



Wissensentdeckung in Datenbanken

Clustering (II)

Nico Piatkowski und Uwe Ligges

Informatik—Künstliche Intelligenz
Computergestützte Statistik
Technische Universität Dortmund

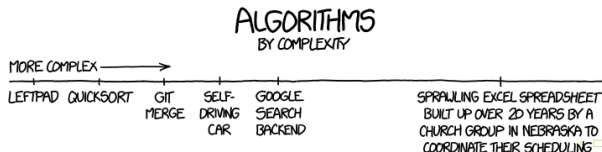
11.07.2017



Zusammenfassung: Merkmalsauswahl

Sowohl bei Regressions- als auch bei Klassifikationsproblemen, kann es helfen unwichtige Variablen **vor dem Lernen** auszuschließen.

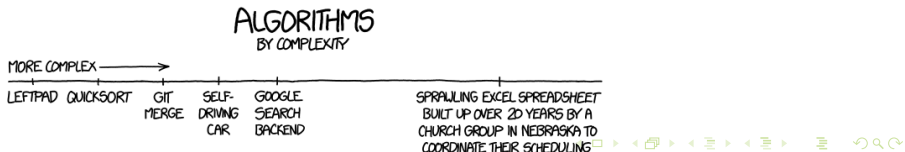
- Eine Greedy-Merkmalsauswahl verwaltet eine Menge von Variablen M
- Das Verfahren paßt die Menge M iterativ an indem Variablen (basierend auf einer **Verlustfunktion**) hinzugeügt/entfernt werden
 - **Forward-Selection:** Die Menge M ist zu Anfang leer und wird in jeder Iteration größer
 - **Backward-Selection:** Die Menge M enthält zu Anfang alle Variablen und wird in jeder Iteration kleiner



Zusammenfassung: Merkmalsauswahl (II)

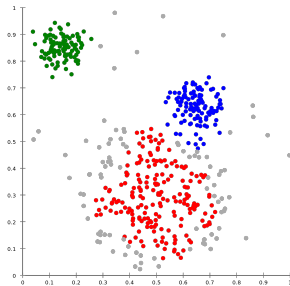
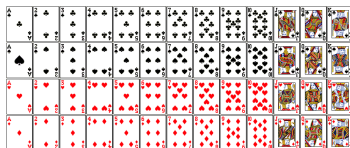
Sowohl bei Regressions- als auch bei Klassifikationsproblemen, kann es helfen unwichtige Variablen **vor dem Lernen** auszuschließen.

- Eine regularisierungsbasierte Merkmalsauswahl bevorzugt Modelle mit kleiner Norm
- Im Falle der l_1 -Norm (LASSO regression) werden Modelle bevorzugt bei denen einige Modellparameter = 0 sind
- Im Falle linearer Modelle entspricht dies direkt einer Auswahl an Variablen
- **Nachdem** die Merkmale ausgewählt wurden, wird das eigentliche Modell gelernt



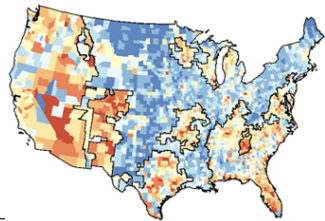
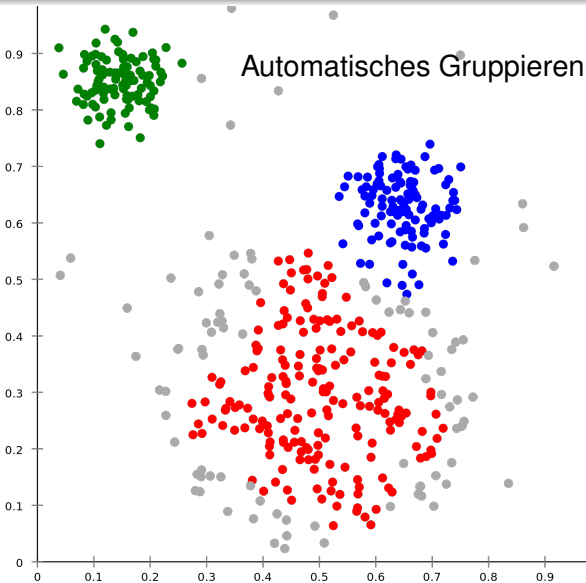
Überblick

