

Masterarbeit

**Kernel k-means Methoden zur
spektralen Clusteranalyse von
Graphen**

Lukas Pradel

14. Februar 2015

Gutachter:

Prof. Dr. Christian Sohler

Dipl.-Inf. Melanie Schmidt

Technische Universität Dortmund

Fakultät für Informatik

Lehrstuhl II - Effiziente Algorithmen und Komplexitätstheorie

<http://ls2-www.cs.tu-dortmund.de/>

Hier kommt eine Danksagung hin!

Inhaltsverzeichnis

1	Einleitung	1
2	Grundlegende Definitionen und Algorithmen	7
2.1	Clustering und k -means	7
2.2	Graphen und Clusteranalyse von Graphen	9
2.3	Kernel-Methoden und spektrales Clustering	12
2.3.1	Kernel- k -means	12
2.3.2	Spektrale Clustering Methoden	16
3	Kernel k-means Methoden für spektrales Graphclustering	20
3.1	Graphschnitte und Graphclustering mit gewichtetem Kernel k -means	20
3.2	Die Kernel Methode für k -means++	26
4	Experimentelle Evaluation	35
4.1	Experimentelle Rahmenbedingungen	35
4.2	Die Kernel Methode	36
4.3	Der Kernel- k -means++-Algorithmus	39
4.4	Kernmengenkonstruktion	43
5	Zusammenfassung	44
	Literatur	46
	Erklärung	51

1 Einleitung

Die Clusteranalyse beschäftigt sich mit Verfahren, die Ähnlichkeitsstrukturen in (großen) Datenmengen identifizieren. Der englische Begriff *cluster* lässt sich in diesem Kontext am ehesten mit Gruppe oder Ansammlung übersetzen. Dem Gebiet liegt die Annahme zu Grunde, dass in gegebenen Datenmengen typischerweise Strukturen existieren, die mit einem geeigneten Ähnlichkeitsmaß erkannt werden können. Diese grundlegende Annahme ist insofern plausibel, als dass die zu analysierenden Datenmengen in der Praxis in der Regel aus einem konkreten Anwendungsfall gesammelt oder aufgezeichnet werden, in dem bestimmte Muster oder sich wiederholende Prozesse existieren. Bei der Analyse von Daten aus natürlichen Prozessen, sind diese beispielsweise oft gemäß einer bestimmten Wahrscheinlichkeitsverteilung verteilt. Hinzu kommt, dass Daten ohne jegliche Strukturen zumeist ohnehin nicht besonders aufschlussreich sind. Sind die vorhandenen Strukturen einmal erkannt, kann der gesamte Datensatz in Gruppen oder *Cluster* gegliedert werden, sodass innerhalb eines Clusters große Ähnlichkeit zwischen Objekten besteht, während Objekte, die unterschiedlichen Clustern zugeordnet sind, geringe Ähnlichkeit aufweisen. Anhand der Cluster, die mit Verfahren der Clusteranalyse gebildet wurden, lassen sich dann häufig aufschlussreiche Aussagen über die Eingabedaten treffen, die beim bloßen Anblick der gesamten unstrukturierten Datenmenge nicht möglich gewesen wären. Das gilt insbesondere, wenn die Daten hochdimensional sind, sodass eine graphische Auswertung auch nicht möglich ist.

Wir betrachten dazu einleitend ein einfaches Beispiel der Gegenwart: die Analyse von Strukturen in sozialen Netzwerken. Die Eingabeobjektmenge besteht in diesem Fall aus Objekten, die Informationen über Beziehungen und Aktivitäten zwischen Nutzern des sozialen Netzwerks enthalten. Mit Hilfe der Clusteranalyse können beispielsweise Gruppen von Benutzern identifiziert werden, die viel miteinander interagieren und somit eine soziale Community, also eine Gruppe von Freunden oder Bekannten, bilden. In diesem Beispiel kann es auch aufschlussreich sein, sogenannte „schwarze Löcher“ zu identifizieren, also Benutzer, die eine große Menge von eingehenden Aktivitäten zahlreicher anderer Benutzer haben, selbst aber nur selten mit einzelnen Benutzern interagieren. Bei Benutzern mit diesem Verhalten handelt es sich häufig um Prominente mit einer großen Fangemeinde. Unabhängig vom konkreten verfolgten Ziel wird schnell klar, dass entsprechende Aussagen gerade bei Netzwerken mit einer großen Nutzeranzahl, wie beispielsweise Facebook, nur dann möglich sind, wenn Verfahren der Clusteranalyse eingesetzt werden. Würde man lediglich die unstrukturierte gesamte Datenmenge betrachten, ließen sich

nur sehr schwer Informationen gewinnen, die keine Aussagen über das globale Netzwerk treffen. Es wird ebenso klar, dass es nötig ist, Verfahren zu entwickeln, die innerhalb der Grenzen gegenwärtiger Hardware Ergebnisse liefern, wenn die Eingabedatenmenge sehr große bis gigantische Größenordnungen annimmt.

In dieser Arbeit beschäftigen wir uns unter anderem mit einem weit verbreiteten Clusteringmodell, in dem das Ähnlichkeitsmaß die (quadrierte) geometrische Distanz ist. Die Eingabe besteht demnach aus einer Menge P von Punkten aus dem euklidischen Raum \mathbb{R}^d . Neben der Punktmenge P enthält unsere Eingabe zusätzlich eine positive ganze Zahl k , welche die Anzahl an Clustern angibt, in welche die Punkte unterteilt werden sollen. Dieses Modell ist unter dem Namen „ k -means-Problem“ bekannt. Den Algorithmen für das k -means-Problem, die wir betrachten wollen, ist gemein, dass sie eine initiale Auswahl von Clustern und zugehörigen *Zentren* iterativ weiter verbessern, bis eine zumindest lokal optimale Lösung ermittelt wurde. Die eigentliche Aufgabe besteht dann darin, die Summe der (quadrierten) Distanzen der Punkte eines Clusters zu ihrem Clusterzentrum zu minimieren.

Da das k -means-Problem schon für $k = 2$ NP-hart ist [ADHP09], sind für praktische Zwecke insbesondere Heuristiken, also Algorithmen ohne feste Schranken für Laufzeit oder Qualität der Lösung) und Approximationsalgorithmen, die eine nicht-optimale Lösung berechnen, welche jedoch garantiert nur um einen gewissen Faktor von der optimalen Lösung abweichen, interessant. Wir betrachten die Konstruktion von *Kernmengen*. Bei Eingabe einer Punktmenge P für das k -means-Problem handelt es sich bei einer Kernmenge um eine kleinere gewichtete Punktmenge S mit der Eigenschaft, dass für *jede* Wahl von k Clusterzentren die Summe der quadrierten Distanzen der Punkte in P zu den Zentren in etwa so groß ist wie die Summe der quadrierten Distanzen der Punkte in S zu den Zentren. Effektiv ist eine Kernmenge also ein hinsichtlich der Clusteringzielfunktion repräsentatives, (teilweise deutlich) kleineres Abbild der ursprünglichen Punktmenge. Üblicherweise bieten Kernmengenkonstruktionen Garantien einer $(1 + \epsilon)$ -Approximation, das heißt, dass ein k -means-Clustering auf der Kernmenge höchstens um einen Faktor $(1 + \epsilon)$ schlechter ist, als das selbe Clustering auf der Ursprungsmenge P . Neben der Approximation des k -means-Problems können Kernmengen auch ein effektives Mittel im Umgang mit großen Datenmengen sein. Wenn die Größe der Kernmenge S nur ein kleiner Anteil an der Größe von P ist, können Clustering-Algorithmen auf S deutlich schneller ausgeführt werden als auf P .

Schließlich betrachten wir in dieser Arbeit noch eine weitere Modellierung von

Clusteringproblemen, die mit der geometrischen Auffassung verwandt ist. Für einige Anwendungen ist es sinnvoll, die zu clusternden Objekte als Knoten eines gewichteten Graphen zu betrachten. Das Ähnlichkeitsmaß ist dann nicht mehr die euklidische Distanz, sondern das Gewicht der Kanten zwischen den Objekten beziehungsweise Knoten. Stellt man sich den gezeichneten Graphen vor, kommt das Kantengewicht der „Distanz“ zwischen den Knoten gleich. Bei der Clusteranalyse von Graphen können wir unter anderem auf diverse Ergebnisse aus der Graphentheorie zurückgreifen und zudem Techniken aus der linearen Algebra anwenden, wie wir später aufzeigen werden.

Anwendungen. Die Clusteranalyse war in der Vergangenheit bereits Gegenstand intensiver Forschung, die bis heute regelmäßig neue theoretische Erkenntnisse liefert. Es handelt sich aber um ein Gebiet mit hoher praktischer Relevanz, da Clustering in zahlreichen Anwendungen entweder selbst durchgeführt werden muss, oder als Teilproblem in einem übergeordneten Kontext auftritt. Wir wollen an dieser Stelle einen kurzen und sicherlich unvollständigen Überblick über einige wichtige Anwendungen der bisher vorgestellten Clusteringmodelle geben.

Eine gegenwärtige Anwendung für k -means-Clustering haben wir mit der Analyse sozialer Netzwerke bereits erwähnt. Eine weitere Anwendung, die gerade heute im Bereich der Online-Shops relevant ist, findet das k -means-Clustering in der Marktforschung. Hier sollen Informationen über das Einkaufsverhalten von Kunden genutzt werden, um diese beispielsweise den Marktsegmenten nach zu unterteilen. Je nachdem, wie detailliert die gesammelten Informationen sind, können diese auch genutzt werden, um Produktplatzierung zu optimieren, die Nachfrage nach neuen Produkten abzuschätzen, oder potenzielle neue Märkte zu erkennen. Neben den wirtschafts- und sozialwissenschaftlichen Anwendungsfällen, kommt die Clusteranalyse auch in geographischen und geologischen Zusammenhängen zum Einsatz. Konkret werden hier zum Beispiel chemische Eigenschaften von gesammelten Proben geclustert, um geographische Profile der untersuchten Substanzen zu erstellen.

Kernmengen werden in drei Szenarien eingesetzt. Wie bereits erwähnt bieten viele Kernmengenkonstruktionen eine garantierte Approximationsgüte, das heißt das anhand der Kernmenge berechnete Clustering weicht für jede beliebige Wahl der Zentren aus \mathbb{R}^d um beispielsweise einen Faktor von höchstens $1 + \epsilon$ von einem Clustering mit den selben Zentren auf der Ursprungsmenge ab. Man spricht in diesem Zusammenhang auch von *starken* Kernmengen. Ist diese Eigenschaft erfüllt, können wir mit einer Kernmenge, die eine deutlich kleinere Größe hat als die Ursprungsmenge, für Approximationsalgorithmen für das k -means-Problem einen

signifikanten *Speedup*, also eine Beschleunigung der Ausführungszeiten, erzielen. Dabei muss selbstverständlich beachtet werden, dass die Konstruktion der Kernmenge selbst ebenfalls Zeit in Anspruch nimmt und nicht zu aufwändig sein darf, wenn ein Speedup beabsichtigt wird.

