

# Data Mining Tutorial

## Clusteranalyse – Teil I

Erich Schubert, Arthur Zimek

Ludwig-Maximilians-Universität München

2013-05-03 — KDD Übung

Gegeben eine Pseudo-Metrik  $d$  auf der Menge  $A$ :

$$d : A \times A \rightarrow \mathbb{R}_0^+.$$

Definiere die Äquivalenzrelation  $\sim$  so dass

$$x \sim y \Leftrightarrow d(x, y) = 0.$$

Sei  $A^\sim$  die Menge der Äquivalenzklassen von  $A$  bzgl.  $\sim$ .

$$d^\sim : A^\sim \times A^\sim \rightarrow \mathbb{R}_0^+ \\ \text{mit } d^\sim(x^\sim, y^\sim) := d(x, y)$$

Eigenschaften?

Gegeben eine Pseudo-Metrik  $d$  auf der Menge  $A$ :  
 $d : A \times A \rightarrow \mathbb{R}_0^+$ .

Definiere die Äquivalenzrelation  $\sim$  so dass  
 $x \sim y \Leftrightarrow d(x, y) = 0$ .

Sei  $A^\sim$  die Menge der Äquivalenzklassen von  $A$  bzgl.  $\sim$ .

$$d^\sim : A^\sim \times A^\sim \rightarrow \mathbb{R}_0^+ \\ \text{mit } d^\sim(x^\sim, y^\sim) := d(x, y)$$

Eigenschaften? Wohldefiniert?

## Aufgabe 2-1

## Induzierte Metrik

## Beispiel

## Aufgabe 2-2

## Aufgabe 2-3

## Lloyd/Forgy

## MacQueen

## MacQueen Alternativ

## Qualität

## Fazit

Zu zeigen: für alle  $z \in x^\sim$ ,  $w \in y^\sim$  gilt  $d(z, w) = d(x, y)$ .

Zu zeigen: für alle  $z \in x^\sim$ ,  $w \in y^\sim$  gilt  $d(z, w) = d(x, y)$ .

Da  $z \in x^\sim$  und  $w \in y^\sim$  haben wir per Definition:

$$\begin{aligned}z^\sim &= x^\sim \text{ und } d(z, x) = 0 \\w^\sim &= y^\sim \text{ und } d(w, y) = 0\end{aligned}$$

Zu zeigen: für alle  $z \in x^\sim$ ,  $w \in y^\sim$  gilt  $d(z, w) = d(x, y)$ .

Da  $z \in x^\sim$  und  $w \in y^\sim$  haben wir per Definition:

$$z^\sim = x^\sim \text{ und } d(z, x) = 0$$

$$w^\sim = y^\sim \text{ und } d(w, y) = 0$$

Durch 1× anwenden Dreiecksungleichung erhalten wir:

$$d(z, w) \leq d(z, x) + d(x, w)$$

$$d(x, y) \leq d(x, z) + d(z, y)$$

Zu zeigen: für alle  $z \in x^\sim$ ,  $w \in y^\sim$  gilt  $d(z, w) = d(x, y)$ .

Da  $z \in x^\sim$  und  $w \in y^\sim$  haben wir per Definition:

$$\begin{aligned}z^\sim &= x^\sim \text{ und } d(z, x) = 0 \\w^\sim &= y^\sim \text{ und } d(w, y) = 0\end{aligned}$$

Durch  $2 \times$  anwenden Dreiecksungleichung erhalten wir:

$$\begin{aligned}d(z, w) &\leq d(z, x) + d(x, y) + d(y, w) \\d(x, y) &\leq d(x, z) + d(z, w) + d(w, y)\end{aligned}$$

Zu zeigen: für alle  $z \in x^\sim$ ,  $w \in y^\sim$  gilt  $d(z, w) = d(x, y)$ .

Da  $z \in x^\sim$  und  $w \in y^\sim$  haben wir per Definition:

$$\begin{aligned}z^\sim &= x^\sim \text{ und } d(z, x) = 0 \\w^\sim &= y^\sim \text{ und } d(w, y) = 0\end{aligned}$$

Durch  $2 \times$  anwenden Dreiecksungleichung erhalten wir:

$$\begin{aligned}d(z, w) &\leq d(z, x) + d(x, y) + d(y, w) = d(x, y) \\d(x, y) &\leq d(x, z) + d(z, w) + d(w, y) = d(z, w)\end{aligned}$$

Zu zeigen: für alle  $z \in x^\sim$ ,  $w \in y^\sim$  gilt  $d(z, w) = d(x, y)$ .

Da  $z \in x^\sim$  und  $w \in y^\sim$  haben wir per Definition:

$$\begin{aligned}z^\sim &= x^\sim \text{ und } d(z, x) = 0 \\w^\sim &= y^\sim \text{ und } d(w, y) = 0\end{aligned}$$

Durch  $2 \times$  anwenden Dreiecksungleichung erhalten wir:

$$\begin{aligned}d(z, w) &\leq d(z, x) + d(x, y) + d(y, w) = d(x, y) \\d(x, y) &\leq d(x, z) + d(z, w) + d(w, y) = d(z, w)\end{aligned}$$

$\Rightarrow$  Distanzberechnungen auf den Äquivalenzklassen wohldefiniert: egal welchen *Repräsentanten* wir wählen, es kommt das gleiche Ergebnis für  $d^\sim$  heraus.

## Aufgabe 2-1

## Induzierte Metrik

## Beispiel

## Aufgabe 2-2

## Aufgabe 2-3

## Lloyd/Forgey

## MacQueen

## MacQueen Alternativ

## Qualität

## Fazit

Zu zeigen (Striktheit):  $d^{\sim}(a^{\sim}, b^{\sim}) = 0 \Rightarrow a^{\sim} = b^{\sim}$

Zu zeigen (Striktheit):  $d^{\sim}(a^{\sim}, b^{\sim}) = 0 \Rightarrow a^{\sim} = b^{\sim}$

$$d^{\sim}(a^{\sim}, b^{\sim}) = 0$$

$$\Rightarrow d(a, b) = 0$$

$$\Rightarrow a \sim b$$

$$\Rightarrow a^{\sim} = b^{\sim}$$

Zu zeigen (Striktheit):  $d^{\sim}(a^{\sim}, b^{\sim}) = 0 \Rightarrow a^{\sim} = b^{\sim}$

$$d^{\sim}(a^{\sim}, b^{\sim}) = 0$$

$$\Rightarrow d(a, b) = 0$$

$$\Rightarrow a \sim b$$

$$\Rightarrow a^{\sim} = b^{\sim}$$

Reflexivität, Symmetrie und Dreiecksungleichung folgen trivial aus den Eigenschaften von  $d$ !

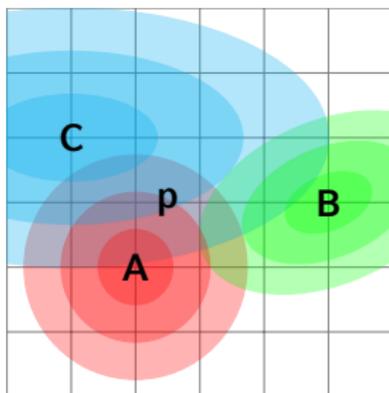
$$\text{euclid}_{xy}((r_1, x_1, y_1), (r_2, x_2, y_2)) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

Datensatz ID	x	y	Datensatz ID	x	y
1	0	1	4	1	1
2	1	1	5	2	2
3	0	1	6	3	3

$$\text{euclid}_{xy}((r_1, x_1, y_1), (r_2, x_2, y_2)) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

Datensatz ID	$x$	$y$	Datensatz ID	$x$	$y$
1	0	1	4	1	1
2	1	1	5	2	2
3	0	1	6	3	3

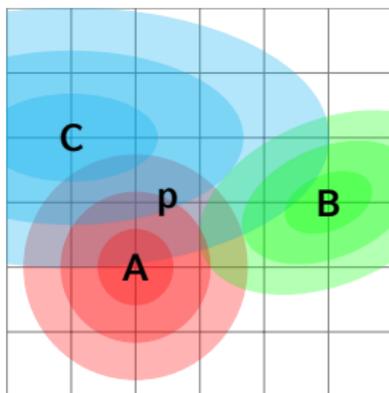
Euklidische Distanz auf  $X \times Y$ . Metrik auf  $\mathbb{R}^2 \sim X \times Y$ ,  
aber nur eine Pseudo-metric auf Datensatz ID  $\times X \times Y$ .  
"Duplikate" haben eine Distanz von 0!