- Clustering
 - Wiederholung: Problemstellung
 - Wiederholung: k -Means
 - DBSCAN
 - LDA



Clusteranalyse

Automatisches Gruppieren von Daten





k -Means / Lloyd's Algorithmus

Eingabe: Daten \mathcal{D} , Anzahl Cluster k , Metrik/Distanzmaß

$$f : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+ \cup \{0\}$$

- (1) Weise jedem Punkt in \mathcal{D} einen zufälligen Cluster zu
- (2) Bestimme Clusterzentrum c ("Mittelpunkt") jedes Clusters
- (3) Weise jedem Punkt x den Cluster zu, dessen Mittelpunkt c am nächsten zu x ist (mittels f)
- (4) Wiederhole Schritte 2 und 3 so lange, bis sich die Clusterzuweisung nicht mehr ändert oder Zeit aufgebraucht

Man kann zeigen: Minimiert Distanzen innerhalb der Cluster
und maximiert Distanz zwischen den Clustern



k -Means / Lloyd's Algorithmus

Eingabe: Daten \mathcal{D} , Anzahl Cluster k , Metrik/Distanzmaß

$$f : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+ \cup \{0\}$$

- (1) Weise jedem Punkt in D einen zufälligen Cluster zu
- (2) Bestimme Clusterzentrum c ("Mittelpunkt") jedes Clusters
- (3) Weise jedem Punkt x den Cluster zu, dessen Mittelpunkt c am nächsten zu x ist (mittels f)
- (4) Wiederhole Schritte 2 und 3 so lange, bis sich die Clusterzuweisung nicht mehr ändert oder Zeit aufgebraucht

Man kann zeigen: Minimiert Distanzen innerhalb der Cluster und maximiert Distanz zwischen den Clustern



k -Means / Lloyd's Algorithmus

Eingabe: Daten \mathcal{D} , Anzahl Cluster k , Metrik/Distanzmaß

$$f : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+ \cup \{0\}$$

- (1) Weise jedem Punkt in D einen zufälligen Cluster zu
- (2) Bestimme Clusterzentrum c ("Mittelpunkt") jedes Clusters
- (3) Weise jedem Punkt x den Cluster zu, dessen Mittelpunkt c am nächsten zu x ist (mittels f)
- (4) Wiederhole Schritte 2 und 3 so lange, bis sich die Clusterzuweisung nicht mehr ändert oder Zeit aufgebraucht

Man kann zeigen: Minimiert Distanzen innerhalb der Cluster
und maximiert Distanz zwischen den Clustern



k -Means / Lloyd's Algorithmus

Eingabe: Daten \mathcal{D} , Anzahl Cluster k , Metrik/Distanzmaß

$$f : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+ \cup \{0\}$$

- (1) Weise jedem Punkt in D einen zufälligen Cluster zu
- (2) Bestimme Clusterzentrum c ("Mittelpunkt") jedes Clusters
- (3) Weise jedem Punkt x den Cluster zu, dessen Mittelpunkt c am nächsten zu x ist (mittels f)
- (4) Wiederhole Schritte 2 und 3 so lange, bis sich die Clusterzuweisung nicht mehr ändert oder Zeit aufgebraucht

Man kann zeigen: Minimiert Distanzen innerhalb der Cluster
und maximiert Distanz zwischen den Clustern



k -Means / Lloyd's Algorithmus

Eingabe: Daten \mathcal{D} , Anzahl Cluster k , Metrik/Distanzmaß

$$f : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+ \cup \{0\}$$

- (1) Weise jedem Punkt in D einen zufälligen Cluster zu
- (2) Bestimme Clusterzentrum c ("Mittelpunkt") jedes Clusters
- (3) Weise jedem Punkt x den Cluster zu, dessen Mittelpunkt c am nächsten zu x ist (mittels f)
- (4) Wiederhole Schritte 2 und 3 so lange, bis sich die Clusterzuweisung nicht mehr ändert oder Zeit aufgebraucht