Ein relativ neues Modell für Algorithmen sind sogenannte *Datenströme*, bei denen die Eingabe nicht initial vollständig vorliegt, sondern kontinuierlich in einem Datenstrom geliefert wird. Algorithmen für Datenströme, also *Datenstromalgorithmen* dürfen die Daten des Stroms nur genau einmal einlesen und sind zusätzlich polylogarithmisch bezüglich der Eingabegröße platzbeschränkt, das heißt, es ist nicht möglich, die gesamte Eingabe zu protokollieren. Beim Clustering in Datenströmen nutzt man aus, dass die Vereinigung einer $(1 + \epsilon_1)$ -Kernmenge und einer $(1 + \epsilon_2)$ -Kernmenge eine $(1 + \epsilon_1)(1 + \epsilon_2)$ -Kernmenge ergibt. Die Eingabe wird in kleinere Teile partitioniert, für die jeweils eine Kernmenge berechnet wird. Die Kernmengen werden dann zu einer einzelnen Kernmenge vereinigt. Da jede Vereinigung einen zusätzlichen Fehler verursacht, darf dies nicht zu oft erfolgen. Man setzt dazu die Merge-and-Reduce-Technik ein, die in einer Baum-Struktur bottom-up die Kernmengen einer Partition miteinander vereinigt. Mit der Technik wird insbesondere sichergestellt, dass jeder einzelne Punkt der Eingabe nur in $\log |P|$ vielen Reduktionsschritten verwendet wird, wobei $|P|$ die Anzahl an bisher betrachteten Eingabepunkten ist. Die Anwendung von Kernmengen im Bereich von Datenstromalgorithmen ist Gegenstand gegenwärtiger Forschung und hat bereits erste Ergebnisse erbracht, wie wir später noch sehen werden.

Schließlich können wir Kernmengen mit einer ähnlichen Idee einsetzen, um Clusteringalgorithmen parallel verteilt auszuführen. Wenn die Eingabepunktmenge verteilt vorliegt, berechnen wir auf den Teilen verteilt oder parallel Kernmengen, übermitteln die berechneten Kernmengen an eine zentrale Instanz und vereinen diese dort zu einer gesamten Kernmenge. Dadurch vermeiden wir es, die gesamte Punktmenge über das Netzwerk übermitteln zu müssen.

Die Clusteranalyse von Graphen findet ihre vielleicht wichtigste Anwendung in der Bildsegmentierung. Sie wird in diesem Kontext eingesetzt, um bei digitalen Bildern Kanten und insbesondere Objekte zu erkennen. Die Kantenerkennung ist eine fundamentale Technik in der digitalen Bildverarbeitung, die als Vorverarbeitung für speziellere Verfahren nötig ist. Das Erkennen von Objekten macht beispielsweise die Erkennung von Gesichtern, Personen oder Orten auf Fotos möglich.

Darüber hinaus ergibt sich eine natürliche Anwendung in der Identifikation relevanter Strukturen und insbesondere der Konnektivitätsanalyse in Netzwerken. Auch die Entdeckung von häufigen Anrufmustern in der Telekommunikation

kann hier von Interesse sein. Schließlich kann es in Netzwerken mit dynamischer Topologie, wie beispielsweise *ad-hoc-Netzwerken*, die in intelligenten Umgebungen mit vielen wechselnden mobilen Geräten vorkommen, hilfreich sein, zu clustern. Lokales Clustering wird hier unter anderem eingesetzt, um Routing-Algorithmen zu verbessern. Das Netzwerk wird dabei als Graph aufgefasst.

Zudem kann die Clusteranalyse von Graphen eingesetzt werden, um Caching-Mechanismen in Datenbanken zu optimieren. Konkret werden in Datenbanksystemen die einzelnen Einträge in Form von *Seiten* auf dem Speichermedium abgelegt, von denen einige im Hauptspeicher gecached werden. Ein gutes Clustering führt in diesem Kontext dazu, dass das Laden eines Clusters, dem ein angefragtes Datum zugeordnet ist, zudem relevante Daten liefert, die für ähnliche zukünftige Anfragen erforderlich sind.

Verwandte Arbeiten. Der Algorithmus von Lloyd [Llo82] ist eine schon früh entwickelte Heuristik für das k -means-Problem, die bis heute von praktischer Relevanz ist, da sie im Allgemeinen gute Clusterings produziert, obwohl der Algorithmus keine Approximationsgüte garantiert. Wir betrachten in dieser Arbeit insbesondere die Verbesserung von Lloyds Algorithmus durch Arthur und Vassilvitskii [AV07], die eine $\mathcal{O}(\log k)$ -Approximation für das k -means-Problem ist.

Das heutige Konzept von *starken* Kernmengen wurde in [HM04] vorgestellt. Wir beschäftigen uns in dieser Arbeit hauptsächlich mit der Konstruktion aus [FSS13, Sch14].

Das Clusteringmodell in Graphen kann bis zu [KL70] zurückverfolgt werden. Die erste Anwendung von sogenannten *spektralen* Methoden für die Clusteranalyse von Graphen und insbesondere die Anwendung im Bereich der Bildsegmentierung stammt von [SM00]. Eine wichtige Weiterentwicklung des Ansatzes geht zurück auf [NJW01]. Die unmittelbare Äquivalenz zum geometrischen k -means-Problem wurde in [DGK04, DGK07] gezeigt. Wir werden sehen, dass diese für unsere Zwecke besonders wichtig ist.

Die Arbeit ist folgendermaßen gegliedert. In Kapitel 2 definieren wir die wichtigsten Begriffe und Konzepte, die wir im Rahmen dieser Arbeit benötigen. Außerdem stellen wir die Algorithmen vor, die wir als Grundlage verwenden, oder weiterentwickeln beziehungsweise verbessern wollen. In Kapitel 3 entwickeln wir zum einen eine Verbesserung des Algorithmus von Dhillon, Guan und Kulis [DGK04, DGK07] und zeigen zum anderen, mit welchen Modifikationen die Kernmengenkonstruktion von Feldman, Sohler und Schmidt [FSS13, Sch14] praktisch umsetzbar wird. Die entwickelten Algorithmen wollen wir dann in Kapitel 4 anhand einer Reihe

von Experimenten empirisch bewerten. Unsere Ergebnisse fassen wir schließlich in Kapitel 5 zusammen.

2 Grundlegende Definitionen und Algorithmen

In diesem Kapitel definieren wir die für unsere Zwecke relevanten Begriffe im Kontext der Clusteranalyse und führen die wichtigen grundlegenden Algorithmen ein, deren Ideen für uns im Folgenden noch von Bedeutung sein werden. Wir gehen dabei nach Themengebieten geordnet vor: In Abschnitt 2.1 skizzieren wir kurz das Themengebiet der Clusteranalyse, definieren die üblichen Zielfunktionen und stellen zwei bedeutende Algorithmen vor. Abschnitt 2.2 führt kurz in die Graphentheorie sowie die Clusteranalyse von Graphen ein. In diesem Abschnitt werden wir zudem verschiedene Optimierungskriterien für die Clusteranalyse von Graphen herausstellen. Schließlich fassen wir in Abschnitt 2.3 die wichtigsten Methoden und Algorithmen aus dem Bereich der spektralen Clusteranalyse zusammen und stellen zudem die wesentlichen Konzepte von Kernel-Methoden vor.

2.1 Clustering und k -means

Clusteranalyse oder „Clustering“ beschäftigt sich mit der Einteilung von Objekten in Gruppen („Cluster“), sodass sich die Objekte innerhalb eines Clusters gemäß eines bestimmten Optimierungskriteriums ähnlich sind und von Objekten eines anderen Clusters unterscheiden. Es existieren zahlreiche grundsätzlich verschiedene Ansätze, Clusteringprobleme zu lösen. Wir beschränken uns in dieser Arbeit auf *partitionierende* Clusteringprobleme und -verfahren. Bei diesen soll eine Menge von Objekten, welche der erste Teil der Eingabe ist, gemäß einer Cluster-Zielfunktion möglichst optimal in genau k Cluster unterteilt werden. Dabei ist k der ganzzahlige zweite Teil der Eingabe, das heißt, die Anzahl an Clustern muss bei diesen Verfahren vorab festgelegt werden. Die initial gewählten k Zentren werden dann iterativ verschoben, bis ein gewisses Abbruchkriterium erfüllt ist.

Für die Zielfunktion, welche die Nähe oder Ferne von Punkten zueinander quantifiziert, sind bei Eingabepunkten aus \mathbb{R}^d Metriken naheliegend. Intuitiv ist dabei die euklidische Distanz, auf der die beiden bekanntesten Clustering-Problemstellungen basieren. Wir notieren die euklidische Distanz zwischen zwei Punkten $x, y \in \mathbb{R}^d$ mit $\|x - y\|$.

Definition 2.1.1 (k -median und k -means). Sei $P \subset \mathbb{R}^d$ und $k \in \mathbb{N}^+$. Das k -median-Problem besteht darin, eine Menge von k (Cluster-)Zentren $C = \{c_1, \dots, c_k\}$ mit $c_i \in \mathbb{R}^d$ zu finden, sodass der folgende Term minimal wird:

$$\sum_{p \in P} \min_{c \in C} \|p - c\|$$

Das k -means-Problem unterscheidet sich nur darin, dass bei diesem die Summe der *quadrierten* euklidischen Distanzen zum jeweils nächstgelegenen Zentrum minimiert werden soll, das heißt, dass der folgende Term minimiert werden soll:

$$\sum_{p \in P} \min_{c \in C} \|p - c\|^2$$

Beim *gewichteten* k -means-Problem werden den Eingabepunkten zusätzlich mit einer Funktion $w : P \rightarrow \mathbb{R}$ Gewichte zugewiesen. Die zu minimierende Zielfunktion lautet dann entsprechend

$$\sum_{p \in P} \min_{c \in C} w(p) \|p - c\|^2$$

Sowohl das k -Median-Problem [MS84] als auch das k -means-Problem [ADHP09] sind optimal NP-schwer lösbar. Typischerweise werden zur Lösung daher approximative oder heuristische Algorithmen eingesetzt. Die bekannteste und bis heute sehr erfolgreiche Heuristik für das k -means-Problem ist der Algorithmus von Lloyd [Llo82]. Wegen seiner großen Popularität wird der Algorithmus häufig auch als „ k -means-Algorithmus“ bezeichnet. Der Algorithmus wählt initial k zufällige Punkte aus der Eingabemenge als initiale Clusterzentren. Darüber hinaus existieren zahlreiche weitere Methoden zur initialen Wahl der Zentren, siehe dazu beispielsweise [CKV13]. Anschließend wird jedem Punkt das am nächsten gelegene Zentrum zugewiesen. Dadurch entstehen die initialen Cluster mit ihren jeweiligen Zentren. Im zweiten Schritt wird das neue Zentrum eines jeden Clusters als der geometrische Zentroid des Clusters gewählt. Die Wahl des geometrischen Zentroiden als neues Zentrum ist für das 1-means-Problem optimal [ORSS12]. Die Zuweisung von Punkten zum nächstgelegenen Clusterzentrum und die Berechnung der neuen Zentroiden werden solange alterniert, bis die Lösung konvergiert, also wenn sich die Zuordnungen der Punkte nicht mehr ändern. In der Praxis wird gelegentlich auch nach einer festen Anzahl von Iterationen terminiert.