Cluster können sein:

- ▶ Kugelförmig (A)
- ▶ Ellipsoid (C)
- ▶ Rotierter Ellipsoid (B)

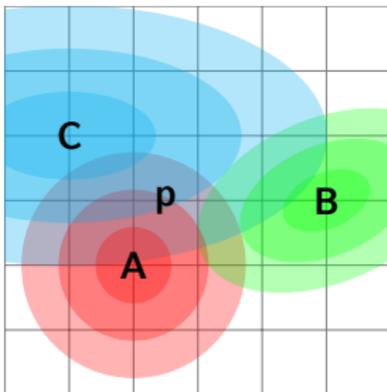
Aber: es ist immer die gleiche Formel!



Cluster können sein:

- ▶ Kugelförmig (A)
- ▶ Ellipsoid (C)
- ▶ Rotierter Ellipsoid (B)

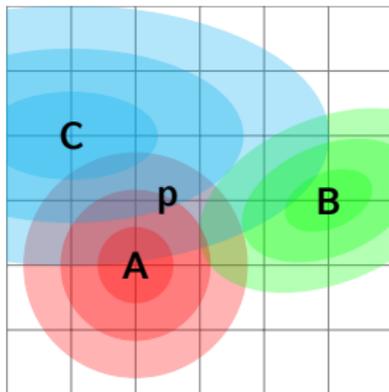
Aber: es ist immer die gleiche Formel!



Wahrscheinlichkeitsdichte (PDF) der multivariaten Normalverteilung:

$$pdf(p, \mu, \Sigma) := \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \cdot e^{-\frac{1}{2}((p-\mu)^T \Sigma^{-1} (p-\mu))}$$

Die wichtigste mehrdimensionale Verteilungsfunktion!



Wahrscheinlichkeitsdichte (PDF) der multivariaten Normalverteilung:

$$pdf(p, \mu, \Sigma) := \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \cdot e^{-\frac{1}{2}((p-\mu)^T \Sigma^{-1} (p-\mu))}$$

Die wichtigste mehrdimensionale  
Verteilungsfunktion!

Das sollten wir uns genauer anschauen!

## Multivariate Normalverteilung:

$$pdf(p, \mu, \Sigma) := \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \cdot e^{-\frac{1}{2}((p-\mu)^T \Sigma^{-1} (p-\mu))}$$

Multivariate Normalverteilung:

$$pdf(p, \mu, \Sigma) := \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \cdot e^{-\frac{1}{2}((p-\mu)^T \Sigma^{-1} (p-\mu))}$$

1-dimensionale Normalverteilung

$$pdf(x, \mu, \sigma) := \frac{1}{\sqrt{(2\pi)\sigma}} \cdot e^{-\frac{1}{2}\left(\frac{(x-\mu)}{\sigma}\right)^2}$$

Multivariate Normalverteilung:

$$pdf(p, \mu, \Sigma) := \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \cdot e^{-\frac{1}{2}((p-\mu)^T \Sigma^{-1} (p-\mu))}$$

1-dimensionale Normalverteilung

$$pdf(x, \mu, \sigma) := \frac{1}{\sqrt{(2\pi)\sigma}} \cdot e^{-\frac{1}{2}\left(\frac{(x-\mu)}{\sigma}\right)^2}$$

Multivariate Normalverteilung:

$$pdf(p, \mu, \Sigma) := \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \cdot e^{-\frac{1}{2}((p-\mu)^T \Sigma^{-1} (p-\mu))}$$

1-dimensionale Normalverteilung

$$pdf(x, \mu, \sigma) := \frac{1}{\sqrt{(2\pi)\sigma}} \cdot e^{-\frac{1}{2}\left(\frac{(x-\mu)}{\sigma}\right)^2}$$

**Normalisierung** (auf Gesamtvolumen 1!) und **quadrierte Abweichung vom Mittelwert**

Multivariate Normalverteilung:

$$pdf(p, \mu, \Sigma) := \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \cdot e^{-\frac{1}{2}((p-\mu)^T \Sigma^{-1} (p-\mu))}$$

1-dimensionale Normalverteilung

$$pdf(x, \mu, \sigma) := \frac{1}{\sqrt{(2\pi)\sigma^2}} \cdot e^{-\frac{1}{2}((x-\mu)\sigma^{-2}(x-\mu))}$$

**Normalisierung** (auf Gesamtvolumen 1!) und **quadrierte Abweichung vom Mittelwert**

Multivariate Normalverteilung:

$$pdf(p, \mu, \Sigma) := \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \cdot e^{-\frac{1}{2}((p-\mu)^T \Sigma^{-1} (p-\mu))}$$

1-dimensionale Normalverteilung

$$pdf(x, \mu, \sigma) := \frac{1}{\sqrt{(2\pi)\sigma^2}} \cdot e^{-\frac{1}{2}((x-\mu)\sigma^{-2}(x-\mu))}$$

Mahalanobis-Distanz:

$$d_{Mahalanobis}(x, \mu, \Sigma) := \sqrt{(x - \mu)^T \Sigma^{-1} (x - \mu)}$$

Multivariate Normalverteilung:

$$pdf(p, \mu, \Sigma) := \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \cdot e^{-\frac{1}{2}((p-\mu)^T \Sigma^{-1} (p-\mu))}$$

1-dimensionale Normalverteilung

$$pdf(x, \mu, \sigma) := \frac{1}{\sqrt{(2\pi)\sigma^2}} \cdot e^{-\frac{1}{2}((x-\mu)\sigma^{-2}(x-\mu))}$$

Mahalanobis-Distanz:

$$d_{Mahalanobis}(x, \mu, \Sigma)^2 := (x - \mu)^T \Sigma^{-1} (x - \mu)$$

Multivariate Normalverteilung:

$$pdf(p, \mu, \Sigma) := \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \cdot e^{-\frac{1}{2}((p-\mu)^T \Sigma^{-1} (p-\mu))}$$

1-dimensionale Normalverteilung

$$pdf(x, \mu, \sigma) := \frac{1}{\sqrt{(2\pi)\sigma^2}} \cdot e^{-\frac{1}{2}((x-\mu)\sigma^{-2}(x-\mu))}$$

Mahalanobis-Distanz:

$$d_{Mahalanobis}(x, \mu, \Sigma)^2 := (x - \mu)^T \Sigma^{-1} (x - \mu)$$

Was ist die Rolle von  $\Sigma^{-1}$ ?

Kovarianzmatrizen sind symmetrisch und auf der Diagonalen nicht negativ, und können daher normalerweise invertiert werden (es gibt degenerierte Fälle, auch die kann man handhaben!)

Kovarianzmatrizen sind symmetrisch und auf der Diagonalen nicht negativ, und können daher normalerweise invertiert werden (es gibt degenerierte Fälle, auch die kann man handhaben!)

Die Matrix kann zerlegt werden:

$$V\Lambda V^{-1} = \Sigma \quad \equiv \quad V\Lambda^{-1}V^{-1} = \Sigma^{-1}$$

wobei  $V$  die Eigenvektoren und  $\Lambda$  die Eigenwerte enthält.

Kovarianzmatrizen sind symmetrisch und auf der Diagonalen nicht negativ, und können daher normalerweise invertiert werden (es gibt degenerierte Fälle, auch die kann man handhaben!)  
Die Matrix kann zerlegt werden:

$$V\Lambda V^{-1} = \Sigma \quad \equiv \quad V\Lambda^{-1}V^{-1} = \Sigma^{-1}$$

wobei  $V$  die Eigenvektoren und  $\Lambda$  die Eigenwerte enthält.  
 $V \cong$  Drehung,  $\Lambda \cong$  Skalierung<sup>2</sup>!  
(Das ist die Kernidee der Hauptachsentransformation=PCA)

Konstruiere  $\Omega$  als  $\omega_i = 1/\sqrt{\lambda_i} = \lambda_i^{-\frac{1}{2}}$ . Dann gilt  $\Omega\Omega = \Lambda^{-1}$ .