Man kann zeigen: Minimiert Distanzen innerhalb der Cluster
und maximiert Distanz zwischen den Clustern



k -Means / Lloyd's Algorithmus

Eingabe: Daten \mathcal{D} , Anzahl Cluster k , Metrik/Distanzmaß

$$f : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+ \cup \{0\}$$

- (1) Weise jedem Punkt in D einen zufälligen Cluster zu
- (2) Bestimme Clusterzentrum c ("Mittelpunkt") jedes Clusters
- (3) Weise jedem Punkt x den Cluster zu, dessen Mittelpunkt c am nächsten zu x ist (mittels f)
- (4) Wiederhole Schritte 2 und 3 so lange, bis sich die Clusterzuweisung nicht mehr ändert oder Zeit aufgebraucht

Man kann zeigen: Minimiert Distanzen innerhalb der Cluster und maximiert Distanz zwischen den Clustern



k -Means Verlustfunktion

Notation:

- Datensatz \mathcal{D} mit $|\mathcal{D}| = N$, n -dimensionalen Datenpunkten
- Metrik/Distanzmaß $f : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+ \cup \{0\}$
- Die Menge $C = \{\mathbf{c}^{(1)}, \mathbf{c}^{(2)}, \dots, \mathbf{c}^{(k)}\}$ enthält die k Clusterzentren

Optimierungsproblem:

$$\min_{C \subset \mathbb{R}^n, |C|=k} \ell(C; \mathcal{D}) = \min_{C \subset \mathbb{R}^n, |C|=k} \sum_{\mathbf{x} \in \mathcal{D}} \min_{\mathbf{c} \in C} f(\mathbf{x}, \mathbf{c})$$



k -Means Verlustfunktion (mit Euklidischem Abstand)

Jetzt:

- Metrik/Distanzmaß $f(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2$

Verlustfunktion:

$$\begin{aligned}\ell(C; \mathcal{D}) &= \sum_{\mathbf{x} \in \mathcal{D}} \min_{\mathbf{c} \in C} \|\mathbf{x} - \mathbf{c}\|_2^2 \\ \frac{\partial}{\partial \mathbf{c}_j^{(i)}} \ell(C; \mathcal{D}) &= \sum_{\mathbf{x} \in \mathcal{D}} \frac{\partial}{\partial \mathbf{c}_j^{(i)}} \min_{\mathbf{c} \in C} \sum_{l=1}^n (\mathbf{x}_l - \mathbf{c}_l)^2 \\ &= \sum_{\mathbf{x} \in \mathcal{D}_i} \frac{\partial}{\partial \mathbf{c}_j^{(i)}} (\mathbf{x}_j - \mathbf{c}_j^{(i)})^2 \\ &= \sum_{\mathbf{x} \in \mathcal{D}_i} 2(\mathbf{c}_j^{(i)} - \mathbf{x}_j)\end{aligned}$$

Kurzschreibweise: \mathcal{D}_i enthält Datenpunkt mit minimalem Abstand zu (Cluster) \mathbf{c}_i



k -Means Verlustfunktion (mit Euklidischem Abstand)

Jetzt:

- Metrik/Distanzmaß $f(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2$

Verlustfunktion:

$$\begin{aligned}\ell(C; \mathcal{D}) &= \sum_{\mathbf{x} \in \mathcal{D}} \min_{\mathbf{c} \in C} \|\mathbf{x} - \mathbf{c}\|_2^2 \\ \frac{\partial}{\partial \mathbf{c}_j^{(i)}} \ell(C; \mathcal{D}) &= \sum_{\mathbf{x} \in \mathcal{D}} \frac{\partial}{\partial \mathbf{c}_j^{(i)}} \min_{\mathbf{c} \in C} \sum_{l=1}^n (\mathbf{x}_l - \mathbf{c}_l)^2 \\ &= \sum_{\mathbf{x} \in \mathcal{D}_i} \frac{\partial}{\partial \mathbf{c}_j^{(i)}} (\mathbf{x}_j - \mathbf{c}_j^{(i)})^2 \\ &= \sum_{\mathbf{x} \in \mathcal{D}_i} 2(\mathbf{c}_j^{(i)} - \mathbf{x}_j)\end{aligned}$$