Algorithmus 1: Algorithmus von Lloyd

Eingabe: $P \subseteq \mathbb{R}^d, k \in \mathbb{N}^+$

Ausgabe: k -means-Clustering von P

- 1 Wähle zufällig k Zentren $c_1^{(0)}, \dots, c_k^{(0)}$ aus P oder \mathbb{R}^d
 - 2 $S_i^{(0)} \leftarrow \{p \in P : \|p - c_i^{(0)}\|^2 \leq \|p - c_{i'}^{(0)}\|^2 \forall i' \in \{1, \dots, k\}\}$
 - 3 **repeat**
 - 4 $c_i^{(t)} \leftarrow \frac{1}{|S_i^{(t-1)}|} \sum_{p_j \in S_i^{(t-1)}} p_j$
 - 5 $S_i^{(t)} \leftarrow \{p \in P : \|p - c_i^{(t)}\|^2 \leq \|p - c_{i'}^{(t)}\|^2 \forall i' \in \{1, \dots, k\}\}$
 - 6 **until** $S_i^{(t)} = S_i^{(t-1)}$
-

Die asymptotische Laufzeit des Algorithmus beträgt $\mathcal{O}(nkdi)$, wobei i die Anzahl an durchgeführten Iterationen ist. Wenn der Algorithmus konvergiert und nicht durch eine feste Anzahl von Iterationen terminiert wird, wurde ein lokales Optimum gefunden, welches jedoch im Allgemeinen kein globales Optimum oder eine Approximation eines globalen Optimums ist.

Der Algorithmus *k-means++* [AV07] sieht eine alternative Wahl der initialen Clusterzentren in Lloyds Algorithmus vor: er wählt diese auf einfache, aber dennoch geschickte Art und Weise und führt anschließend die übrigen Schritte von Lloyds Algorithmus aus. Dazu wird zunächst ein einzelnes Clusterzentrum c_1 zufällig gleichverteilt aus der Eingabe-Punktmenge P gewählt. Alle weiteren Clusterzentren werden sukzessive nach der folgenden Vorschrift gewählt, bis insgesamt k Zentren gewählt wurden. Im Weiteren bezeichnen wir mit $D(x)$ für einen Punkt x aus der Eingabe-Punktmenge P die geringste Distanz von x zum nächstgelegenen bereits gewählten Zentrum. In jeder Iteration wird als nächstes Zentrum c_i der Punkt $x' \in P \setminus \{c_1, \dots, c_{i-1}\}$ mit Wahrscheinlichkeit $\frac{D(x')^2}{\sum_{x \in P} D(x)^2}$ gewählt.

Algorithmus 2: *k-means++*

Eingabe: : $P \subseteq \mathbb{R}^d, k \in \mathbb{N}^+$

Ausgabe: : k initiale Clusterzentren für P

- 1 Wähle c_1 zufällig gleichverteilt aus P
 - 2 **for** $i \leftarrow 2$ **to** k **do**
 - 3 Wähle den Punkt $x' \in P \setminus \{c_1, \dots, c_{i-1}\}$ als Zentrum c_i mit
 Wahrscheinlichkeit $\frac{D(x')^2}{\sum_{x \in P} D(x)^2}$
 - 4 Führe Lloyds Algorithmus mit den initialen Clusterzentren c_1, \dots, c_k aus.
-

Die k Zentren, die von *k-means++* ausgewählt werden, sind eine $\mathcal{O}(\log k)$ -Approximation für das *k-means*-Problem, die durch die anschließende Ausführung von Lloyds Algorithmus noch zu einem lokalen Optimum verbessert werden.

Im nächsten Abschnitt betrachten wir das Clusteringproblem auf Graphen, das wir im später noch genauer untersuchen wollen.

2.2 Graphen und Clusteranalyse von Graphen

Der *Graph* ist in der Informatik eine vielseitig einsetzbare und gut erforschte Datenstruktur. Wir können diverse Clusteringprobleme nicht nur für eine Eingabepunktmenge formulieren, sondern auch für eine Eingabe, die aus einem Graphen besteht. Wir beginnen mit der folgenden Definition des Graphen.

Definition 2.2.1 (Graph). Sei V eine endliche Menge und $E \subseteq \{\{u, v\} \mid u, v \in V, u \neq v\}$. Dann heißt das Tupel $G = (V, E)$ ein (endlicher) *Graph* mit *Knotenmenge* V und *Kantenmenge* E . Ist $e = \{u, v\} \in E$, dann sagen wir, dass die Kante e des Graphen G die Knoten u und v *verbindet*. In diesem Fall sind u und v die *Endknoten* von e . Ein Knoten $u \in V$ und eine Kante $e \in E$ heißen *inzident* genau dann, wenn $u \in e$. Wir sagen, dass u und ein weiterer Knoten $v \in V$ *adjazent* sind genau dann, wenn es eine Kante $e' = \{u, v\}$ in E gibt. Typischerweise bezeichnet $n = |V|$ die *Knotenzahl* von G und $m = |E|$ die *Kantenzahl* von G .

Bei einem 3-Tupel $G' = (V', E', w')$ mit $w' : E \rightarrow \mathbb{N}$ spricht man von einem *gewichteten* Graphen, dessen Kanten über ein Gewicht verfügen, das von der Gewichtsfunktion w' abgebildet wird.

Für die konkrete Datenhaltung von Graphen haben sich im Wesentlichen zwei Ansätze als praktikabel erwiesen: Bei den sogenannten *Adjazenzlisten* wird für jeden Knoten v im Graphen eine Liste Adj_v gehalten, in der für jeden zu v inzidenten Knoten ein Eintrag in Adj_v enthalten ist, der den entsprechenden adjazenten Knoten referenziert, sowie gegebenenfalls das Gewicht der Kante zwischen den beiden Knoten. Alternativ wird in einer *Adjazenzmatrix* Adj^G der Größe $|V| \times |V|$ an der Stelle $Adj_{u,v}^G$ das Gewicht der Kante zwischen den Knoten u und v eingetragen, sofern die beiden Knoten durch eine Kante miteinander verbunden sind. Anderenfalls wird zumeist -1 oder 0 eingetragen. Bei ungewichteten Graphen wird dementsprechend lediglich 0 oder 1 eingetragen.

Wenn die Eingabe keine Punktmenge, sondern ein (gewichteter) Graph ist, können wir die Ideen der bereits vorgestellten Verfahren weiterhin anwenden. Wir interessieren uns in diesem Fall für die „Ähnlichkeit“ von Knotenmengen im Graphen.

Definition 2.2.2 (Graph-Schnitt). Sei $G = (V, E, w)$ ein gewichteter Graph. Ein *Schnitt* $C = (S, T)$ von G ist eine Partitionierung von V in die beiden Mengen S und T , das heißt, dass $V = S \cup T$ und gleichzeitig $S \cap T = \emptyset$. Die *Schnittmenge* von C sind die Kanten in E , die einen Endpunkt in S und den anderen Endpunkt in T haben, das heißt formal ist die Schnittmenge definiert als

$$\{(u, v) \in E \mid u \in S, v \in T\}.$$

Das Gewicht oder der Wert eines Schnittes ist die Summe der Kantengewichte der Schnittmenge. Wir verwenden die folgende Notation:

$$w_{cut}(S, T) = \sum_{u \in S, v \in T} w((u, v))$$

Falls der Graph in Form einer Adjazenzmatrix Adj^G vorliegt, lautet die Berechnungsvorschrift entsprechend:

$$w_{cut}(S, T) = \sum_{u \in S, v \in T} Adj_{u,v}^G$$

Für eine direkte Analogie zum k -means-Problem wäre ein Schnitt- beziehungsweise Partitionierungs-Begriff wünschenswert, der eine k -fache Partitionierung der Knotenmenge erlaubt. Diese lautet wie folgt.

Definition 2.2.3 (k -Graphpartitionierung). Sei $G = (V, E, w)$ ein gewichteter Graph. Für zwei Mengen $A, B \subseteq V$ definieren wir:

$$w(A, B) = \sum_{u \in A, v \in B} w((u, v))$$

Das k -Graphpartitionierungsproblem besteht darin, eine Partitionierung der Knotenmenge V in k disjunkte Teilmengen V_1, \dots, V_k mit $\bigcup_{i \in \{1, \dots, k\}} V_i = V$ zu ermitteln, sodass sich die Knoten innerhalb einer Partition bezüglich einer Ähnlichkeitsrelation möglichst ähnlich sind und die Knoten unterschiedlicher Partitionen bezüglich der Ähnlichkeitsrelation möglichst stark voneinander unterscheiden [KL70].

Beim Clustering von Punktmengen lagen für die Unähnlichkeitsrelation Metriken nahe, im Falle der Graphpartitionierung existiert eine Reihe von Optimierungskriterien, von denen wir im Folgenden die verbreitetsten einführen.

1. **Ratio Association.** Bei der Ratio Association sollen die Intra-Clusterabstände relativ zur jeweiligen Clustergröße maximiert werden:

$$\max \sum_{c=1}^k \frac{w(V_c, V_c)}{|V_c|}$$

2. **Ratio Cut.** Beim Ratio Cut wird das Gewicht des Schnitts zwischen jeweils einem Cluster und allen anderen Punkten im Graphen minimiert:

$$\min \sum_{c=1}^k \frac{w(V_c, V \setminus V_c)}{|V_c|}$$

3. **Kernighan-Lin.** Bei dem in [KL70] vorgestellten Optimierungskriterium werden die Intra-Clusterabstände ähnlich der Ratio Association minimiert, allerdings müssen alle Partitionen hier zusätzlich die selbe Größe haben. Die

hier vorgestellte Zielfunktion ist von 2 auf k Partitionen verallgemeinert:

$$\min \sum_{c=1}^k \frac{w(V_c, V \setminus V_c)}{|V_c|} \text{ s.t. } |V_c| = \frac{|V|}{k} \forall c \in \{1, \dots, k\}$$

4. **Normalized Cut.** Ziel des Normalized Cut ist wie beim Ratio Cut die Minimierung des Schnitts von einem Cluster mit den übrigen Punkten im Graphen, jedoch relativ zum Schnitt des Clusters mit *allen* Knoten im Graphen. Der Normalized Cut oder kurz NCut ist in der Bild-Segmentierung recht verbreitet [SM00]. Wir betrachten hier ebenfalls die auf k Partitionen verallgemeinerte Zielfunktion.

$$\min \sum_{c=1}^k \frac{w(V_c, V \setminus V_c)}{w(V_c, V)}$$

Im nächsten Abschnitt betrachten wir zwei Techniken, mit denen Limitierungen hinsichtlich der Qualität der erzielten Clusterings sowie der Komplexität der Berechnungen verbessert werden können.