Konstruiere  $\Omega$  als  $\omega_i = 1/\sqrt{\lambda_i} = \lambda_i^{-\frac{1}{2}}$ . Dann gilt  $\Omega\Omega = \Lambda^{-1}$ .

$$\Sigma^{-1} = V\Lambda^{-1}V^{-1} = V\Omega\Omega^T V^T = V\Omega(V\Omega)^T$$

Konstruiere  $\Omega$  als  $\omega_i = 1/\sqrt{\lambda_i} = \lambda_i^{-\frac{1}{2}}$ . Dann gilt  $\Omega\Omega = \Lambda^{-1}$ .

$$\Sigma^{-1} = V\Lambda^{-1}V^{-1} = V\Omega\Omega^T V^T = V\Omega(V\Omega)^T$$

$$d_{\text{Mahalanobis}}^2 = (x - \mu)^T \Sigma^{-1} (x - \mu)$$

Konstruiere  $\Omega$  als  $\omega_i = 1/\sqrt{\lambda_i} = \lambda_i^{-\frac{1}{2}}$ . Dann gilt  $\Omega\Omega = \Lambda^{-1}$ .

$$\Sigma^{-1} = V\Lambda^{-1}V^{-1} = V\Omega\Omega^T V^T = V\Omega(V\Omega)^T$$

$$d_{\text{Mahalanobis}}^2 = (x - \mu)^T V\Omega(V\Omega)^T (x - \mu)$$

Konstruiere  $\Omega$  als  $\omega_i = 1/\sqrt{\lambda_i} = \lambda_i^{-\frac{1}{2}}$ . Dann gilt  $\Omega\Omega = \Lambda^{-1}$ .

$$\Sigma^{-1} = V\Lambda^{-1}V^{-1} = V\Omega\Omega^T V^T = V\Omega(V\Omega)^T$$

$$\begin{aligned}d_{\text{Mahalanobis}}^2 &= (x - \mu)^T V\Omega(V\Omega)^T (x - \mu) \\ &= \langle (V\Omega)^T (x - \mu), (V\Omega)^T (x - \mu) \rangle\end{aligned}$$

Konstruiere  $\Omega$  als  $\omega_i = 1/\sqrt{\lambda_i} = \lambda_i^{-\frac{1}{2}}$ . Dann gilt  $\Omega\Omega = \Lambda^{-1}$ .

$$\Sigma^{-1} = V\Lambda^{-1}V^{-1} = V\Omega\Omega^T V^T = V\Omega(V\Omega)^T$$

$$\begin{aligned}d_{\text{Mahalanobis}}^2 &= (x - \mu)^T V\Omega(V\Omega)^T (x - \mu) \\ &= \langle (V\Omega)^T (x - \mu), (V\Omega)^T (x - \mu) \rangle \\ &= L_2((V\Omega)^T (x - \mu))^2\end{aligned}$$

Konstruiere  $\Omega$  als  $\omega_i = 1/\sqrt{\lambda_i} = \lambda_i^{-\frac{1}{2}}$ . Dann gilt  $\Omega\Omega = \Lambda^{-1}$ .

$$\Sigma^{-1} = V\Lambda^{-1}V^{-1} = V\Omega\Omega^T V^T = V\Omega(V\Omega)^T$$

$$\begin{aligned} d_{\text{Mahalanobis}}^2 &= (x - \mu)^T V\Omega(V\Omega)^T (x - \mu) \\ &= \langle (V\Omega)^T (x - \mu), (V\Omega)^T (x - \mu) \rangle \\ &= L_2((V\Omega)^T (x - \mu))^2 \end{aligned}$$

$L_2$  ist die  $L_2$ -Norm (Euclidische Distanz  $d(x, y) = L_2(x - y)$ !)

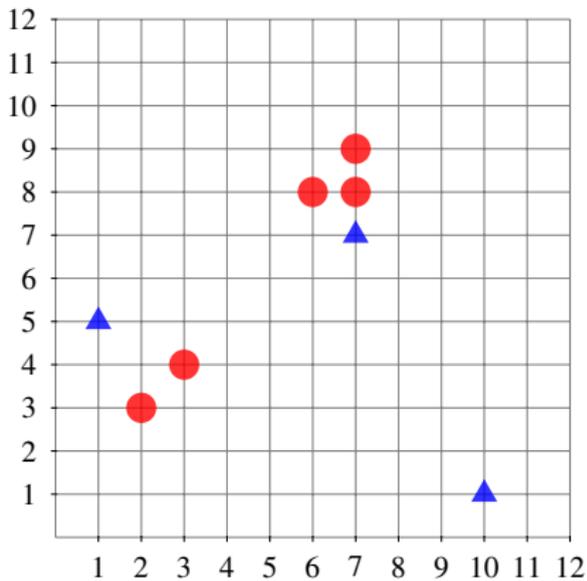
Konstruiere  $\Omega$  als  $\omega_i = 1/\sqrt{\lambda_i} = \lambda_i^{-\frac{1}{2}}$ . Dann gilt  $\Omega\Omega = \Lambda^{-1}$ .

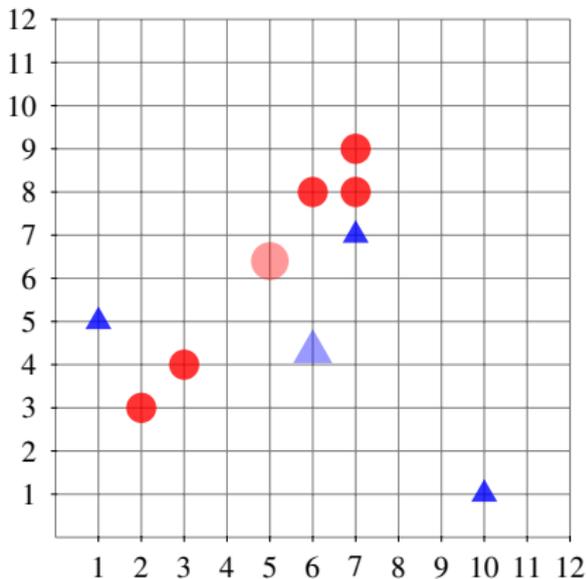
$$\Sigma^{-1} = V\Lambda^{-1}V^{-1} = V\Omega\Omega^T V^T = V\Omega(V\Omega)^T$$

$$\begin{aligned} d_{\text{Mahalanobis}}^2 &= (x - \mu)^T V\Omega(V\Omega)^T (x - \mu) \\ &= \langle (V\Omega)^T (x - \mu), (V\Omega)^T (x - \mu) \rangle \\ &= L_2((V\Omega)^T (x - \mu))^2 \end{aligned}$$

$L_2$  ist die  $L_2$ -Norm (Euclidische Distanz  $d(x, y) = L_2(x - y)$ !)

Mahalanobis  $\approx$  Euclidische Distanz nach PCA!

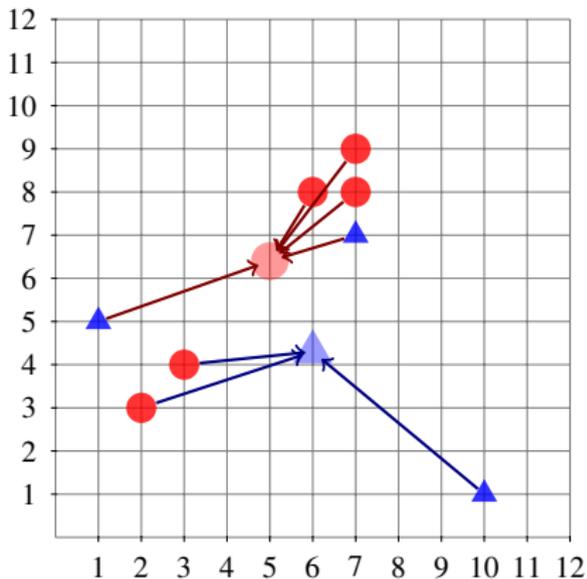




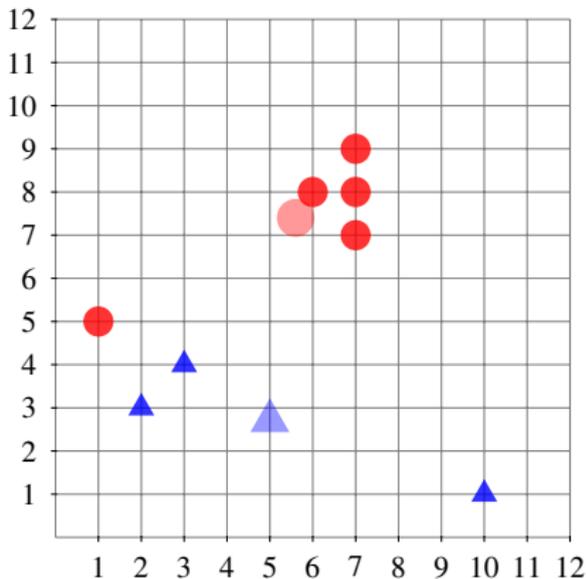
Zentroide neu berechnen:

$$\mu \approx (6.0, 4.3)$$

$$\mu \approx (5.0, 6.4)$$



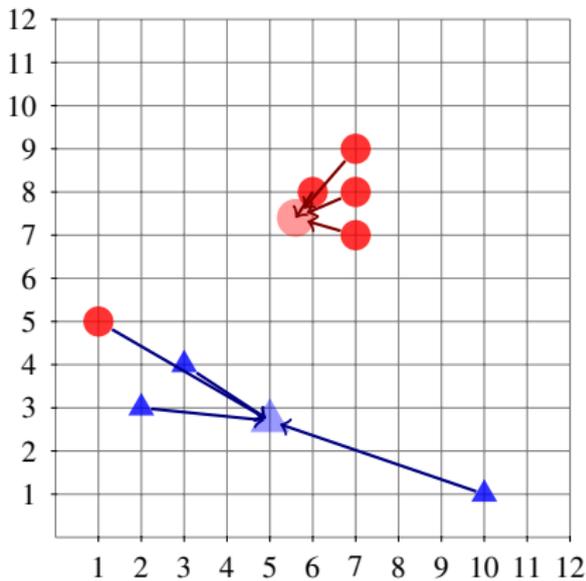
Punkte neu zuordnen



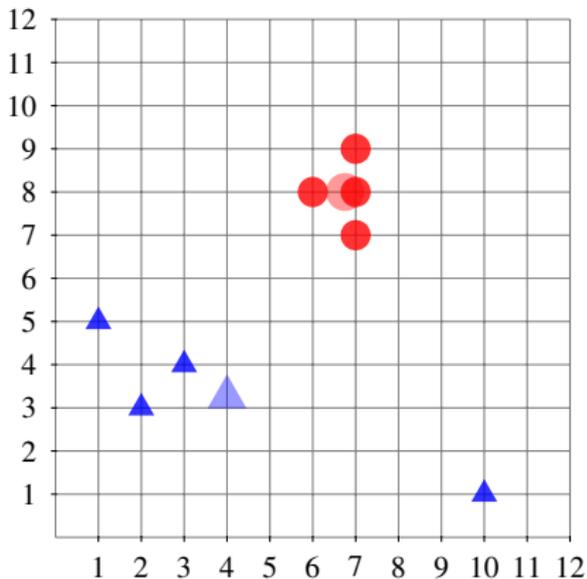
Zentroide neu berechnen:

$$\mu \approx (5.0, 2.7)$$

$$\mu \approx (5.6, 7.4)$$



Punkte neu zuordnen

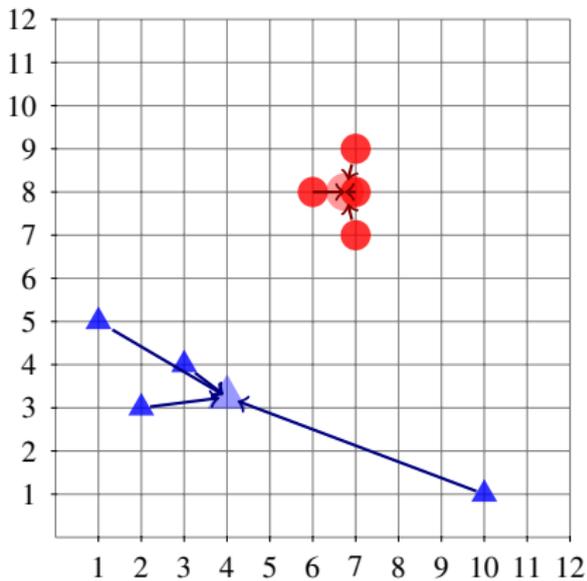


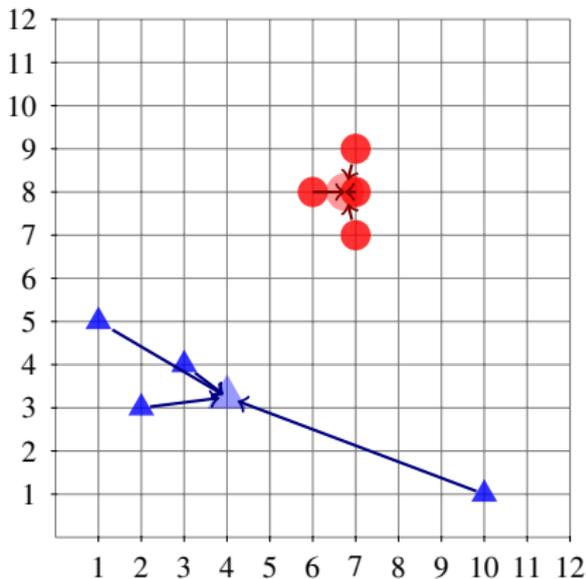
Zentroide neu berechnen:

$$\mu \approx (4.0, 3.25)$$

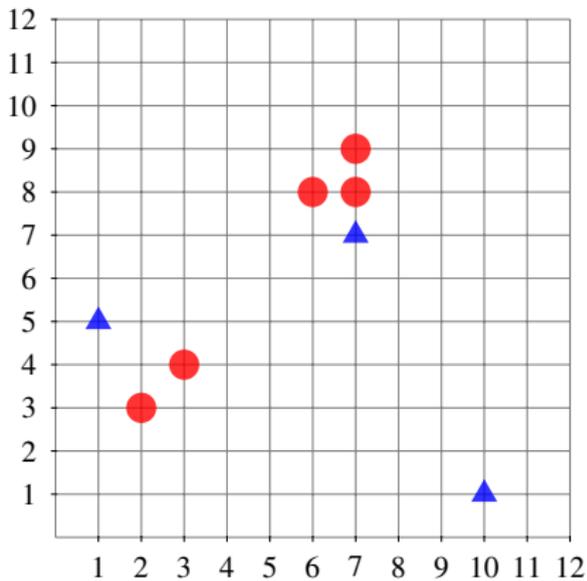
$$\mu \approx (6.75, 8.0)$$

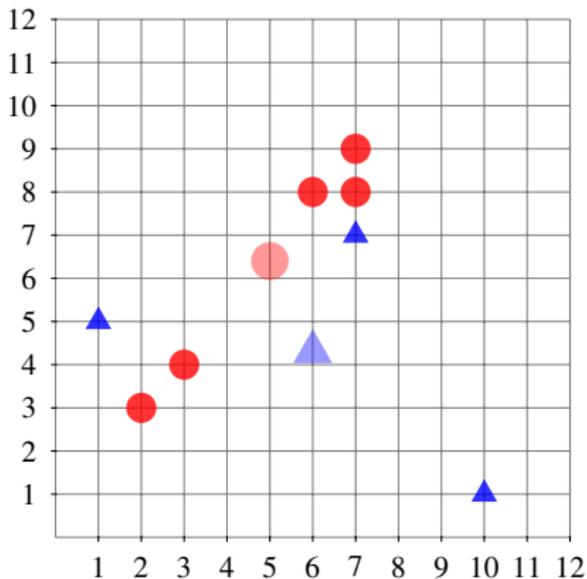
Punkte neu zuordnen





Punkte neu zuordnen  
Keine Änderung  
Konvergenz!

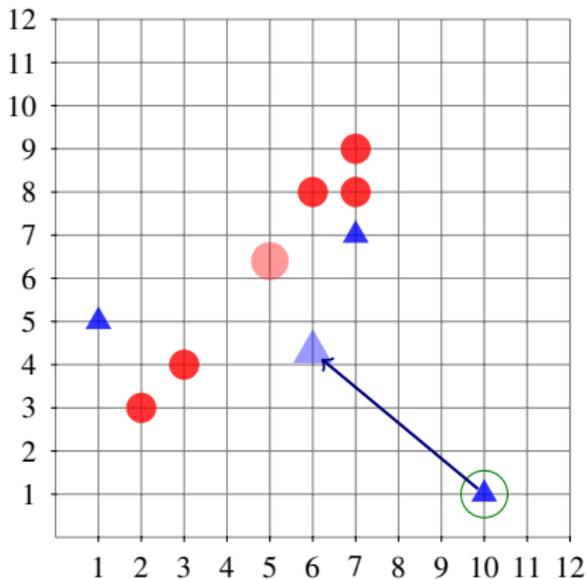




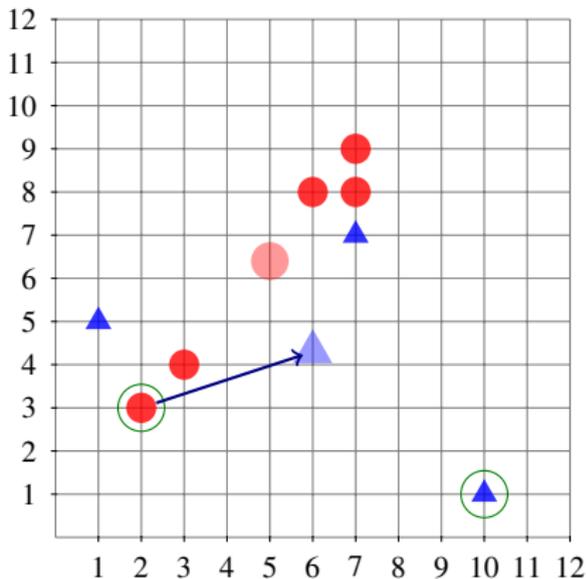
Zentroide  
(z.B.: aus  
vorheriger Iteration):