Kurzschreibweise: \mathcal{D}_i enthält Datenpunkt mit minimalem Abstand zu (Cluster) \mathbf{c}_i



k -Means Verlustfunktion (mit Euklidischem Abstand) (II)

Es gilt im Optimum:

$$\ell(C; \mathcal{D}) = 0 = \sum_{\mathbf{x} \in \mathcal{D}_i} 2(\mathbf{c}_j^{(i)} - \mathbf{x}_j)$$

$$\Leftrightarrow 0 = |\mathcal{D}_i| 2\mathbf{c}_j^{(i)} - 2 \sum_{\mathbf{x} \in \mathcal{D}_i} \mathbf{x}_j$$

$$\frac{1}{|\mathcal{D}_i|} \sum_{\mathbf{x} \in \mathcal{D}_i} \mathbf{x}_j = \mathbf{c}_j^{(i)}$$

Also:

- k -Means ist ein Optimierungsverfahren erster Ordnung (wie Gradientenabstieg!)
- **Aber:** Zielfunktion ist nicht-konvex. Keine Konvergenz zum globalen Optimum!



k -Means Verlustfunktion (mit Euklidischem Abstand) (II)

Es gilt im **Optimum**:

$$\ell(C; \mathcal{D}) = 0 = \sum_{\mathbf{x} \in \mathcal{D}_i} 2(\mathbf{c}_j^{(i)} - \mathbf{x}_j)$$

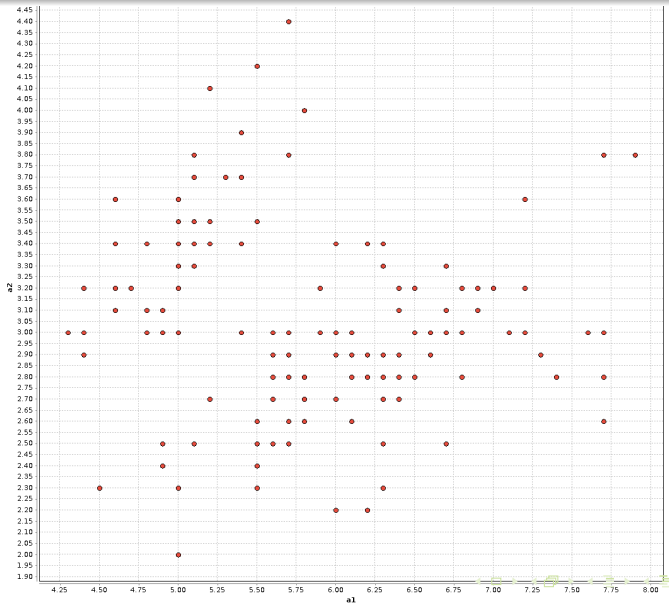
$$\Leftrightarrow 0 = |\mathcal{D}_i| 2\mathbf{c}_j^{(i)} - 2 \sum_{\mathbf{x} \in \mathcal{D}_i} \mathbf{x}_j$$

$$\frac{1}{|\mathcal{D}_i|} \sum_{\mathbf{x} \in \mathcal{D}_i} \mathbf{x}_j = \mathbf{c}_j^{(i)}$$

Also:

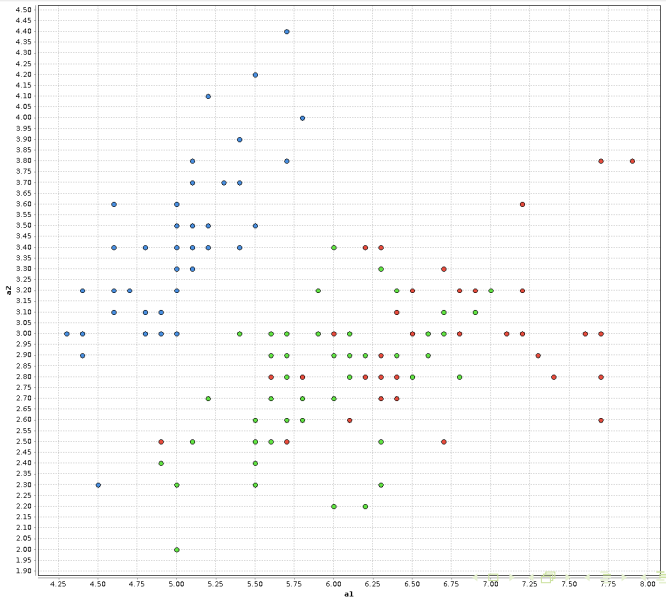
- k -Means ist ein Optimierungsverfahren erster Ordnung (wie Gradientenabstieg!)
- **Aber:** Zielfunktion ist nicht-konvex. Keine Konvergenz zum globalen Optimum!

Beispiel: Daten (x)

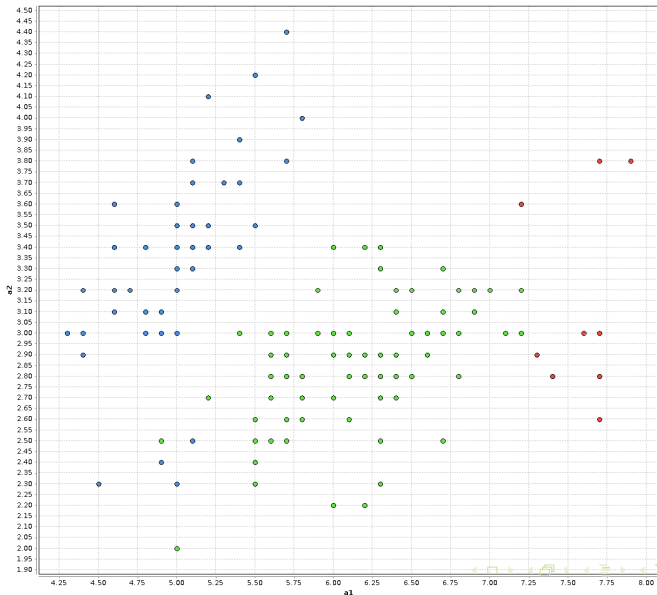




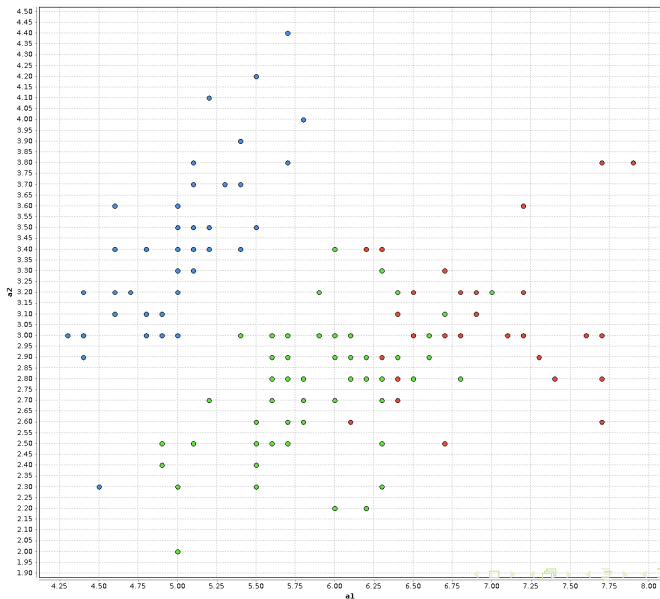
Beispiel: Daten mit Klassen (x, y)