2.3 Kernel-Methoden und spektrales Clustering

Legt man die k -means Zielfunktion zu Grunde, können im Ursprungsraum \mathbb{R}^d nur Cluster berechnet werden, die *linear separierbar* sind. Ein illustratives Beispiel im \mathbb{R}^2 sind zwei konzentrisch angeordnete ringförmige Punktmengen. Die intuitiv optimalen Cluster sind die beiden jeweiligen Ringe, es ist jedoch im \mathbb{R}^2 nicht möglich Clusterzentren anzugeben, die den zweidimensionalen Raum in Ringe partitionieren. Die Beschränkung liegt darin, dass wir die Partitionierung oder Separierung durch Hyperebenen vornehmen. Im Ursprungsraum können wir daher nur linear separierbare Cluster bestimmen. Wir betrachten nun ein Verfahren, mit dem sich diese Beschränkung umgehen lässt.

2.3.1 Kernel- k -means

Unsere fundamentale Motivation für diesen Abschnitt ist Covers Theorem [Cov65], welches informell zusammengefasst besagt, dass eine zufällige Eingabepunktmenge im d -dimensionalen Raum mit hoher Wahrscheinlichkeit in einen höherdimensionalen Raum (zum Beispiel der Dimension $D \gg d$) abgebildet linear separierbar ist.

Für Klassifizierungs- und insbesondere Clusteringprobleme benötigen wir sowohl im Ursprungsraum als auch im abgebildeten Raum ein algebraisches Maß für

Unähnlichkeit. Üblicherweise wird dabei eine auf *Skalarprodukten* basierende Metrik eingesetzt. Wir beziehen uns im Weiteren auf die folgende Definition von *euklidischen Vektorräumen*, also Vektorräumen mit einem reellen Skalarprodukt.

Definition 2.3.1 (Euklidischer Vektorraum). Sei V ein Vektorraum über \mathbb{R} . Ein reelles Skalarprodukt auf V ist eine Abbildung $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{R}$, die den folgenden Eigenschaften genügt:

1. $\langle a, b \rangle = \langle b, a \rangle \forall a, b \in V$
2. $\langle \lambda a + \mu b, c \rangle = \lambda \langle a, c \rangle + \mu \langle b, c \rangle \forall a, b, c \in V$ und $\lambda, \mu \in \mathbb{R}$
3. $\langle a, a \rangle \geq 0$ und $\langle a, a \rangle = 0 \Leftrightarrow a = 0 \forall a \in V$

Ein \mathbb{R} -Vektorraum mit einem reellen Skalarprodukt heißt *euklidischer Vektorraum*.

Die k -means-Zielfunktion kann auf euklidische Vektorräume erweitert werden, indem wir für einen euklidischen Vektorraum \mathcal{V} eine Norm $\|\cdot\| : \mathcal{V} \rightarrow \mathbb{R}$ über das Skalarprodukt gemäß der folgenden Abbildungsvorschrift für jedes $\chi \in \mathcal{V}$ definieren: $\|\chi\| = \sqrt{\langle \chi, \chi \rangle}$. Das auf Skalarprodukten basierte Analog zu Definition 2.1.1 des k -means-Problems lautet dann wie folgt:

Definition 2.3.2 (k -means in euklidischen Vektorräumen [Sch14]). Sei \mathcal{V} ein euklidischer Vektorraum mit der Norm $\|\cdot\| : \mathcal{V} \rightarrow \mathbb{R}$ und sei $\mathcal{P} \subset \mathcal{V}$ sowie $k \in \mathbb{N}^+$. Das k -means-Problem in \mathcal{V} besteht darin, eine Menge von k (Cluster-)Zentren $\mathcal{C} = \{c_1, \dots, c_k\}$ mit $c_i \in \mathcal{V}$ zu finden, sodass der folgende Term minimal wird:

$$\sum_{p \in \mathcal{P}} \min_{c \in \mathcal{C}} \|p - c\|^2$$

Die gewichtete Variante des k -means-Problems in euklidischen Vektorräumen lautet entsprechend wie folgt:

$$\sum_{p \in \mathcal{P}} \min_{c \in \mathcal{C}} w(p) \|p - c\|^2$$

Damit ist es nun möglich, ein intuitives Verfahren zum Bestimmen von Clusterings mit nicht-linearen Separatoren anzugeben. Wir bilden zunächst mit einer Abbildung $\phi : \mathcal{X} \rightarrow \mathcal{V}$ unsere Eingabepunkte aus ihrem Ursprungsraum \mathcal{X} (also beispielsweise \mathbb{R}^d) auf einen D -dimensionalen (üblicherweise $D \gg d$) euklidischen Vektorraum \mathcal{V} ab, und berechnen in diesem mit Lloyds Algorithmus oder einem anderen Clusteringverfahren ein Clustering, welches im Ursprungsraum \mathcal{X} im Allgemeinen

nicht-linear separierten Clustern entspricht. Bei diesem Verfahren wird entsprechend der folgende Term minimiert:

$$\sum_{p \in P} \min_{c \in C} \|\phi(p) - c\|^2$$

Die Zentroid-Berechnung zur Bestimmung eines Clusterzentrums c_i im Sinne von Lloyds Algorithmus lautet dann

$$c_i = \frac{\sum_{p \in S_i} \phi(p)}{|S_i|}$$

wobei S_i die Menge aller dem Clusterzentrum c_i zugeordneten Punkte ist. Zunächst haben wir uns damit nicht-linear separierte Cluster „erkauft“, der Preis, den wir dafür zahlen, besteht jedoch in höherem Rechenaufwand, da wir einerseits die Punkte in den jeweils höherdimensionalen Raum abbilden müssen und andererseits die Berechnung der (euklidischen) Distanzen aufwändiger wird; diese benötigt im Ursprungsraum Zeit $\mathcal{O}(d)$ und im höherdimensionalen Raum $\mathcal{O}(D)$. Der Mehraufwand für die Distanzberechnung fällt dabei stärker ins Gewicht, da wir die Distanzberechnung beispielsweise bei Lloyds Algorithmus in jeder Schleifeniteration durchführen.

Dieses Problem wurde erstmals im Zusammenhang mit *Support Vector Machines* (kurz SVMs) gelöst [BGV92]. Wir übertragen die Lösung mittels des sogenannten *Kernel-Tricks* in den Clustering-Kontext und skizzieren diese kurz. Zunächst halten wir eine wichtige Beobachtung hinsichtlich der (quadrierten) euklidischen Distanz $\|\phi(p) - c\|^2$ zwischen dem abgebildeten Punkt und einem Clusterzentrum fest. Diese lässt sich wie folgt ausschließlich über Skalarprodukte berechnen (wobei auch hier wieder S_c die dem Clusterzentrum c zugeordnete Punktmenge ist):

$$\begin{aligned} \|\phi(p) - c\|^2 &= \left\| \phi(p) - \frac{1}{|S_c|} \sum_{x \in S_c} \phi(x) \right\|^2 \\ &= \left\langle \phi(p) - \frac{1}{|S_c|} \sum_{x \in S_c} \phi(x), \phi(p) - \frac{1}{|S_c|} \sum_{x \in S_c} \phi(x) \right\rangle \\ &= \langle \phi(p), \phi(p) \rangle - \frac{1}{|S_c|} \sum_{x \in S_c} \langle \phi(p), \phi(x) \rangle - \frac{1}{|S_c|} \sum_{x \in S_c} \langle \phi(x), \phi(p) \rangle \\ &\quad + \frac{1}{|S_c|^2} \sum_{x, y \in S_c} \langle \phi(x), \phi(y) \rangle \end{aligned} \tag{2.1}$$

Diese Beobachtung allein hilft uns noch nicht weiter, da auch die Berechnung des

Linear-Kernel	$\kappa(x_i, x_j) = x_i \cdot x_j + \gamma$
Polynom-Kernel	$\kappa(x_i, x_j) = (x_i \cdot x_j + \gamma)^\delta$
Gauss-/RBF-Kernel	$\kappa(x_i, x_j) = \exp\left(-\frac{\ x_i - x_j\ ^2}{2\sigma^2}\right)$
Sigmoid-Kernel	$\kappa(x_i, x_j) = \tanh(\alpha(x_i \cdot x_j) + \theta)$

Tabelle 1: Beispiele für häufig eingesetzte Kernelfunktionen im Maschinellen Lernen [HSS08].

Skalarproduktes im höherdimensionalen Raum einen asymptotischen Aufwand von $\mathcal{O}(D)$ hätte. An dieser Stelle kommt nun der in [BGV92] vorgestellte *Kernel-Trick* zum Tragen. Durch den Einsatz von sogenannten *Kernelfunktionen* lässt sich dieser Mehraufwand umgehen. Dabei handelt es sich um positiv definite Abbildungen $\kappa : X \times X \rightarrow \mathbb{R}$, die bei Eingabe von zwei Punkten aus dem Ursprungsraum X das Skalarprodukt der in den höherdimensionalen Raum abgebildeten Punkte berechnen:

$$\kappa(x, y) = \langle \phi(x), \phi(y) \rangle$$

Der Kernel-Trick besteht also darin, dass sich die Distanzberechnung auf Skalarprodukte von Eingabepunkten reduzieren lässt, für die keine explizite Kenntnis der Abbildung ϕ oder der abgebildeten Punkte $\phi(p)$ nötig ist. Dies können wir umsetzen, indem wir insgesamt $\mathcal{O}(n^2)$ Skalarprodukte berechnen oder direkt eine Kernelfunktion κ verwenden. Mit der Kernelfunktion κ können wir die Distanzberechnung aus 2.1 weiter vereinfachen:

$$\|\phi(p) - c\|^2 = \kappa(p, p) - \frac{1}{|S_c|} \sum_{x \in S_c} \kappa(p, x) - \frac{1}{|S_c|} \sum_{x \in S_c} \kappa(x, p) + \frac{1}{|S_c|^2} \sum_{x, y \in S_c} \kappa(x, y) \quad (2.2)$$

Tabelle 1 bietet eine Übersicht über häufig eingesetzte Kernelfunktionen im Maschinellen Lernen [HSS08]. Insbesondere die Gauss-Kernelfunktion und die Polynom-Kernelfunktion werden auch bei Support-Vektor-Maschinen gerne verwendet.

Eine auf Kernelfunktionen basierende Variante von (Lloyds) Algorithmus 1 folgt mit 2.2 unmittelbar. Der Algorithmus wird üblicherweise kurz „Kernel- k -means-Algorithmus“ genannt. Der Vorfaktor für gewichtetes k -means hat auf die Berechnungsvorschriften keine weiteren Auswirkungen, Algorithmus 3 kann dementsprechend für gewichtetes k -means erweitert werden. Der Algorithmus lässt sich mit einer asymptotischen Laufzeit von $\mathcal{O}(n^2(\tau + d))$ implementieren [DGK04, DGK07],

wobei τ die Anzahl an Iterationen der äußeren Schleife ist.