$$\mu \approx (6.0, 4.3)$$

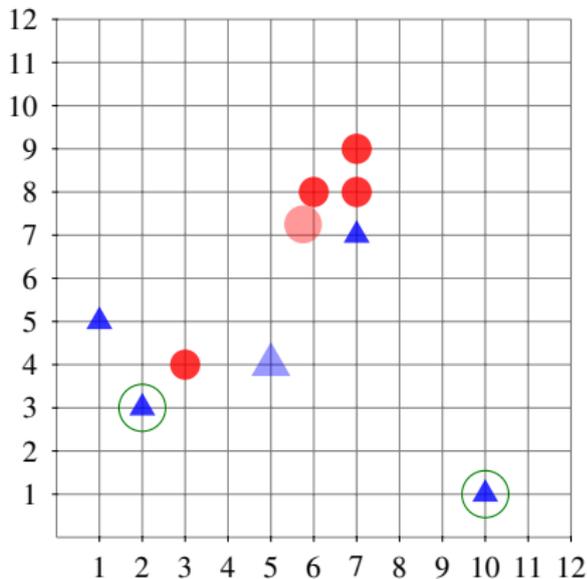
$$\mu \approx (5.0, 6.4)$$



Ersten Punkt zuordnen



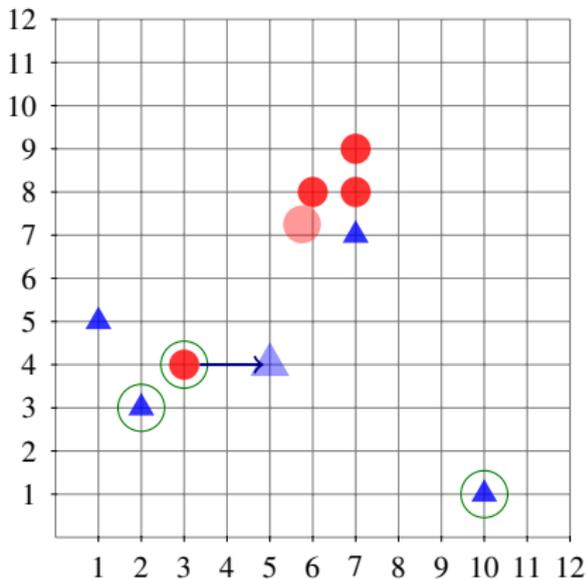
Zweiten Punkt zuordnen



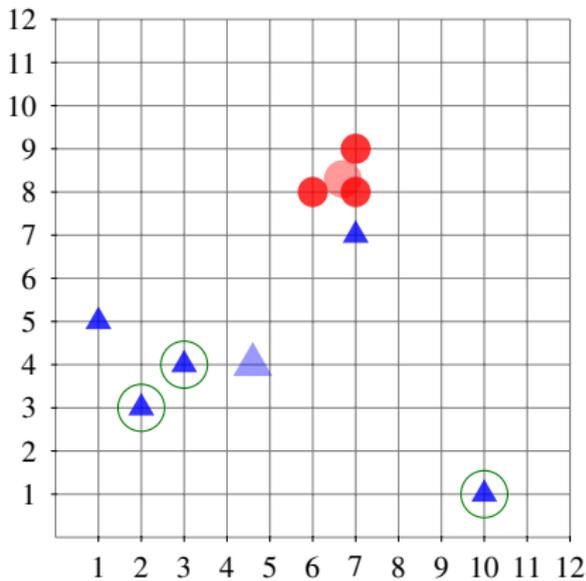
Zentroide aktualisieren:

$$\mu \approx (5.0, 4.0)$$

$$\mu \approx (5.75, 7.25)$$



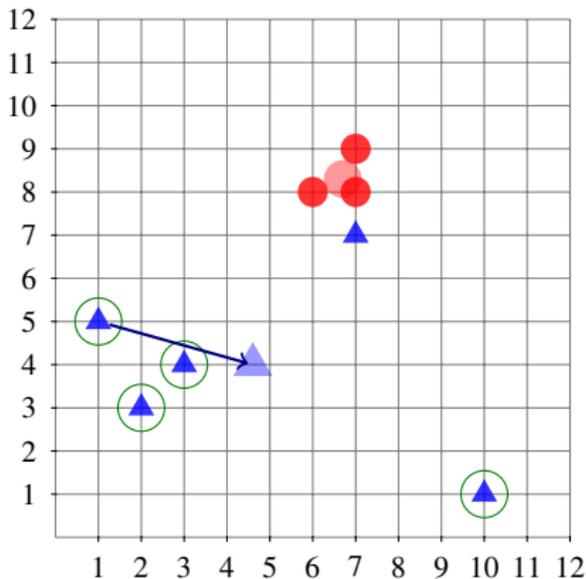
Dritten Punkt zuordnen



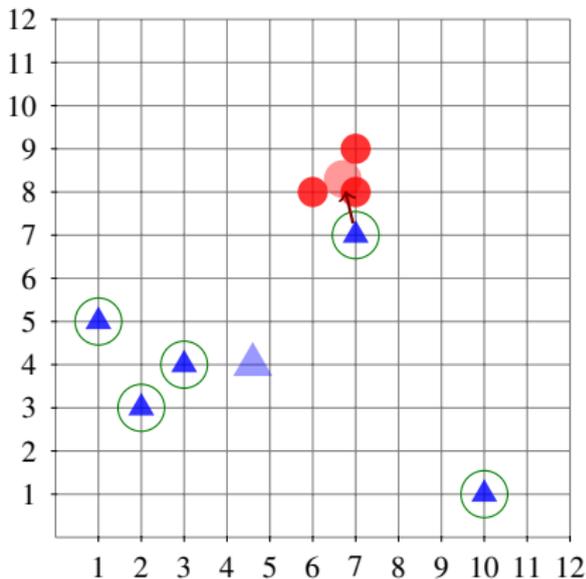
Zentroide aktualisieren:

$$\mu \approx (4.6, 4.0)$$

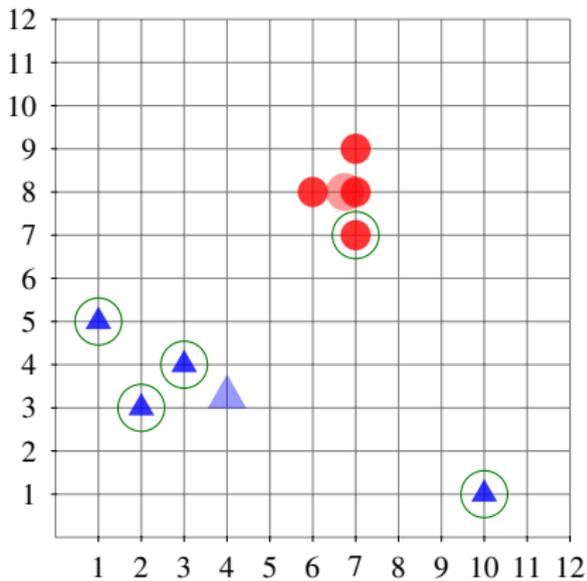
$$\mu \approx (6.7, 8.3)$$



Vierten Punkt neu  
zuordnen



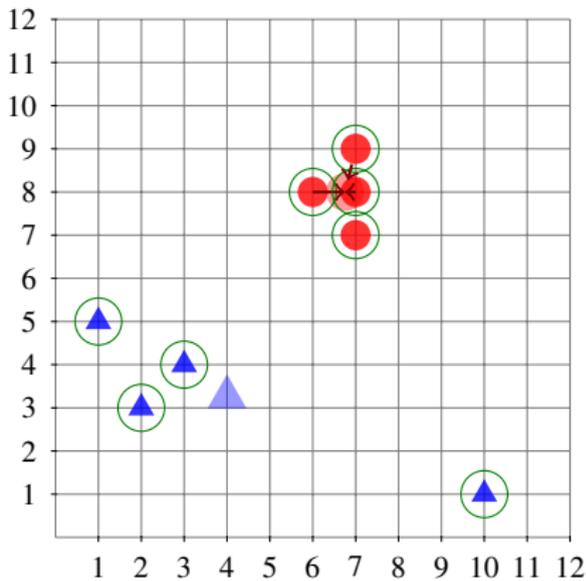
Fünften Punkt neu  
zuordnen



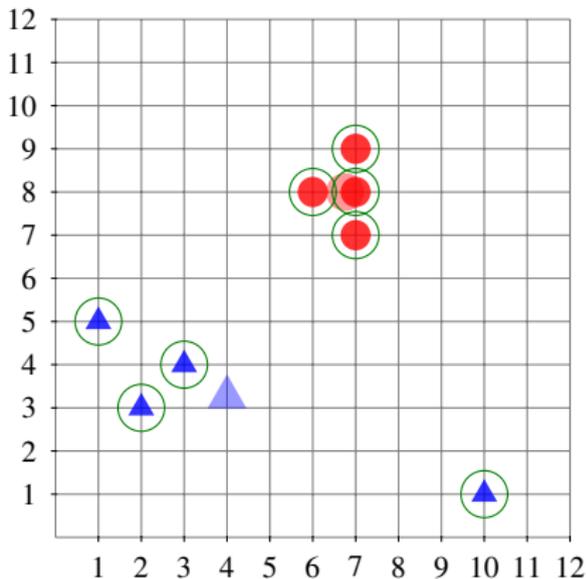
Zentroide aktualisieren:

$$\mu \approx (4.0, 3.25)$$

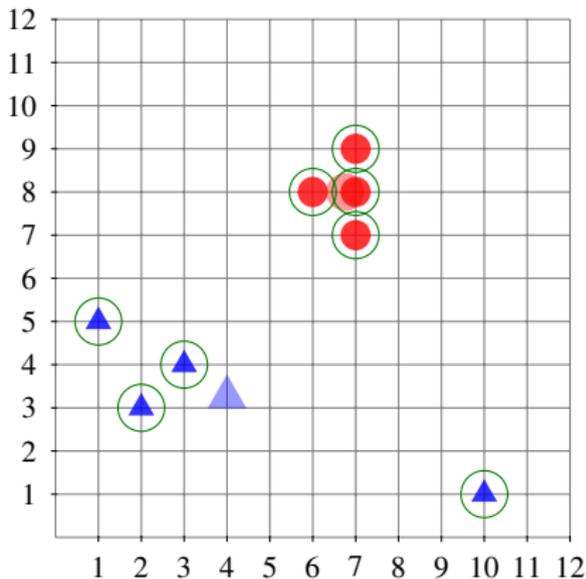
$$\mu \approx (6.75, 8.0)$$



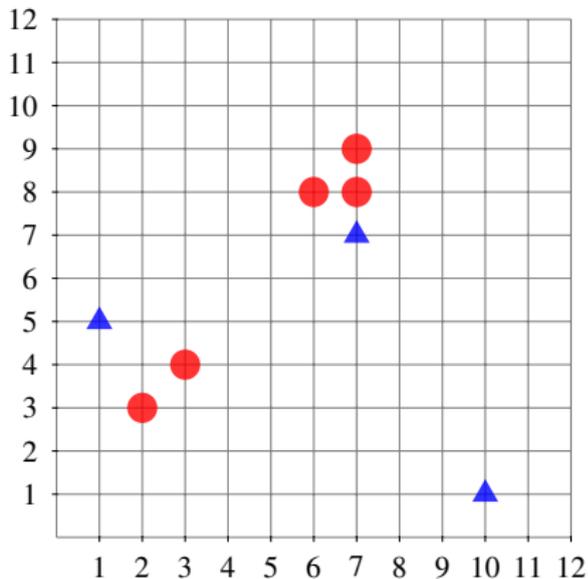
Weitere Punkte neu  
zuordnen

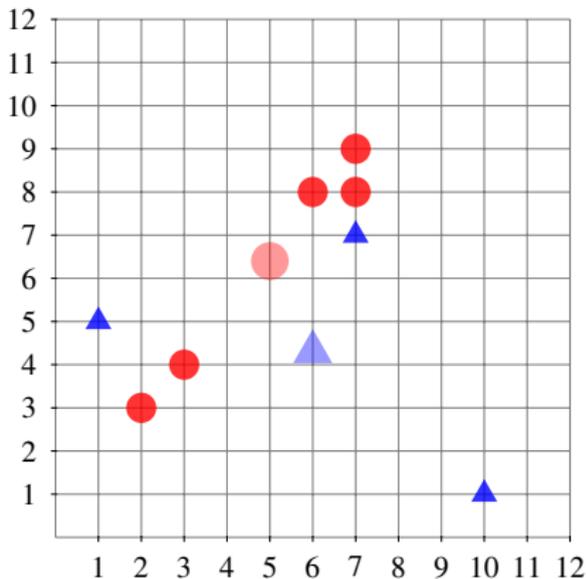


Weitere Punkte neu  
zuordnen  
ggf. Weitere Iterationen



Weitere Punkte neu  
zuordnen  
ggf. Weitere Iterationen  
Konvergenz

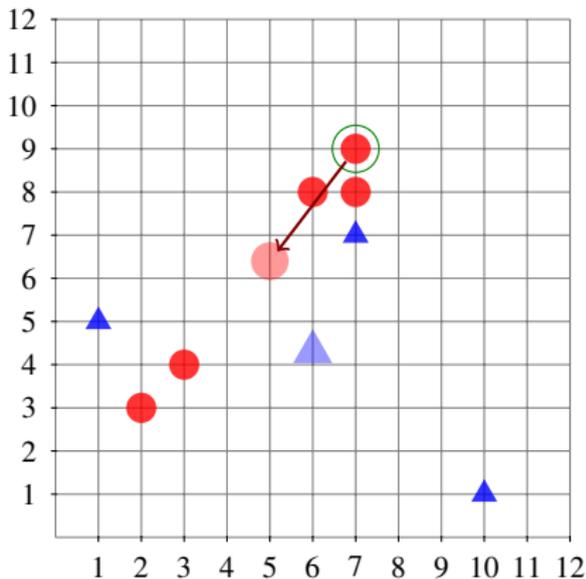




Zentroide  
(z.B.: aus  
vorheriger Iteration):

