

Tutorium 3

Montag, 25. Mai 2020 15:37

- 1 a) Warum ist es sinnvoll, die Variablen zu standardisieren (zentrieren und skalieren), oder zu normalisieren (Maximum abziehen und durch Spannweite/Range teilen) bevor man die euklidische Distanzmatrix berechnet?

Standardisieren:

$$X^* = \frac{x - \bar{x}}{\sqrt{\text{Var}(x)}}$$

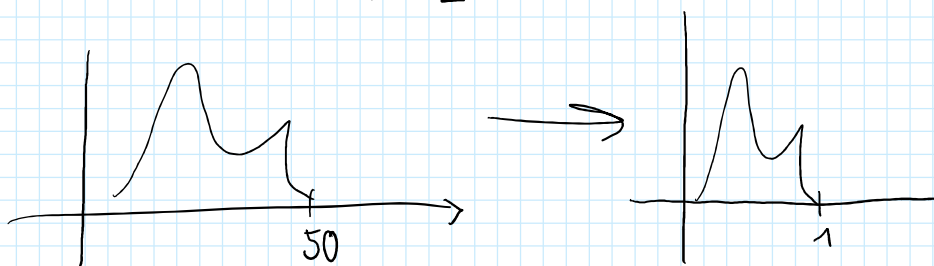
→ danach $E(X^*) = 0$ & $\text{Var}(X^*) = 1$

Normalisieren:

$$X^{**} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

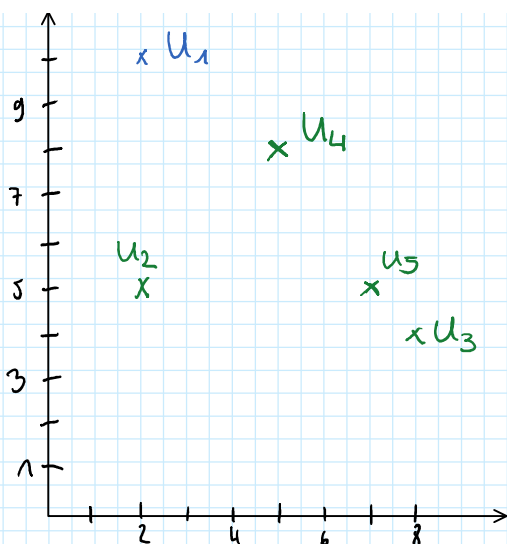
→ behält Form der Verteilung

→ $X^{**} \in [0, 1]$



2a)

- a) Geben Sie $C^{(0)}$ an und berechnen Sie die neuen Schwerpunkte für jede Klasse.



$$C^0 = \{ \{u_1\}, \{u_4, u_2, u_3, u_5\} \}$$

neue Schwerpunkte: $x_1^{(1)} = \begin{pmatrix} 2 \\ 10 \end{pmatrix}, x_2^{(1)} = \begin{pmatrix} \frac{2+8+5+7}{4} \\ \frac{5+4+8+5}{4} \end{pmatrix} = \begin{pmatrix} 5.5 \\ 5.5 \end{pmatrix}$

2b)

b) Führen Sie das k-means Clusterverfahren zu Ende.

$$D^{(1)} = \begin{pmatrix} \overset{\text{min}}{\circ} & 2.5 & 7.2 & 1.3 & 5.0 \\ 32.5 & \circ & 8.5 & 6.5 & 2.5 \end{pmatrix} \begin{matrix} \leftarrow x_1 \text{ Distanz zu Schwerpunkt 1} \\ \leftarrow x_2 \text{ Distanz zu Schwerpunkt 2} \end{matrix}$$

$\uparrow \quad \uparrow \quad \uparrow \quad \uparrow \quad \uparrow$
 $u_1 \quad u_2 \quad u_3 \quad u_4 \quad u_5$

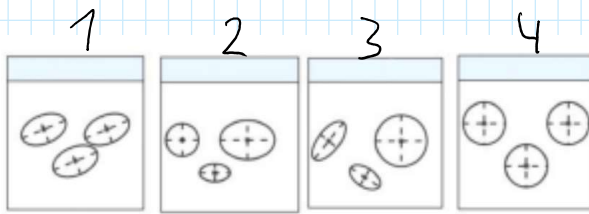
$$d(u_1, x_2^{(1)}) = (2 - 5.5)^2 + (10 - 5.5)^2 = 12.25 + 20.25 = 32.5$$

$$C^{(1)} = \{ \{u_1\}, \{u_2, u_3, u_4, u_5\} \} = C^{(0)}$$

neue Partition = alte Partition → Konvergenz

3b) i)

i) Lesen Sie sich die Hilfe zu der Funktion `mclustModelNames` aus dem Paket `mclust` durch und ordnen Sie den unten stehenden Abbildungen jeweils ein Modell zu.



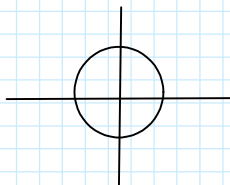
1: EEE ellipsoidal, equal volume, shape

2: VII diagonal, varying volume, shape

3: VVV ellipsoidal, varying volume, shape, & orientation

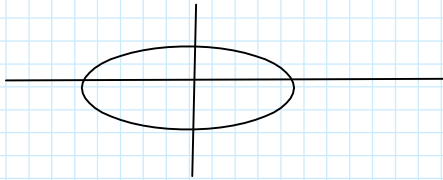
4: EII spherical, equal volume (→ ähnliche Ergebnisse wie k-means)

spherical: keine Korrelationen unabhängige Variablen
gleiche Varianzen der Variablen



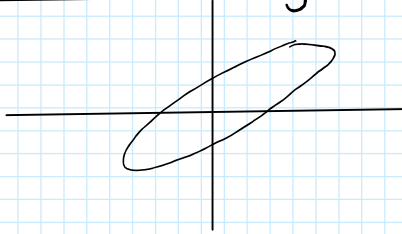
diagonal: Streckung in x & y Richtung → unterschiedliche Varianzen
→ keine Korrelationen

→ keine Korrelationen
(↔ unabhängige Variablen)



ellipsoidal:

streckung in beliebige Richtung → Korrelationen möglich
→ unterschiedliche Varianzen



ii) Beschreiben Sie anhand der Kovarianzmatrix, wie Distribution, Volume, Shape und Orientation ausgeprägt sind.

spherical	$\Sigma_i = \begin{pmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_1^2 \end{pmatrix}$	} Kovarianz für <u>ein</u> Cluster
diagonal	$\Sigma_i = \begin{pmatrix} \sigma_1^2 & & 0 \\ & \sigma_2^2 & \\ 0 & & \ddots \end{pmatrix}$	
ellipsoidal	$\Sigma_i = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots \\ \sigma_{21} & \sigma_2^2 & \ddots \\ \vdots & \vdots & \ddots \end{pmatrix}$	

$\Sigma_k = \lambda_k D_k A_k D_k^T \rightarrow$ Zerlegung der Kovarianzmatrix

λ_k : Volumen

A_k : Diagonalmatrix mit Determinante 1 → shape

D_k : Orthogonalmatrix → orientation