Im nächsten Unterabschnitt führen wir Verfahren ein, die mit Techniken aus der

Algorithmus 3: Kernel- k -means

Eingabe: $: P \subseteq \mathbb{R}^d, k \in \mathbb{N}^+, \kappa : X \times X \rightarrow \mathbb{R}$

Ausgabe: $: k$ -means-Clustering von P

```

1 Wähle zufällig  $k$  Zentren  $c_1^{(0)}, \dots, c_k^{(0)}$  aus  $P$  oder  $\mathbb{R}^d$ 
2  $t \leftarrow 0$ 
3 repeat
4   Berechne  $\|\phi(p) - c_i^{(t)}\|^2$  für alle  $p \in P$  und  $i \in \{1, \dots, k\}$  mit 2.2
5   foreach  $p \in P$  do
6      $c^*(p) \leftarrow \arg \min_i \left( \left\| \phi(p) - c_i^{(t)} \right\| \right)$ 
7    $t \leftarrow t + 1$ 
8   for  $i \leftarrow 1$  to  $k$  do
9      $S_i^{(t)} \leftarrow \{p \mid c^*(p) = i\}$ 
10     $c_i^{(t)} \leftarrow \frac{1}{|S_i^{(t)}|} \sum_{p_j \in S_i^{(t)}} p_j$ 
11 until Konvergenz oder  $t > t_{max}$ 

```

linearen Algebra das Graphclusteringproblem lösen und in Kombination mit dem hier vorgestellten Algorithmus eine robuste Grundlage für unsere Zwecke bilden.

2.3.2 Spektrale Clustering Methoden

Beim spektralen Clustering besteht die Eingabe aus einem Graphen, dessen Knoten geclustert werden sollen. Informationen über die „Ähnlichkeit“ von Objekten werden in Form eines Kantengewichts der Kante zwischen den jeweiligen Objekten realisiert. Wir werden im Weiteren insbesondere Eigenvektoren und zugehörige Eigenwerte von Matrizen betrachten und rufen uns daher kurz die Definition dieser in Erinnerung.

Definition 2.3.3 (Eigenvektor und Eigenwert). Sei A eine quadratische Matrix. Das *Eigenwertproblem* besteht darin, eine Zahl λ und einen dazugehörigen Vektor $\vec{x} \neq \vec{0}$ zu finden, für die gilt:

$$A\vec{x} = \lambda\vec{x}$$

Die reelle oder komplexe Zahl λ heißt *Eigenwert* von A , der Vektor \vec{x} heißt Eigenvektor von A .

Wenn G der Graph ist, durch den die Eingabemenge repräsentiert wird, dann ist die Adjazenzmatrix Adj^G des Graphen der Ausgangspunkt der spektralen

Clusteranalyse. Bevor wir mit konkreten Verfahren der spektralen Clusteranalyse beginnen können, benötigen wir einige grundlegende Definitionen der (spektralen) Graphentheorie.

Definition 2.3.4 (Grad-Matrix). Sei $G = (V, E)$ ein Graph mit $|V| = n$ Knoten. Die *Grad-Matrix* von G ist die $n \times n$ -Diagonalmatrix D , die folgendermaßen definiert ist:

$$D_{i,j} = \begin{cases} \deg(v_i) & \text{falls } i = j, \\ 0 & \text{sonst} \end{cases}$$

Dabei ist $\deg(v_i)$ die Anzahl an Kanten, die zum Knoten v_i inzident sind. Diese Anzahl wird auch *Grad* des Knotens v_i genannt.

Mit Hilfe der Grad-Matrix eines Graphen lässt sich die sogenannte *Laplace-Matrix* des Graphen berechnen, welche eine Matrix-Repräsentation von Graphen ist, die aufschlussreiche Informationen über die Struktur und Beschaffenheit des Graphen erlaubt. Sie ist wie folgt definiert:

Definition 2.3.5 (Laplace-Matrix [Lux07]). Sei $G = (V, E)$ ein Graph ohne Schleifen mit Grad-Matrix D und (gewichteter) Adjazenzmatrix Adj^G . Die (nicht-normalisierte) *Laplace-Matrix* L von G ist definiert als:

$$L = D - Adj^G$$

Die Einträge der Matrix sind dementsprechend folgendermaßen definiert:

$$L_{i,j} = \begin{cases} \deg(v_i) & \text{falls } i = j, \\ -1 & \text{falls } i \neq j \text{ und } v_i \text{ inzident zu } v_j \text{ ist,} \\ 0 & \text{sonst} \end{cases}$$

Es existieren zwei normalisierte Formen der Laplace-Matrix, die *normalisierte* Laplace-Matrizen genannt werden. Sie sind wie folgt definiert:

$$\begin{aligned} \mathcal{L}_{sym} &= D^{-\frac{1}{2}} L D^{-\frac{1}{2}} \\ \mathcal{L}_{rw} &= D^{-1} L \end{aligned}$$

Die Eigenwerte der Laplace-Matrizen eines Graphen geben beispielsweise Aufschluss über die *Zusammenhangskomponenten* eines Graphen. Eine Zusammenhangskomponente definiert ist als Teilgraph, in dem jedes Paar von Knoten über einen Pfad miteinander verbunden ist. Das folgende Lemma illustriert die Zusammenhangsinformationen, die sich aus der Laplace-Matrix ablesen lassen:

Lemma 2.3.1 (Zusammenhang und Laplace-Matrix [Lux07]). Gegeben sei $G = (V, E)$ ein schleifenloser Graph mit nicht-negativen Gewichten und Laplace-Matrix L . Die Anzahl an Zusammenhangskomponenten von G ist gleich der Anzahl an Eigenwerten von L , die gleich Null sind.

Da die Eigenwerte einer Matrix auch „Spektrum“ genannt werden und man bei den hier vorgestellten Methoden die Eigenwerte der Laplace-Matrix nutzt, wird diese Technik „spektrale Clusteranalyse“ genannt.

Die wichtigsten Algorithmen, die auf dem Spektrum der Laplace-Matrix der Eingabe basieren, wollen wir im Folgenden vorstellen. Wir beginnen mit den Algorithmen 4 und 5, welche die normalisierten Laplace-Matrizen verwenden und in den jeweils zitierten Papieren vorgestellt wurden.

Algorithmus 4: Nicht-normalisiertes spektrales Clustering [Lux07]

Eingabe: : Gewichtete Adjazenzmatrix $Adj^G \in \mathbb{R}^{n \times n}$, $k \in \mathbb{N}^+$

Ausgabe: : k Cluster C_1, \dots, C_k

- 1 Berechne die nicht-normalisierte Laplace-Matrix L des Graphen G mit Adj^G
 - 2 Bestimme Eigenvektoren u_1, \dots, u_k zu den k kleinsten Eigenwerten von L
 - 3 Sei $U \in \mathbb{R}^{n \times k}$ die Matrix mit u_1, \dots, u_k als Spaltenvektoren
 - 4 Für $i \leftarrow \{1, \dots, n\}$ sei $y_i \in \mathbb{R}^k$ der Vektor der i -ten Reihe von U
 - 5 Clustere die Punkte $(y_i)_{i \in \{1, \dots, n\}}$ mit einem k -means-Algorithmus in C_1, \dots, C_k
-

Algorithmus 5: Normalisiertes spektrales Clustering [SM00]

Eingabe: : Gewichtete Adjazenzmatrix $Adj^G \in \mathbb{R}^{n \times n}$, $k \in \mathbb{N}^+$

Ausgabe: : k Cluster C_1, \dots, C_k

- 1 Berechne die normalisierte Laplace-Matrix \mathcal{L}_{rw} des Graphen G mit Adj^G
 - 2 Bestimme Eigenvektoren u_1, \dots, u_k zu den k kleinsten Eigenwerten von \mathcal{L}_{rw}
 - 3 Sei $U \in \mathbb{R}^{n \times k}$ die Matrix mit u_1, \dots, u_k als Spaltenvektoren
 - 4 Für $i \leftarrow \{1, \dots, n\}$ sei $y_i \in \mathbb{R}^k$ der Vektor der i -ten Reihe von U
 - 5 Clustere die Punkte $(y_i)_{i \in \{1, \dots, n\}}$ mit einem k -means-Algorithmus in C_1, \dots, C_k
-

Der in [SM00] vorgestellte Algorithmus sieht in seiner ursprünglichen Form zunächst eine Bipartitionierung anhand des Eigenvektors mit dem zweitkleinsten Eigenwert vor (dies entspricht einer Approximation des *sparsest cut*), die anschließend gegebenenfalls weiter partitioniert wird. Dieses Vorgehen ist jedoch analog zu

dem hier vor uns angegebenen Algorithmus.

Algorithmus 6: Normalisiertes spektrales Clustering [NJW01]

Eingabe: : Gewichtete Adjazenzmatrix $Adj^G \in \mathbb{R}^{n \times n}$, $k \in \mathbb{N}^+$

Ausgabe: : k Cluster C_1, \dots, C_k

- 1 Berechne die normalisierte Laplace-Matrix \mathcal{L}_{sym} des Graphen G mit Adj^G
 - 2 Bestimme Eigenvektoren u_1, \dots, u_k zu den k kleinsten Eigenwerten von \mathcal{L}_{sym}
 - 3 Sei $U \in \mathbb{R}^{n \times k}$ die Matrix mit u_1, \dots, u_k als Spaltenvektoren
 - 4 Sei $T \in \mathbb{R}^{n \times k}$ die Matrix mit den Einträgen $T_{i,j} = \frac{U_{i,j}}{(\sum_k U_{i,k}^2)^{\frac{1}{2}}}$
 - 5 Für $i \leftarrow \{1, \dots, n\}$ sei $y_i \in \mathbb{R}^k$ der Vektor der i -ten Reihe von T
 - 6 Clustere die Punkte $(y_i)_{i \in \{1, \dots, n\}}$ mit einem k -means-Algorithmus in C_1, \dots, C_k
-

Der Algorithmus aus [NJW01] ist dem von [SM00] sehr ähnlich, verwendet jedoch die normalisierte Laplace-Matrix \mathcal{L}_{sym} . Außerdem ist für diesen Algorithmus eine zusätzliche Normalisierung nötig, deren Details in [Lux07] erläutert werden.

Das Berechnen aller Eigenvektoren benötigt asymptotischen Rechenaufwand von $\mathcal{O}(n^3)$, was für viele Anwendungen deutlich zu viel Rechenzeit ist. Beschränkt man sich auf die Berechnung von wenigen oder sogar nur einem einzelnen Eigenvektor, lässt sich die Komplexität signifikant verringern, beispielsweise unter Einsatz von Lanczos Algorithmus [Lan50]. Da für viele spektrale Clusteringalgorithmen die Berechnung der Eigenvektoren jedoch nach wie vor ein großer Flaschenhals ist, betrachten wir im nächsten Kapitel einen Ansatz, mit dem die Berechnung von Eigenvektoren für Graphclusterings prinzipiell vermieden werden kann.

3 Kernel k -means Methoden für spektrales Graphclustering

In diesem Kapitel wollen wir Verfahren zur spektralen Clusteranalyse von Graphen vorstellen und verbessern. Als Grundlage wollen wir den Kernel- k -means Algorithmus verwenden. Lloyds Algorithmus ist nach wie vor eine Heuristik, die in der Praxis mit geringen Ausführungszeiten Clusterings von akzeptabler Qualität berechnet. Unter Einsatz der Kernel Methode können wir die Qualität der berechneten Clusterings noch einmal steigern, ohne dabei signifikante Performanzeinbußen in Kauf nehmen zu müssen. Der Algorithmus ist zudem leicht zu implementieren. Gerade im Bereich der Clusteranalyse sollte dies nicht unterschätzt werden, da eine praktische Evaluation der Algorithmen in diesem anwendungsreichen Gebiet unerlässlich ist.