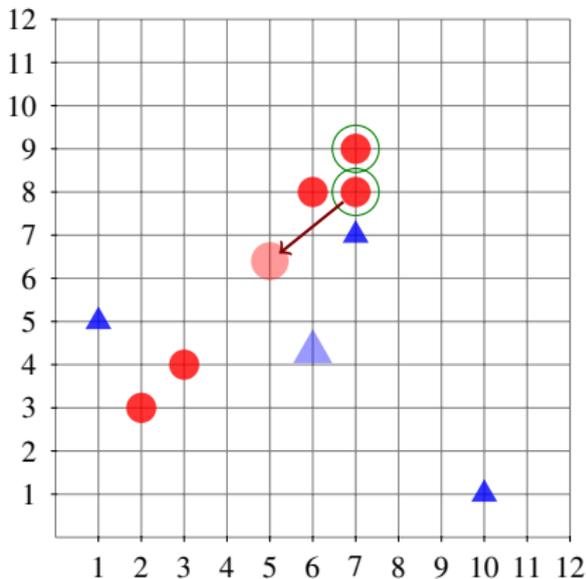
$$\mu \approx (6.0, 4.3)$$

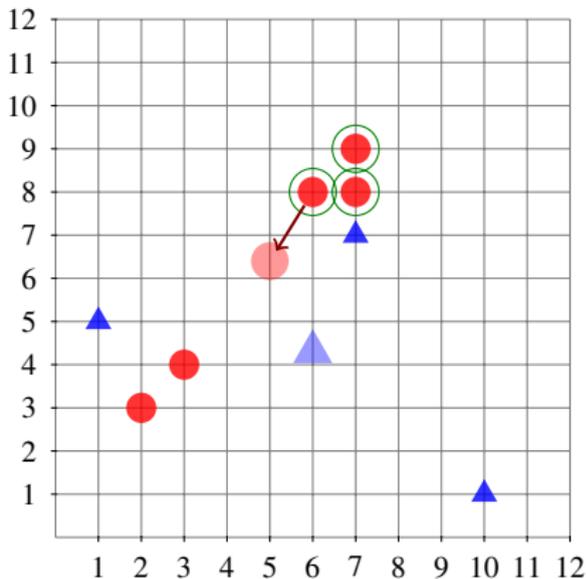
$$\mu \approx (5.0, 6.4)$$



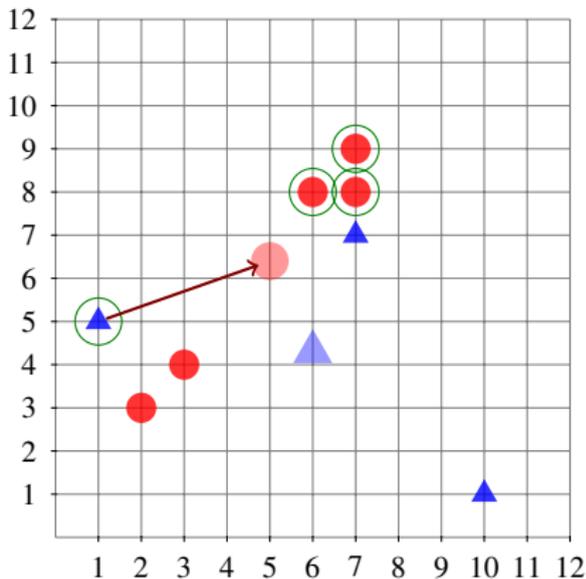
Ersten Punkt zuordnen

Zweiten Punkt zuordnen

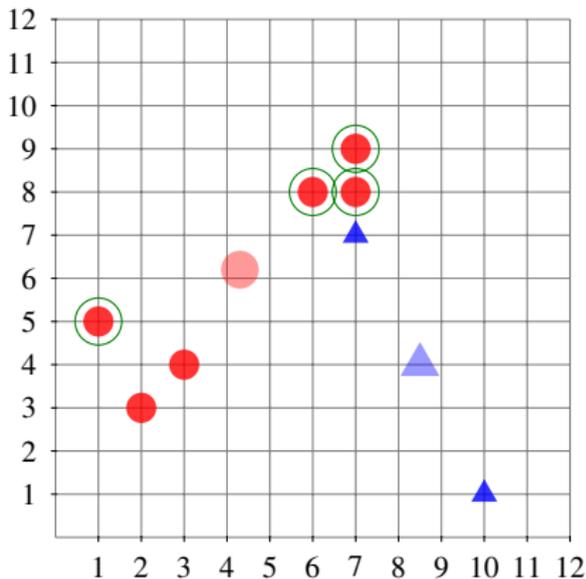




Dritten Punkt zuordnen



Vierten Punkt zuordnen



Zentroide aktualisieren:

$$\mu \approx (4.0, 8.5)$$

$$\mu \approx (4.3, 6.2)$$

# k-Means Clustering – MacQueen Algorithmus

Alternativer Ablauf – andere Reihenfolge

Data Mining  
Tutorial

E. Schubert,  
A. Zimek

Aufgabe 2-1

Induzierte Metrik  
Beispiel

Aufgabe 2-2

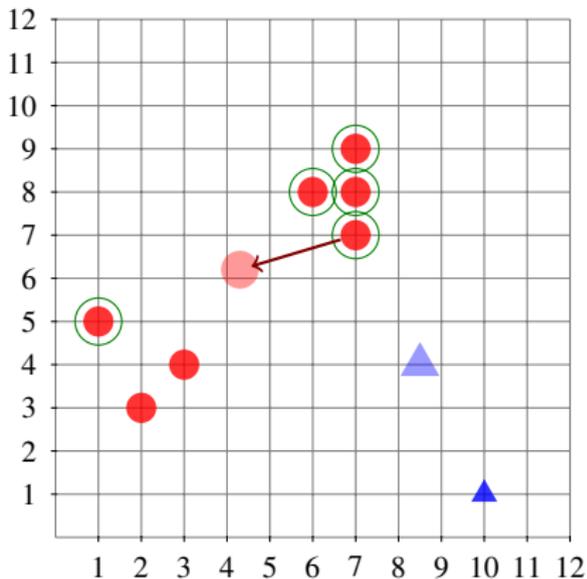
Aufgabe 2-3

Lloyd/Forgey  
MacQueen

MacQueen Alternativ

Qualität  
Fazit

Fünften Punkt zuordnen



# k-Means Clustering – MacQueen Algorithmus

Alternativer Ablauf – andere Reihenfolge

## Data Mining Tutorial

E. Schubert,  
A. Zimek

### Aufgabe 2-1

Induzierte Metrik  
Beispiel

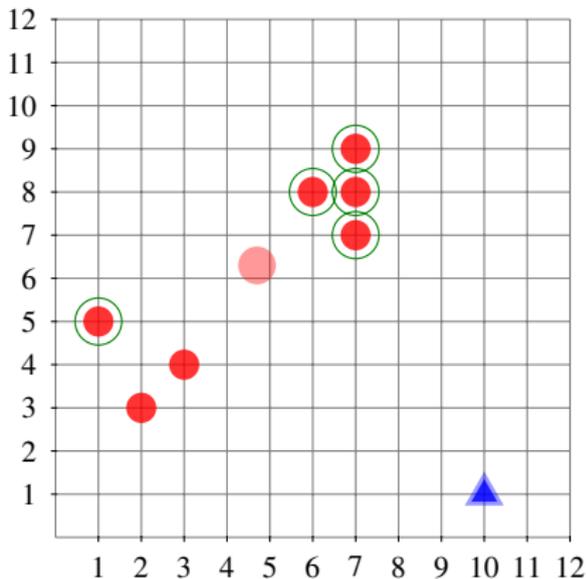
### Aufgabe 2-2

### Aufgabe 2-3

Lloyd/Forgey  
MacQueen

MacQueen Alternativ

Qualität  
Fazit



Zentroide aktualisieren:

