

Big Data / Data Mining

Einführung

Aufgabe 1

Datentypen

- a) Sie haben einige Datentypen kennengelernt, mit denen Data Mining konfrontiert sein kann. Ordnen Sie den Attributen den jeweiligen Datentyp zu:

- Postleitzahl **Nominal**
- ISBN **Nominal**
- Kraftstoffverbrauch **metrisch, kontinuierlich**
- Fahrzeugleistung **metrisch, kontinuierlich**
- Hausnummer **metrisch, diskret**

- b) Finden Sie für alle Datentypen mindestens zwei weitere Beispiele.

Nominal - Gesichtsausdrücke, Vetter
Ordinal - Farben, Schulnoten
metrisch diskret - Anzahl Fußball Tore, Anzahl Pixel
kont. - Distanz - Zeit

Aufgabe 2

Abstands- und Ähnlichkeitsmaße

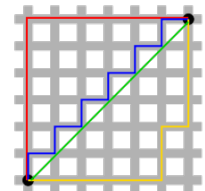
Wir betrachten einige typische Distanzfunktionen:

Hamming-Distanz $\text{dist}_H(v, w) = \text{count}_i(v_i \neq w_i) \rightarrow$ vergleicht elementweise $011001, 001000 \Rightarrow 2$

Euklidische Distanz $\text{dist}_E(v, w) = \sqrt{\sum_i (v_i - w_i)^2} \rightarrow$ Abstand zweier Punkte im n -dim Raum

Manhattan-Distanz $\text{dist}_{Man}(v, w) = \sum_i |v_i - w_i| \rightarrow$ Summe der Längen aller Achsen / Dimensionen

Maximum-Distanz $\text{dist}_{Max}(v, w) = \max(|v_i - w_i|) \rightarrow$ größtmöglicher Abstand
 Tschebyscheff-Abstand



- a) Informieren Sie sich zu diesen Distanzfunktionen in Cleve und Lämmel 2014, S. 43–46.

- b) Berechnen Sie die Distanz zwischen den Punkten (2, 3) und (8, 7). Verwenden Sie alle 4 aufgeführten Distanzfunktionen (vgl. Abbildung 1).

$$u = (2, 3) \quad v = (8, 7)$$

$$\text{dist}_H(u, v) = 2 \quad \text{dist}_{Man}(u, v) = 10$$

$$\text{dist}_E(u, v) = 7,2111 \quad \text{dist}_{Max}(u, v) = 6$$

c)

$$\text{Dist}_H(u, v) = 3$$

$$(u, w) = 3$$

$$(v, w) = 2$$

$$\text{Dist}_E(u, v) = \sqrt{17}$$

$$(u, w) = \sqrt{10}$$

$$(v, w) = \sqrt{52}$$

$$\text{Dist}_M(u, v) = 6$$

$$(u, w) = 8$$

$$(v, w) = 10$$

$$\text{Dist}_{\text{max}}(u, v) = 4$$

$$(u, w) = 4$$

$$(v, w) = 7$$

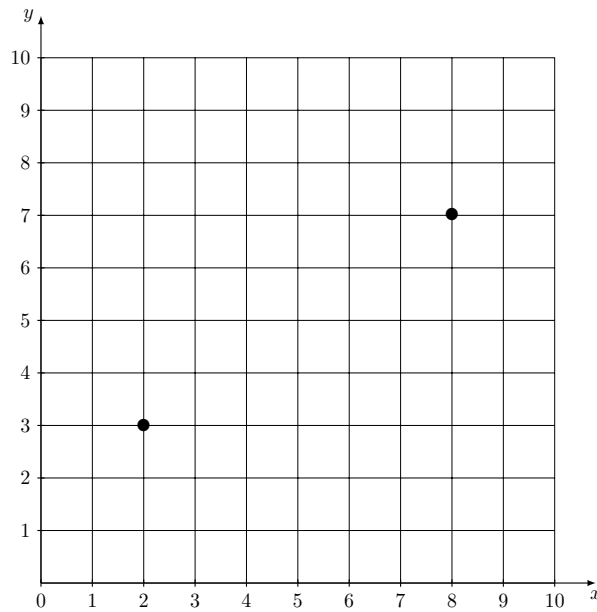


Abbildung 1: Distanzen im zweidimensionalen Raum

$u =$ $v =$ $w =$

- c) Berechnen Sie jeweils die Distanz zwischen den Punkten $(0, 1, 2)$, $(1, 5, 3)$ und $(4, -2, 3)$ für alle 4 Distanzfunktionen.
- d) Begründen Sie, wieso die Hamming-Distanz unempfindlich gegen Ausreißer ist.

Weil Elemente nur paarweise auf Gleichheit geprüft wurden

Aufgabe 3

Erste Schritte in RapidMiner

- a) Registrieren Sie sich bei <https://rapidminer.com/> und beziehen Sie eine Educational-Lizenz für RapidMiner Studio. Starten Sie RapidMiner Studio und synchronisieren Sie Ihre Lizenz im Lizenzmanager mit Ihrem Account.

Hinweis: Benutzen Sie während der Registrierung unbedingt Ihre Hochschul-E-Mail-Adresse. Sie können die Software ohne Installation starten. Ihr Dozent wird Ihnen die benötigten Dateien an einem geeigneten Ort bereitstellen. Auf privaten Geräten können Sie auch die Setups von der RapidMiner-Website verwenden und installieren.

- b) Machen Sie sich mit dem Dashboard vertraut.

Literatur

Cleve, Jürgen und Uwe Lämmel (2014). *Data mining*. Studium. München: De Gruyter Oldenbourg. 306 S. ISBN: 978-3-486-72034-1 978-3-486-71391-6.