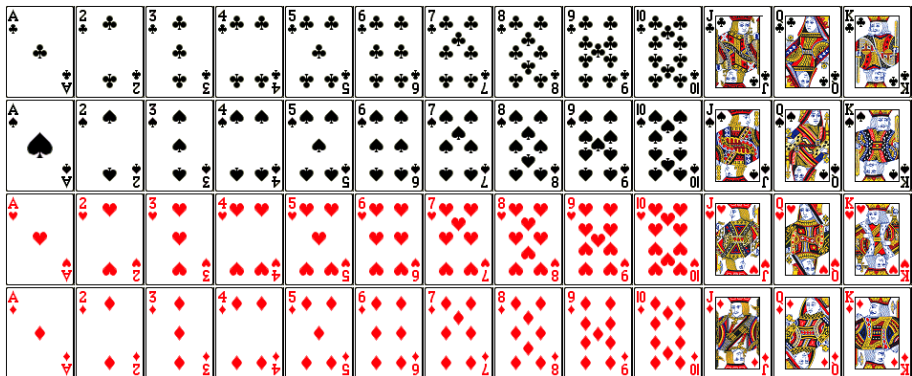
Beispiel: 3-means, 1 Iteration



Beispiel: 3-means, 10 Iterationen



Intuition: Auswahl von k und Distanzmaß (II)



DBSCAN: Vermeidung von k

Wie kann die Wahl von k vermieden werden?

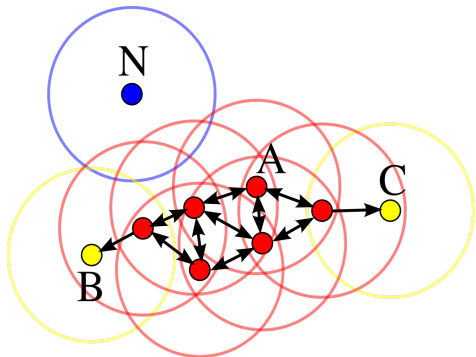
- Neue Verlustfunktion(?) = neuer Algorithmus

Kategorisierung der Punkte eines Datensatzes:

- Kernpunkte
- Dichte-erreichbar
- Rauschen

Problemstellung mittels

- Nachbarschaftsgröße ϵ
- Mindestanzahl an Nachbarn $minPts$



DBSCAN: Vermeidung von k

Wie kann die Wahl von k vermieden werden?

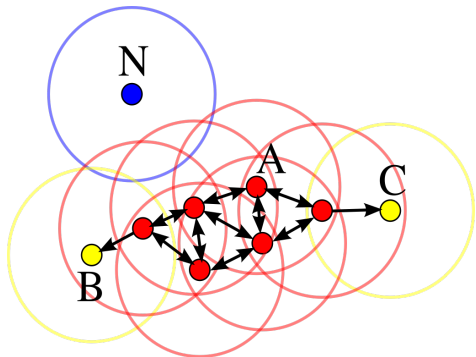
- Neue Verlustfunktion(?) = neuer Algorithmus

Kategorisierung der Punkte eines Datensatzes:

- Kernpunkte
- Dichte-erreichbar
- Rauschen

Problemstellung mittels

- Nachbarschaftsgröße ϵ
- Mindestanzahl an Nachbarn $minPts$



DBSCAN: Vermeidung von k

Wie kann die Wahl von k vermieden werden?

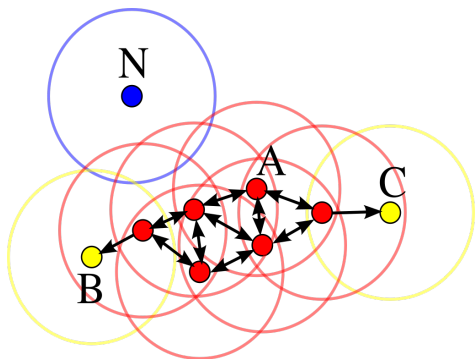
- Neue Verlustfunktion(?) = neuer Algorithmus

Kategorisierung der Punkte eines Datensatzes:

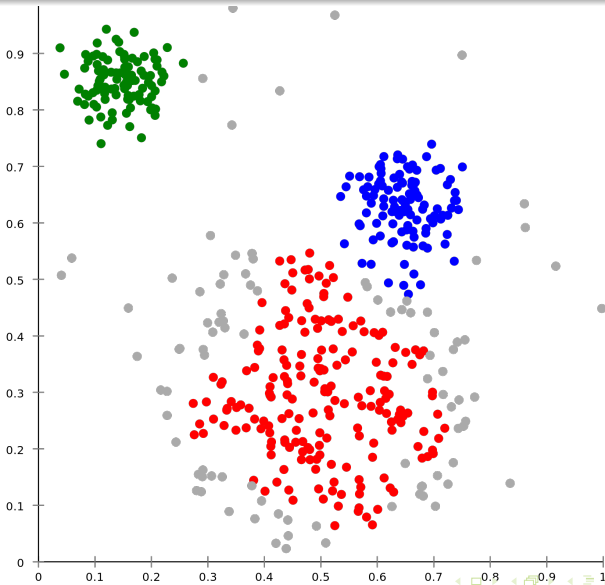
- Kernpunkte
- Dichte-erreichbar
- Rauschen

Problemstellung mittels

- Nachbarschaftsgröße ϵ
- Mindestanzahl an Nachbarn $minPts$



DBSCAN Clustering





Clustern von Text

Menge von Dokumenten = Korpus

Verschiedene Darstellungen denkbar:

- Vorhandensein von Worten (Binärvektoren, Mengen)
- Anzahl von Worten (Bag-of-Words)
- Term-Frequency-Inverse-Document-Frequency (TF-IDF)

$$tfidf_{w,d} = \frac{\text{Häufigkeit Wort } w \text{ in Dokument } d}{\max_{d'} \text{ Häufigkeit Wort } w \text{ in Dokument } d'} \log \frac{N}{N_w}$$

N_w ist Anzahl Dok. mit Wort w ; N ist Anzahl aller Dok.

- k -means Clustering der obigen Darstellungen möglich
- **Aber:** Bei großem Vokabular sind alle Dokumente weit weg (“Fluch der hohen Dimensionen”)



Clustern von Text

Menge von Dokumenten = Korpus

Verschiedene Darstellungen denkbar:

- Vorhandensein von Worten (Binärvektoren, Mengen)
- Anzahl von Worten (Bag-of-Words)
- Term-Frequency-Inverse-Document-Frequency (TF-IDF)

$$tfidf_{w,d} = \frac{\text{Häufigkeit Wort } w \text{ in Dokument } d}{\max_{d'} \text{Häufigkeit Wort } w \text{ in Dokument } d'} \log \frac{N}{N_w}$$

N_w ist Anzahl Dok. mit Wort w ; N ist Anzahl aller Dok.

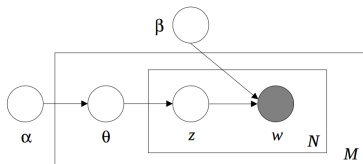
- k -means Clustering der obigen Darstellungen möglich
- **Aber:** Bei großem Vokabular sind alle Dokumente weit weg (“Fluch der hohen Dimensionen”)



Clustern von Text mit Graphischen Modellen (Topic Models)

Idee: Datengenerierender Prozess von Dokument d ist:

- Wähle Länge des Dokuments N_d aus $\mathbb{P}(N_d | \lambda)$ [Poisson]
- Wähle Themenverteilung θ_d aus $\mathbb{P}(\theta_d | \alpha)$ [Dirichlet]
- Erzeuge die Worte $w = 1 \dots N_d$:
 - Wähle ein Thema z_w aus $\mathbb{P}(z_w | \theta)$ [Kategorisch]
 - Wähle ein Wort aus $\mathbb{P}(w | z_w, \beta)$ [Kategorisch]



Lernen: Bestimme λ, α, β , und θ via Expectation-Maximization
(Maximum-Likelihood für unvollständige Daten)