$$\mu \approx (10.0, 1.0)$$

$$\mu \approx (4.7, 6.3)$$

# k-Means Clustering – MacQueen Algorithmus

## Alternativer Ablauf – andere Reihenfolge

Data Mining  
Tutorial

E. Schubert,  
A. Zimek

Aufgabe 2-1

Induzierte Metrik  
Beispiel

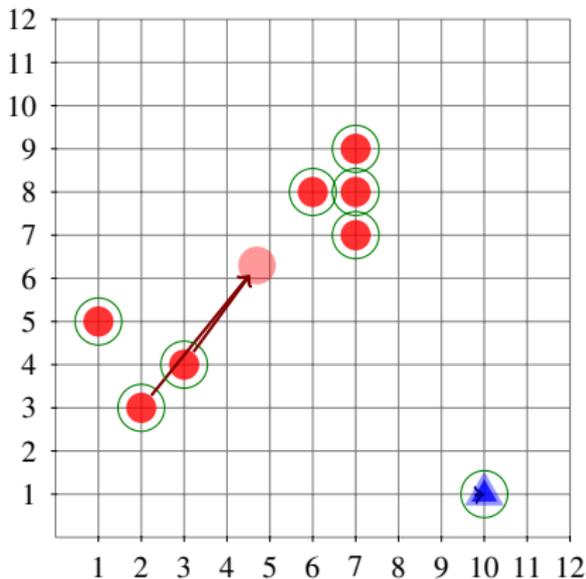
Aufgabe 2-2

Aufgabe 2-3

Lloyd/Forgey  
MacQueen

MacQueen Alternativ

Qualität  
Fazit



Weitere Punkte neu  
zuordnen

# k-Means Clustering – MacQueen Algorithmus

## Alternativer Ablauf – andere Reihenfolge

Data Mining  
Tutorial

E. Schubert,  
A. Zimek

Aufgabe 2-1

Induzierte Metrik  
Beispiel

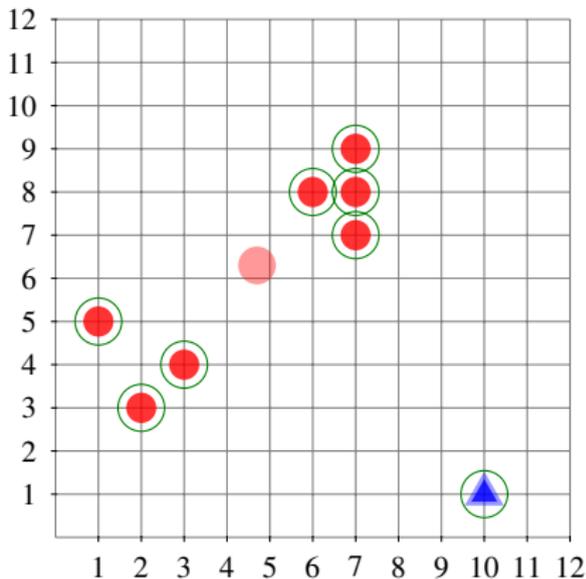
Aufgabe 2-2

Aufgabe 2-3

Lloyd/Forgy  
MacQueen

MacQueen Alternativ

Qualität  
Fazit



Weitere Punkte neu  
zuordnen  
ggf. Weitere Iterationen

# k-Means Clustering – MacQueen Algorithmus

## Alternativer Ablauf – andere Reihenfolge

Data Mining  
Tutorial

E. Schubert,  
A. Zimek

Aufgabe 2-1

Induzierte Metrik  
Beispiel

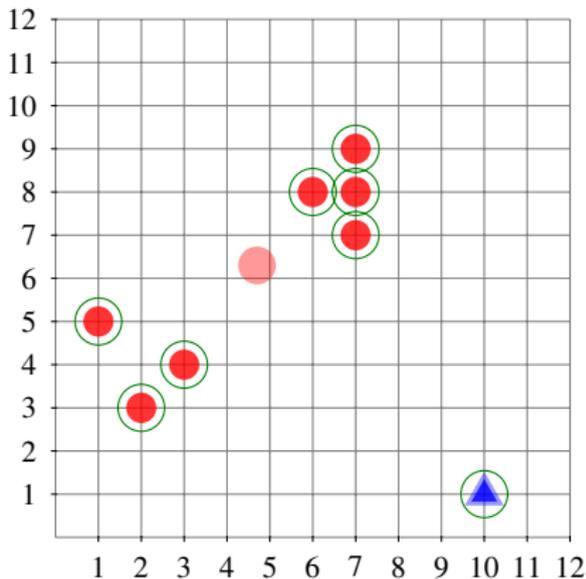
Aufgabe 2-2

Aufgabe 2-3

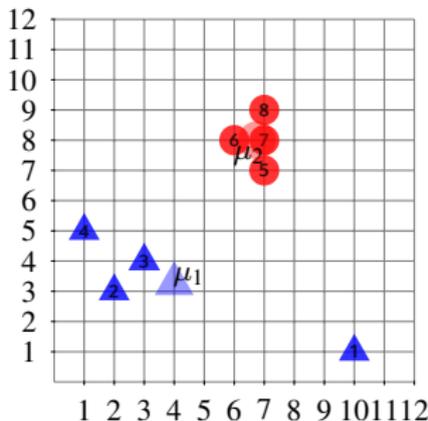
Lloyd/Forgy  
MacQueen

MacQueen Alternativ

Qualität  
Fazit



Weitere Punkte neu  
zuordnen  
ggf. Weitere Iterationen  
Konvergenz



Erste Lösung:  $TD^2 = 61\frac{1}{2}$

$$d^2(\mu_1, p_1) = |4 - 10|^2 + |3.25 - 1|^2 = 36 + 5\frac{1}{16} = 41\frac{1}{16}$$

$$d^2(\mu_1, p_2) = |4 - 2|^2 + |3.25 - 3|^2 = 4 + \frac{1}{16} = 4\frac{1}{16}$$

$$d^2(\mu_1, p_3) = |4 - 3|^2 + |3.25 - 4|^2 = 1 + \frac{9}{16} = 1\frac{9}{16}$$

$$d^2(\mu_1, p_4) = |4 - 1|^2 + |3.25 - 5|^2 = 9 + 3\frac{1}{16} = 12\frac{1}{16}$$

$$TD^2(C_1) = 58\frac{3}{4}$$

$$d^2(\mu_2, p_5) = |6.75 - 7|^2 + |8 - 7|^2 = \frac{1}{16} + 1 = 1\frac{1}{16}$$

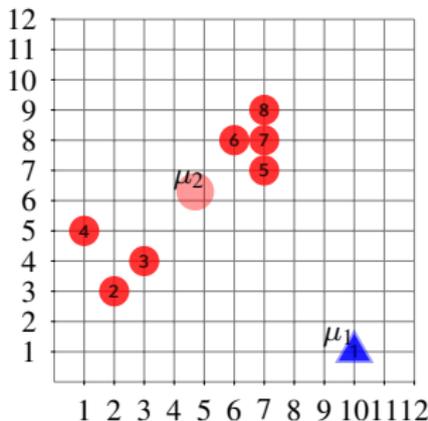
$$d^2(\mu_2, p_6) = |6.75 - 6|^2 + |8 - 8|^2 = \frac{9}{16} + 0 = \frac{9}{16}$$

$$d^2(\mu_2, p_7) = |6.75 - 7|^2 + |8 - 8|^2 = \frac{1}{16} + 0 = \frac{1}{16}$$

$$d^2(\mu_2, p_8) = |6.75 - 7|^2 + |8 - 9|^2 = \frac{1}{16} + 1 = 1\frac{1}{16}$$

$$TD^2(C_2) = 2\frac{3}{4}$$

(Hier: sum-of-squares  $\equiv$  quadrierte Euklidische Distanz – mit Manhattan<sup>2</sup> kommen andere aber ähnliche Zahlen heraus)



$$d^2(\mu_1, p_1) = |10 - 10|^2 + |1 - 1|^2 = 0$$

$$TD^2(C_1) = 0$$

$$d^2(\mu_2, p_2) \approx |4.7 - 2|^2 + |6.3 - 3|^2 \approx 18.2$$

$$d^2(\mu_2, p_3) \approx |4.7 - 3|^2 + |6.3 - 4|^2 \approx 8.2$$

$$d^2(\mu_2, p_4) \approx |4.7 - 1|^2 + |6.3 - 5|^2 \approx 15.4$$

$$d^2(\mu_2, p_5) \approx |4.7 - 7|^2 + |6.3 - 7|^2 \approx 5.7$$

$$d^2(\mu_2, p_6) \approx |4.7 - 6|^2 + |6.3 - 8|^2 \approx 4.6$$

$$d^2(\mu_2, p_7) \approx |4.7 - 7|^2 + |6.3 - 8|^2 \approx 8.2$$

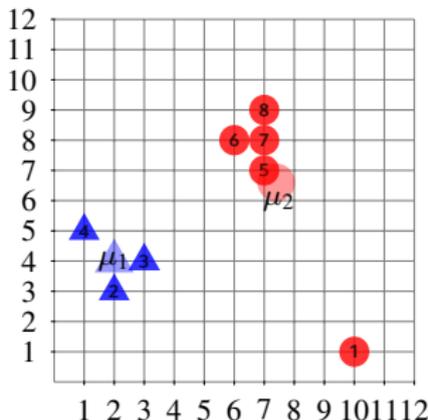
$$d^2(\mu_2, p_8) \approx |4.7 - 7|^2 + |6.3 - 9|^2 \approx 12.6$$

$$TD^2(C_2) \approx 72.86$$

Erste Lösung:  $TD^2 = 61\frac{1}{2}$

Zweite Lösung:  $TD^2 \approx 72.68$

(Hier: sum-of-squares  $\equiv$  quadrierte Euklidische Distanz – mit Manhattan<sup>2</sup> kommen andere aber ähnliche Zahlen heraus)



$$d^2(\mu_1, p_2) = |2 - 2|^2 + |4 - 3|^2 = 0 + 1 = 1$$

$$d^2(\mu_1, p_3) = |2 - 3|^2 + |4 - 4|^2 = 1 + 0 = 1$$

$$d^2(\mu_1, p_4) = |2 - 1|^2 + |4 - 5|^2 = 1 + 1 = 2$$

$$TD^2(C_1) = 4$$

$$d^2(\mu_2, p_1) = |7.4 - 10|^2 + |6.6 - 1|^2 = 6\frac{19}{25} + 31\frac{9}{25} = 38\frac{3}{25}$$

$$d^2(\mu_2, p_5) = |7.4 - 7|^2 + |6.6 - 7|^2 = \frac{4}{25} + \frac{4}{25} = \frac{8}{25}$$

$$d^2(\mu_2, p_6) = |7.4 - 6|^2 + |6.6 - 8|^2 = 1\frac{24}{25} + 1\frac{24}{25} = 3\frac{23}{25}$$

$$d^2(\mu_2, p_7) = |7.4 - 7|^2 + |6.6 - 8|^2 = \frac{4}{25} + 1\frac{24}{25} = 2\frac{3}{25}$$

$$d^2(\mu_2, p_8) = |7.4 - 7|^2 + |6.6 - 9|^2 = \frac{4}{25} + 5\frac{19}{25} = 5\frac{23}{25}$$

$$TD^2(C_2) = 50\frac{2}{5}$$

Erste Lösung:  $TD^2 = 61\frac{1}{2}$

Zweite Lösung:  $TD^2 \approx 72.68$

Optimale Lösung:  $TD^2 = 54\frac{2}{5}$

(Hier: sum-of-squares  $\equiv$  quadrierte Euklidische Distanz – mit Manhattan<sup>2</sup> kommen andere aber ähnliche Zahlen heraus)

## Merke:

- ▶ K-means konvergiert nur gegen ein lokales Minimum
- ▶ K-means ist abhängig von den Startparametern
- ▶ K-means nach MacQueen ist reihenfolgeabhängig
- ▶ K-means ist anfällig gegen Rauschen
  - ▶ Degenerierte 1-Element "Cluster"
  - ▶ Dadurch Reduktion von effektivem  $k$
- ▶ K-means minimiert Varianzen, ist also eigentlich nur für Euklidische Distanz korrekt (ggf. aber auch garantierte Konvergenz bei anderen Distanzen)
- ▶ K-means (nach Lloyd) ist dennoch das beliebteste Verfahren, da es sehr einfach und schnell ist und mit geringem Aufwand implementiert werden kann!