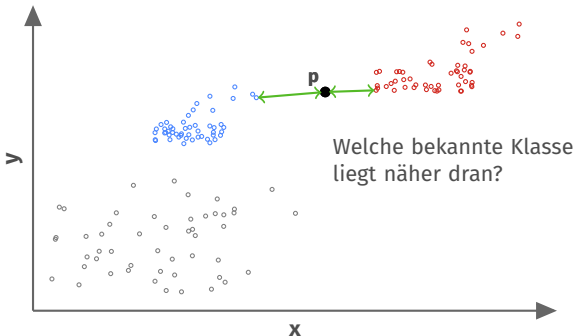
Betrachte 2-dimensionalen Raum: \mathbb{R}^2



2-dimensionaler Raum: Jeder Punkt \mathbf{p} besteht aus 2 Koordinaten:

$$\mathbf{p} = (p_1, p_2)$$

Betrachte 2-dimensionalen Raum: \mathbb{R}^2



Idee: Wir nutzen den Abstand als **Ähnlichkeit** und sagen die Klasse vorher, die am nächsten ist!

Distanzen funktionieren nur auf metrischen, skalierten Variablen

Datensätze bisher:

- Iris-Daten, Attribute waren Maße in Zentimetern
- Ball-Daten, Attribute in Gramm und Zentimetern

Frage: Was bedeutet $dist(p, q) = 10$ bei den Ball-Daten?

$$\sqrt{\underbrace{(p_1 - q_1)^2}_{\text{Umfang}} + \underbrace{(p_2 - q_2)^2}_{\text{Gewicht}}} = 10$$

0 cm

10 g

Gleicher Umfang, 10g schwerer

10 cm

0 g

10cm größer, gleiches Gewicht

Distanzen funktionieren nur auf metrischen, skalierten Variablen

Datensätze bisher:

- Iris-Daten, Attribute waren Maße in Zentimetern
- Ball-Daten, Attribute in Gramm und Zentimetern

Frage: Was bedeutet $dist(p, q) = 10$ bei den Ball-Daten?

$$\sqrt{\underbrace{(p_1 - q_1)^2}_{\text{Umfang}} + \underbrace{(p_2 - q_2)^2}_{\text{Gewicht}}} = 10$$

0 cm

10 g

Gleicher Umfang, 10g schwerer

10 cm

0 g

10cm größer, gleiches Gewicht

Wertebereich *Umfang*: 48,57 cm bis 83,97 cm

Wertebereich *Gewicht*: **315,64** g bis **686,33** g

Frage: Was bedeutet $dist(p, q) = 10$ bei den Ball-Daten?

$$\sqrt{\underbrace{(p_1 - q_1)^2}_{\text{Umfang}} + \underbrace{(p_2 - q_2)^2}_{\text{Gewicht}}} = 10$$

0 cm

10 g

Gleicher Umfang, 10g schwerer

10 cm

0 g

10cm größer, gleiches Gewicht

Frage: Was bedeutet $dist(p, q) = 10$ bei den Ball-Daten?

$$\sqrt{\underbrace{(p_1 - q_1)^2}_{\text{Umfang}} + \underbrace{(p_2 - q_2)^2}_{\text{Gewicht}}} = 10$$

0 cm

10 g

Gleicher Umfang, 10g schwerer

10 cm

0 g

10cm größer, gleiches Gewicht

Wertebereich *Umfang*: 48,57 cm bis 83,97 cm

Wertebereich *Gewicht*: **315,64** g bis **686,33** g

Frage: Was bedeutet $dist(p, q) = 10$ bei den Ball-Daten?

$$\sqrt{\underbrace{(p_1 - q_1)^2}_{\text{Umfang}} + \underbrace{(p_2 - q_2)^2}_{\text{Gewicht}}} = 10$$

0 cm

10 g

Gleicher Umfang, 10g schwerer

10 cm

0 g

10cm größer, gleiches Gewicht

Wertebereich *Umfang*: 48,57 cm bis 83,97 cm

35.40

Wertebereich *Gewicht*: **315,64** g bis **686,33** g

370.69

**Die Metrik behandelt beide Variablen gleich.
Das macht eventuell nicht immer so viel Sinn!**

Bezug z.B. zur **Wirtschaftsstatistik**

Charakterisierung von Attributen/Merkmalen/Variablen durch

- Minimum, Maximum
- Mittelwert, Standardabweichung

In welcher Relation steht $\text{dist}(p,q) = 10$ zu Umfang/Gewicht?

Idee: **Normalisierung der Attribute/Variablen**

- Anpassung der Werte auf gleichen Wertebereich
- z.B. Skalierung jeder Spalte auf [0,1]

Min-Max-Normalisierung einer Variablen X

$$X' = \frac{X - \min(X)}{\max(X) - \min(X)}$$

```
zaehler = df['Umfang'] - min(df['Umfang'])  
nenner = max(df['Umfang']) - min(df['Umfang'])
```

```
df['Umfang'] = zaehler / nenner
```


Frage: Was ist mit der Verteilung der Attribute?

- Wertebereiche mit Min/Max auf [0,1] normalisiert
- Variablen haben aber ggf. unterschiedliche Mittelwerte/Std-Abweichung?

Z-Normalisierung einer Variablen X

$$X'' = \frac{X - \mu(X)}{\sigma(X)}$$

ergibt eine Variable X'' mit Mittelwert 0 und Standardabweichung etwa bei 1.