Das Kapitel ist folgendermaßen gegliedert: in Abschnitt 3.1 diskutieren wir, wie wir mit dem Kernel- k -means Algorithmus Graphclusteringprobleme (und umgekehrt) lösen können. In Abschnitt 3.2 zeigen wir, wie sich die Kernel Methode auch auf den Algorithmus k -means++ anwenden lässt und wie wir damit den Kernel- k -means Algorithmus für Graphclustering verbessern können.

3.1 Graphschnitte und Graphclustering mit gewichtetem Kernel k -means

Um die Graphclusteringprobleme aus Definition 2.2.3 mit dem Kernel- k -means Algorithmus zu lösen, müssen wir uns zunächst überlegen, wie wir diese in eine Eingabe für den Algorithmus transformieren können. Dass eine solche Transformation (bidirektional) möglich ist, haben bereits Dhillon, Guan und Kulis gezeigt [DGK04, DGK07]. Wir fassen daher zunächst den für uns relevanten Teil ihrer Ergebnisse zusammen.

Die grundlegende Idee besteht darin, zu zeigen, dass sowohl die Zielfunktion für Kernel- k -means als auch die (diversen) Zielfunktionen für Graphpartitionierungsprobleme als *Spurmaximierungsprobleme* formuliert werden können. Die *Spur* (engl. *trace*) einer quadratischen $n \times n$ -Matrix M ist dabei die Summe der Elemente auf der Hauptdiagonalen von M :

$$\text{trace}(M) = \sum_{i=1}^n M_{i,i}$$

Weiterhin gilt für zwei quadratische Matrizen M, M' :

1. $\text{trace}(MM') = \text{trace}(M'M)$
2. $\text{trace}(M^T M) = \|M\|_F^2$
3. $\text{trace}(M + M') = \text{trace}(M) + \text{trace}(M')$

Wir formen zunächst die Kernel- k -means Zielfunktion um und widmen uns dann den Graphpartitionierungsproblemen.

Gewichtetes Kernel- k -means. Für die Umformung des k -means-Problems legen wir das *gewichtete* Kernel- k -means Problem (engl. *weighted kernel k -means*, kurz WKKM) zu Grunde. Dieses lautet analog zu Definition 2.3.2 wie folgt. Unsere Eingabe besteht aus der gewichteten Punktmenge als $n \times d$ -Matrix A , einer Gewichtsfunktion w , die einen Punkt-Zeilenvektor \mathbf{p} aus A auf ein reelles Gewicht abbildet, und der Clusteranzahl k . Die Zielfunktion besteht wie gewohnt darin, die quadrierten euklidischen Distanzen der Punkte eines Clusters zu ihrem Clusterzentrum zu minimieren:

$$\min D(\{S_i\}_{i=1}^k) = \sum_{i=1}^k \sum_{\mathbf{p} \in S_i} w(\mathbf{p}) \|\phi(\mathbf{p}) - \mathbf{c}_i\|^2 \quad (3.1)$$

Dabei sind die S_1, \dots, S_k die Cluster, wobei ein Cluster S_i das Zentrum

$$\mathbf{c}_i = \frac{\sum_{\mathbf{p} \in S_i} w(\mathbf{p}) \phi(\mathbf{p})}{\sum_{\mathbf{p} \in S_i} w(\mathbf{p})}$$

hat. Wir setzen auch hier wieder den Kernel-Trick für die Distanzberechnungen ein. Dazu berechnen wir initial n^2 viele Skalarprodukte $\langle \phi(\mathbf{p}_i), \phi(\mathbf{p}_j) \rangle$ und speichern diese in einer $n \times n$ -Kernelmatrix K mit $K_{i,j} = \kappa(\mathbf{p}_i, \mathbf{p}_j) = \langle \phi(\mathbf{p}_i), \phi(\mathbf{p}_j) \rangle$. Die Distanzberechnung von einem Punkt zu seinem Clusterzentrum können wir ebenfalls analog zu 2.2 mit Hilfe der Kernelmatrix umformulieren. Sie unterscheidet sich von dieser nur darin, dass die Punkte gewichtet sind:

$$\|\phi(\mathbf{p}_j) - \mathbf{c}_i\|^2 = K_{i,i} - \frac{2 \sum_{\mathbf{p}_j \in S_i} w(\mathbf{p}_j) K_{i,j}}{\sum_{\mathbf{p}_j \in S_i} w(\mathbf{p}_j)} + \frac{\sum_{\mathbf{p}_j, \mathbf{p}_l \in S_i} w(\mathbf{p}_j) w(\mathbf{p}_l) K_{j,l}}{\left(\sum_{\mathbf{p}_j \in S_i} w(\mathbf{p}_j) \right)^2} \quad (3.2)$$

Die Erweiterung von Algorithmus 3 zu Algorithmus 7 ergibt sich ebenfalls unmittelbar.

Algorithmus 7: Weighted Kernel- k -means

Eingabe: : Kernel-Matrix K , $k \in \mathbb{N}^+$, Punktgewichte w

Ausgabe: : k -means-Clustering der Eingabepunktmenge

```
1 Wähle  $k$  initiale Cluster  $S_1^{(0)}, \dots, S_k^{(0)}$ ;  $t \leftarrow 0$ 
2 foreach  $\mathbf{p}$  do
3   Bestimme  $c^*(\mathbf{p}) \leftarrow \arg \min_{i=1}^k \|\phi(\mathbf{p}) - \mathbf{c}_i\|^2$  mit 3.2
4 for  $i \leftarrow 1$  to  $k$  do
5    $S_i^{(t+1)} \leftarrow \{\mathbf{p} : c^*(\mathbf{p}) = i\}$ 
6  $t \leftarrow t + 1$ 
7 if  $t < t_{\max}$  and nicht konvergiert then
8   go to 2
9 else
10  return  $S_1^{(t)}, \dots, S_k^{(t)}$ 
```

Gewichtetes Kernel- k -means als Spurmaximierung. Wir formen als nächstes die Zielfunktion des gewichteten Kernel- k -means Problems in ein Spurmaximierungsproblem um. Dazu gehen wir zunächst davon aus, dass die Punktgewichte in einer $n \times n$ -Diagonalmatrix W gespeichert sind, sodass $W_{l,l} = w(\mathbf{p}_j)$ für alle Punkte \mathbf{p}_j und die Punktgewichte des Clusters S_i in der Diagonalmatrix W^i mit $W_{m,m}^i = w(\mathbf{p}_m)$ für alle Punkte $\mathbf{p}_m \in S_i$ gespeichert sind.

Sei nun s_i die Summe der Gewichte der Punkte in Cluster i , also

$$s_i = \sum_{\mathbf{p}_j \in S_i} w(\mathbf{p}_j) = \sum_{\mathbf{p}_j \in S_i} W_{j,j}^i$$

Wir betrachten die $n \times k$ -Matrix Z mit den Einträgen

$$Z_{j,i} = \begin{cases} \frac{1}{\sqrt{s_i}}, & \text{falls } \mathbf{p}_j \in S_i, \\ 0 & \text{sonst} \end{cases}$$

sowie die (nicht explizit berechnete) Matrix $\Phi_i = [\phi(\mathbf{p}_1), \dots, \phi(\mathbf{p}_m)]$ von abgebildeten Punkten des Clusters $S_i = \{\mathbf{p}_1, \dots, \mathbf{p}_m\}$. Mit diesen können wir die Zielfunktion des gewichteten Kernel- k -means Problems (3.1) umformen. Wir erinnern uns, dass diese die Minimierung der Summe der Intracusterabstände ist. Notieren wir den Intracusterabstand eines Clusters S_i mit $d(S_i) = \sum_{\mathbf{p} \in S_i} w(\mathbf{p}) \|\phi(\mathbf{p}) - \mathbf{c}_i\|^2$, können wir für den Zielfunktionswert $D(\{S_i\}_{i=1}^k)$ verkürzt schreiben:

$$D(\{S_i\}_{i=1}^k) = \sum_{i=1}^k d(S_i)$$

Der Zentroidvektor c_i des Clusters S_i berechnet sich in unserer ursprünglichen Problemformulierung durch

$$\mathbf{c}_i = \frac{\sum_{\mathbf{p} \in S_i} w(\mathbf{p}) \phi(\mathbf{p})}{\sum_{\mathbf{p} \in S_i} w(\mathbf{p})}$$

Diese Berechnungsvorschrift ist äquivalent zu

$$\mathbf{c}_i = \Phi_i \frac{W_j \mathbf{e}}{s_j},$$

wobei \mathbf{e} der Einheitsvektor der passenden Dimension ist. Wir können alternativ unter Einsatz der Matrix Z auch direkt eine Matrix C aller Clusterzentren bestimmen, wenn wir für die gesamte Eingabepunktmenge die (ebenfalls nicht explizit berechnete) Abbildungsmatrix $\Phi = [\Phi_1, \dots, \Phi_k] = [\phi(\mathbf{p}_1), \dots, \phi(\mathbf{p}_n)]$ verwenden:

$$C = \Phi W Z Z^T$$

In diesem Fall entspricht das Zentrum eines einzelnen Clusters S_i

$$\mathbf{c}_i = (\Phi W Z Z^T)_{\cdot i}$$

wobei wir für eine Matrix M mit $M_{\cdot i}$ die i -te Spalte von M meinen. Damit ist es uns nun möglich, den Zielfunktionswert so zu formulieren, dass wir später eine Äquivalenz zur Graphpartitionierung sehen werden. Mit $Y = W^{\frac{1}{2}} Z$ gilt:

$$\begin{aligned} D(\{S_i\}_{i=1}^k) &= \sum_{i=1}^k \sum_{\mathbf{p} \in S_i} w(\mathbf{p}) \|\phi(\mathbf{p}) - \mathbf{c}_i\|^2 \\ &= \sum_{j=1}^n W_{j,j} \left\| \Phi_{\cdot j} - (\Phi W Z Z^T)_{\cdot j} \right\|^2 \\ &= \sum_{j=1}^n W_{j,j} \left\| \Phi_{\cdot j} - \left(\Phi W^{\frac{1}{2}} Y Y^T W^{-\frac{1}{2}} \right)_{\cdot j} \right\|^2 \\ &= \sum_{j=1}^n \left\| \Phi_{\cdot j} (W_{j,j})^{\frac{1}{2}} - \left(\Phi W^{\frac{1}{2}} Y Y^T \right)_{\cdot j} \right\|^2 \\ &= \left\| \Phi W^{\frac{1}{2}} - \Phi W^{\frac{1}{2}} Y Y^T \right\|_F^2 \end{aligned}$$

Der ursprüngliche Zielfunktionswert ist damit auf die (quadrierte) Frobenius-Norm von $\Phi W^{\frac{1}{2}} - \Phi W^{\frac{1}{2}} Y Y^T$ transformiert. Wir können nun die oben genannte Verbindung zwischen der Frobenius-Norm und der Matrix-Spur ausnutzen, um zur

beabsichtigten Spur-Maximierung zu gelangen. Wir erinnern uns, dass für eine Matrix M gilt, dass $\text{trace}(M^T M) = \|M\|_F^2$. Wir können somit weiter umformen:

$$\begin{aligned} D(\{S_i\}_{i=1}^k) &= \text{trace}(W^{\frac{1}{2}} \Phi^T \Phi W^{\frac{1}{2}} - W^{\frac{1}{2}} \Phi^T \Phi W^{\frac{1}{2}} Y Y^T \\ &\quad - Y Y^T W^{\frac{1}{2}} \Phi^T \Phi W^{\frac{1}{2}} + Y Y^T W^{\frac{1}{2}} \Phi^T \Phi W^{\frac{1}{2}} Y Y^T) \\ &= \text{trace}(W^{\frac{1}{2}} \Phi^T \Phi W^{\frac{1}{2}}) - \text{trace}(Y^T W^{\frac{1}{2}} \Phi^T \Phi W^{\frac{1}{2}} Y) \end{aligned}$$

In der Kernel-Matrix K speichern wir sämtliche Skalarprodukte der projizierten Punkte, demnach ist per Definition $K = \Phi^T \Phi$. Die Kernel-Matrix ist genau wie die Gewichts-Matrix W konstant, da beide Teil der Eingabe sind. Dementsprechend ist auch $\text{trace}(W^{\frac{1}{2}} K W^{\frac{1}{2}})$ eine Konstante. Die Minimierungszielfunktion 3.1 ist daher gleichzusetzen mit der Maximierungszielfunktion des folgenden Spurmaximierungsproblems:

$$\max_Y \text{trace}(Y^T W^{\frac{1}{2}} K W^{\frac{1}{2}} Y) \quad (3.3)$$

Da das eigentliche Clustering von Algorithmus 7 offenbar einer spurmaximierenden Wahl der Matrix Y gleichkommt, wollen wir abschließend noch kurz einen genaueren Blick auf diese werden. Die Matrix $Y = W^{\frac{1}{2}} Z$ ist eine orthonormale Diagonalmatrix, das heißt $Y^T Y = I$. Ihre Einträge sehen dementsprechend wie folgt aus:

$$Y = \begin{pmatrix} \frac{W_1^{\frac{1}{2}} \mathbf{e}}{\sqrt{s_1}} & & & \\ & \frac{W_2^{\frac{1}{2}} \mathbf{e}}{\sqrt{s_2}} & & \\ & & \dots & \\ & & & \frac{W_k^{\frac{1}{2}} \mathbf{e}}{\sqrt{s_k}} \end{pmatrix}$$

Graphpartitionierung als Spurmaximierung. Es ist möglich, für alle in Definition 2.2.3 vorgestellten Problemstellungen ein äquivalentes Spurmaximierungsproblem anzugeben. Wir beschränken uns hier zunächst auf die beiden Zielfunktionen mit praktischer Relevanz für das Evaluations-Kapitel dieser Arbeit, nämlich die Ratio Association und den Normalized Cut. Wir beginnen mit der Ratio Association.

Wir führen zunächst einen Indikator-Vektor \mathbf{x}_c der Dimension n für jede Partition beziehungsweise jedes Cluster c ein, sodass $\mathbf{x}_c(i) = 1$ genau dann, wenn das Cluster c den Knoten i enthält. Die Zielfunktion bei der Ratio Association lautet

$$\max \sum_{c=1}^k \frac{w(V_c, V_c)}{|V_c|}$$

Wir nehmen an, dass A die (gewichtete) Adjazenzmatrix des Eingabegraphen ist und definieren zusätzlich

$$\tilde{\mathbf{x}}_c = \frac{\mathbf{x}_c}{(\mathbf{x}_c^T \mathbf{x}_c)^{\frac{1}{2}}}$$

Das Produkt $\mathbf{x}_c^T \mathbf{x}_c$ entspricht der Größe der Partition c :

$$|V_c| = \mathbf{x}_c^T \mathbf{x}_c$$

Die Summe der Kantengewichte der Partition $w(V_c, V_c)$ können wir über folgendes Produkt abbilden:

$$w(V_c, V_c) = \mathbf{x}_c^T A \mathbf{x}_c$$

Damit können wir eine erste Transformation der Zielfunktion vornehmen:

$$\max \sum_{c=1}^k \frac{w(V_c, V_c)}{|V_c|} \Leftrightarrow \max \sum_{c=1}^k \frac{\mathbf{x}_c^T A \mathbf{x}_c}{\mathbf{x}_c^T \mathbf{x}_c} \Leftrightarrow \max \sum_{c=1}^k \tilde{\mathbf{x}}_c^T A \tilde{\mathbf{x}}_c$$

Betrachten wir \tilde{X} als Matrix deren c -te Spalte der Vektor $\tilde{\mathbf{x}}_c$ ist, können wir die Zielfunktion umformen zu

$$\max_{\tilde{X}} \text{trace}(\tilde{X}^T A \tilde{X}) \quad (3.4)$$

Wenn wir in 3.3 für die Gewichtsmatrix W die Einheitsmatrix wählen, also $W = \mathbb{1}_n$ und die Adjazenzmatrix als Kernel-Matrix setzen, sodass $K = A$, dann sind 3.4 und 3.3 äquivalent. Auf diese Weise können wir das Ratio Association Problem mit unserem Algorithmus für gewichtetes Kernel- k -means lösen.

Für den Normalized Cut gehen wir ähnlich vor. Unsere ursprüngliche Zielfunktion lautet

$$\min \sum_{c=1}^k \frac{w(V_c, V \setminus V_c)}{w(V_c, V)}$$

Zunächst wollen wir die Funktion in ein Maximierungsproblem umwandeln. Wir beobachten, dass die Zielfunktion äquivalent ist zu

$$\max \sum_{c=1}^k \frac{w(V_c, V_c)}{w(V_c, V)}$$

Wir benötigen zur Umformung des Normalized Cut die diagonale Gradmatrix D des Graphen, der der Adjazenzmatrix A zu Grunde liegt. Wir verwenden wieder

die Indikator-Vektoren und definieren zusätzlich

$$\hat{\mathbf{x}}_c = \frac{\mathbf{x}_c}{(\mathbf{x}_c^T D \mathbf{x}_c)^{\frac{1}{2}}}$$

Unter Verwendung der Gradmatrix können wir die Variante des Normalized Cut als Maximierungsproblem wie folgt umformen:

$$\max \sum_{c=1}^k \frac{w(V_c, V_c)}{w(V_c, V)} \Leftrightarrow \max \sum_{c=1}^k \frac{\mathbf{x}_c^T A \mathbf{x}_c}{\mathbf{x}_c^T D \mathbf{x}_c} \Leftrightarrow \max \sum_{c=1}^k \hat{\mathbf{x}}_c^T A \hat{\mathbf{x}}_c$$

Wir nehmen analog zur Ratio Association an, dass die Matrix \hat{X} die Spaltenmatrix der Vektoren $\hat{\mathbf{x}}_c$ ist. Mit $\tilde{Y} = D^{\frac{1}{2}} \hat{X}$ lautet das Spurmaximierungsproblem für den Normalized Cut folgendermaßen:

$$\max_{\tilde{Y}} \text{trace}(\tilde{Y}^T D^{-\frac{1}{2}} A D^{-\frac{1}{2}} \tilde{Y}) \quad (3.5)$$

Das Problem in 3.5 ist äquivalent zu 3.3, wenn wir als Gewichtsmatrix $W = D$ und als Kernel-Matrix $K = D^{-1} A D^{-1}$ setzen.

Wir wollen abschließend noch anmerken, dass die Adjazenzmatrix A positiv definit sein muss, damit sie als Kernel-Matrix in Algorithmus 7 garantiert zu einer iterativen Verringerung des gewichteten Kernel- k -means Zielfunktionswertes führt. Um dies in der Praxis sicherzustellen ist bei der Wahl der Kernel-Matrix gegebenenfalls das Aufaddieren einer „Verschiebung“ auf die ursprüngliche Kernel-Matrix nötig. Wir gehen auf dieses Detail im Experimente-Kapitel noch genauer ein.

3.2 Die Kernel Methode für k -means++

Ein wichtiges Merkmal von Algorithmus 7 ist, dass er grundsätzlich ohne Methoden der spektralen Clusteranalyse auskommt. Wir müssen uns jedoch wie auch bei Lloyds Algorithmus ohne Kernel Methode mit der geeigneten initialen Wahl der Clusterzentren beschäftigen. Zunächst einmal wählen wir diese wie schon zuvor zufällig aus der Eingabepunktmenge. Diese Vorgehensweise ist durchaus praktikabel, die experimentelle Evaluation aus [DGK04, DGK07] zeigt jedoch, dass sich mit einer geschickteren Wahl die Qualität des letztlich berechneten Clusterings noch signifikant steigern lässt. Die Autoren verwenden für die Initialisierung die spektrale Clusteranalyse, was wiederum den Aufwand der Eigenvektor-Berechnung mit sich bringt.

Konkret wird mit einem Verfahren aus [GL96] eine Problemrelaxion von 3.3/ 3.4 beziehungsweise von 3.3/ 3.5 gelöst. Relaxieren wir die Wahl von Y beziehungsweise \tilde{Y} auf eine beliebige orthonormale Matrix, dann sind diese von der Form $Y = V_k Q$ beziehungsweise $\tilde{Y} = \tilde{V}_k Q$. Dabei ist Q eine beliebige orthogonale $k \times k$ -Matrix, V_k die Matrix der k größten Eigenvektoren von $W^{\frac{1}{2}} K W^{\frac{1}{2}}$ und \tilde{V}_k die Matrix der k größten Eigenvektoren von $D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$.

Die Berechnung der jeweiligen k größten Eigenvektoren wird mit wachsender Clusteranzahl k aufwändig, wie wir schon in Abschnitt 2.3 beobachtet haben. Intuitiv ist daher die Fragestellung, ob es möglich ist, eine Initialisierung ohne Eigenvektorberechnungen vorzunehmen, die dennoch deutlich verbesserte Clusterings mit dem Kernel- k -means-Algorithmus ermöglicht.

In Abschnitt 2.1 haben wir bereits den Algorithmus k -means++ von Arthur und Vassilvitskii [AV07] kennengelernt. Die initiale Wahl der Cluster erfolgt hier gemäß einer Wahrscheinlichkeitsverteilung der Punkte als mögliche Clusterzentren, die proportional zur quadrierten euklidischen Distanz zum jeweils nächstgelegenen bereits gewählten Zentrum ist. Wir zeigen nun, dass die Kernel Methode auch auf den Algorithmus k -means++ angewendet werden kann. Damit ist es uns möglich, Algorithmus 3 (und damit auch Algorithmus 7) in eine Kernel-basierte Variante von k -means++ umzuwandeln, die wir Kernel- k -means++ nennen.

Wir erinnern uns, dass wir im Zusammenhang mit k -means++ für einen Punkt x aus der Eingabe-Punktmenge P mit $D(x)$ die geringste Distanz von x zum nächstgelegenen bereits gewählten Zentrum notieren. Der Punkt x wird als nächstes Zentrum mit der Wahrscheinlichkeit

$$\frac{D(x)^2}{\sum_{y \in P} D(y)^2}$$

gewählt. Wir können demnach die Kernel Methode auf den Algorithmus k -means++ anwenden, wenn wir für jeden Eingabepunkt x den Wert $D(\phi(x))$ beziehungsweise $D(\phi(x))^2$ berechnen können. Dazu zeigen wir nun, dass wir diese Werte wie schon zuvor beim Kernel- k -means-Algorithmus ausschließlich über die Kernelfunktion κ ohne explizite Kenntnis der abgebildeten Punkte oder der Abbildung ϕ selbst berechnen können.

Wir nehmen zunächst an, dass die Zentren c_1, \dots, c_l mit $l \in \{1, \dots, k-1\}$ bereits gewählt worden sind. Dann können wir $D(x)$ formal notieren als das Minimum über alle quadrierten euklidischen Distanzen von x zu jedem der bisher gewählten Zentren c_i :

$$D(x) = \min_{i=1 \dots l} \|x - c_i\|^2$$

Wenden wir die Kernel Methode an, berechnet sich $D(\phi(x))$ dementsprechend folgendermaßen:

$$D(\phi(x)) = \min_{i=1..l} \|\phi(x) - \phi(c_i)\|^2$$

Demnach genügt es zu zeigen, dass wir $\|\phi(x) - \phi(c_i)\|^2$ ausschließlich mit der Kernelfunktion berechnen können. Wir haben bereits gesehen, dass die Kernelfunktion gerade dafür gedacht ist. Es gilt:

$$\|\phi(x) - \phi(c_i)\|^2 = \kappa(x, x) - 2\kappa(x, c_i) + \kappa(c_i, c_i) \quad (3.6)$$

Mit den selben Überlegungen können wir auch zeigen, dass sich der Wert von $D(\phi(x))^2$ mit der Kernel Methode berechnen lässt. Zunächst lässt sich auch $D(\phi(x))^2$ als Minimum formulieren:

$$D(\phi(x))^2 = \left(\min_{i=1..l} \|\phi(x) - \phi(c_i)\|^2 \right)^2$$

Dazu müssen wir effektiv $(\|\phi(x) - \phi(c_i)\|^2)^2$ berechnen. Auch dies lässt sich mit der Kernelfunktion umsetzen:

$$\begin{aligned} (\|\phi(x) - \phi(c_i)\|^2)^2 &= (\kappa(x, x) - 2\kappa(x, c_i) + \kappa(c_i, c_i))^2 \\ &= (\kappa(x, x))^2 - 4\kappa(x, x)\kappa(x, c_i) + 2\kappa(x, x)\kappa(c_i, c_i) \\ &\quad + 4\kappa(x, c_i)\kappa(x, c_i) - 4\kappa(x, c_i)\kappa(c_i, c_i) + (\kappa(c_i, c_i))^2 \end{aligned} \quad (3.7)$$

Der Algorithmus Kernel- k -means++ wird damit die logische Weiterentwicklung von k -means++ auf Basis der Kernel Methode.

Algorithmus 8: Kernel- k -means++

Eingabe: : $P \subseteq \mathbb{R}^d, k \in \mathbb{N}^+$, Kernelfunktion κ

Ausgabe: : k initiale Clusterzentren für P

- 1 Wähle c_1 zufällig gleichverteilt aus P
 - 2 **for** $i \leftarrow 2$ **to** k **do**
 - 3 Berechne $D(\phi(x'))^2$ mit 3.7 für alle $x' \in P \setminus \{c_1, \dots, c_{i-1}\}$
 - 4 Wähle den Punkt $x' \in P \setminus \{c_1, \dots, c_{i-1}\}$ als Zentrum c_i mit
 Wahrscheinlichkeit $\frac{D(\phi(x'))^2}{\sum_{x \in P} D(\phi(x))^2}$
 - 5 Führe den Kernel- k -means-Algorithmus mit den initialen Clusterzentren
 c_1, \dots, c_k aus.
-

Wir zeigen nun, dass auch Kernel- k -means++ eine $\mathcal{O}(\log k)$ -Approximation für das k -means-Problem ist. Strukturell gehen wir dabei analog zu dem Papier von Arthur und Vassilvitskii [AV07] vor, da sich durch die Kernel Methode keine wesentlichen

Änderungen ergeben. Für ein Clustering $C = \{C_1, \dots, C_k\}$ mit den jeweiligen Clusterzentren c_1, \dots, c_k bezeichnen wir mit dem *Potenzial* θ den Zielfunktionswert des Clusterings bezüglich der k -means-Zielfunktion, das heißt

$$\theta(C) = \sum_{i=1}^k \theta(C_i) = \sum_{i=1}^k \sum_{p \in C_i} \|p - c_i\|^2$$

Das Potenzial des optimalen Clusterings C_{OPT} notieren wir mit θ_{OPT} . Bevor wir mit der Analyse beginnen, benötigen wir folgendes Lemma, das einen beliebigen Punkt mit einer Punktmenge und ihrem Zentroiden in Relation setzt:

Lemma 3.2.1. Sei S eine Punktmenge mit dem geometrischen Zentroiden $z(S)$. Weiterhin sei \mathbf{q} ein beliebiger Punkt. Es gilt:

$$\sum_{\mathbf{x} \in S} \|\mathbf{x} - \mathbf{q}\|^2 - \sum_{\mathbf{x} \in S} \|\mathbf{x} - z(S)\|^2 = |S| \cdot \|z(S) - \mathbf{q}\|^2$$

Wir zeigen nun das eigentliche Ergebnis.

Satz 3.2.1 (Kernel- k -means++ ist $\mathcal{O}(\log k)$ -kompetitiv). Wenn das Clustering C mit Kernel- k -means++ berechnet wird, dann gilt für das Potenzial θ von C :

$$E[\theta(C)] \leq 8(\ln k + 2)\theta_{\text{OPT}}$$

Der Beweis erfolgt in drei Schritten. Wir zeigen zunächst, dass das erste zufällig gleichverteilt gewählte Clusterzentrum die Approximationseigenschaft erfüllt. Die selbe Eigenschaft zeigen wir im zweiten Schritt für die übrigen Zentren, die gemäß der D^2 -Gewichtung gewählt werden. Im dritten und letzten Schritt beweisen wir mittels Induktion das Gesamtergebnis.

Lemma 3.2.2 (Das erste Zentrum ist kompetitiv). Sei A ein beliebiges Cluster von C_{OPT} und sei C ein Clustering mit genau einem Zentrum, das zufällig gleichverteilt aus A gewählt wird. Dann gilt:

$$E[\theta(A)] = 2\theta_{\text{OPT}}(A)$$

Beweis. Sei $z(A)$ der geometrische Zentroid der Punkte in A , das heißt

$$z(A) = \frac{1}{|A|} \sum_{\mathbf{p} \in A} \mathbf{p}$$

Mit Lemma 3.2.1 gilt:

$$\begin{aligned}
& \sum_{\mathbf{p} \in A} \frac{1}{|A|} \cdot \left(\sum_{\mathbf{p}' \in A} \|\mathbf{p} - \mathbf{p}'\|^2 \right) \\
&= \frac{1}{|A|} \sum_{\mathbf{p} \in A} \left(\sum_{\mathbf{p}' \in A} \|\mathbf{p}' - z(A)\|^2 + |A| \cdot \|\mathbf{p} - z(A)\|^2 \right) \\
&= 2 \sum_{\mathbf{p} \in A} \|\mathbf{p} - z(A)\|^2 \\
&= E[\theta(A)] \quad \square
\end{aligned}$$

Wir zeigen nun, dass die Approximationseigenschaft auch für die übrigen Zentren, die gemäß der D^2 -Gewichtung gewählt werden, gilt.

Lemma 3.2.3 (Die übrigen Zentren sind kompetitiv). Sei A ein beliebiges Cluster von C_{OPT} und sei C ein beliebiges Clustering. Wenn ein Zentrum aus A zufällig nach der D^2 -Gewichtung zu C hinzugefügt wird, dann gilt:

$$E[\theta(A)] \leq 8\theta_{\text{OPT}}(A)$$

Beweis. Ein Punkt $\mathbf{p} \in A$ wird mit der Wahrscheinlichkeit

$$\frac{D(\phi(\mathbf{p}))^2}{\sum_{\mathbf{p}' \in A} D(\phi(\mathbf{p}'))^2}$$

als Zentrum gewählt. Wenn der Punkt \mathbf{p} als Zentrum gewählt wird, beträgt der Potenzial-Wert eines Punktes \mathbf{p}' , der dem Zentrum \mathbf{p} zugeordnet ist, genau

$$\min\{D(\phi(\mathbf{p}')), \|\phi(\mathbf{p}') - \phi(\mathbf{p})\|^2\}$$

Der Erwartungswert für das Potenzial des gesamten Clusters A ist dementsprechend:

$$E[\theta(A)] = \sum_{\mathbf{p} \in A} \frac{D(\phi(\mathbf{p}))^2}{\sum_{\mathbf{p}' \in A} D(\phi(\mathbf{p}'))^2} \sum_{\mathbf{p}' \in A} \min\{D(\phi(\mathbf{p}')), \|\phi(\mathbf{p}') - \phi(\mathbf{p})\|^2\}$$

Für die Berechnung von $D(\phi(\mathbf{p}'))$ und der Distanzen $\|\phi(\mathbf{p}') - \phi(\mathbf{p})\|$ verwenden wir im Kernel- k -means+-Algorithmus die Kernelfunktion κ . Wenn die Kernelfunktion, beziehungsweise die zugehörige Kernelmatrix *positiv definit* ist (für Details siehe beispielsweise [SC04]), gilt für diese die Cauchy-Schwarz-Ungleichung, also $|\kappa(x, y)|^2 \leq \kappa(x, x) \cdot \kappa(y, y)$. Aus der Cauchy-Schwarz-Ungleichung folgen die Dreiecks-Ungleichung, sowie die Power-Mean-Ungleichung bezüglich der Kernel-

funktion. Die üblichen Kernelfunktionen wie beispielsweise in Tabelle 1 sowie die in unseren Experimenten eingesetzten Kernel haben eine positiv definite zugehörige Kernelmatrix. Wir können daher an dieser Stelle die Dreiecks-Ungleichung anwenden und erhalten für alle \mathbf{p}, \mathbf{p}' :

$$D(\phi(\mathbf{p})) \leq D(\phi(\mathbf{p}')) + \|\phi(\mathbf{p}') - \phi(\mathbf{p})\|$$

Für die D^2 -Gewichtung gilt zudem mit der Power-Mean-Ungleichung:

$$D(\phi(\mathbf{p}))^2 \leq 2D(\phi(\mathbf{p}'))^2 + 2\|\phi(\mathbf{p}') - \phi(\mathbf{p})\|^2$$

Summiert über alle Punkte $\mathbf{p} \in A$ erhalten wir somit:

$$D(\phi(\mathbf{p}))^2 \leq \frac{2}{|A|} \sum_{\mathbf{p}' \in A} D(\phi(\mathbf{p}'))^2 + \frac{2}{|A|} \sum_{\mathbf{p}' \in A} \|\phi(\mathbf{p}') - \phi(\mathbf{p})\|^2$$

Wir schätzen nun $E[\theta(A)]$ nach oben ab:

$$\begin{aligned} E[\theta(A)] &\leq \frac{2}{|A|} \cdot \sum_{\mathbf{p} \in A} \frac{\sum_{\mathbf{p}' \in A} D(\phi(\mathbf{p}'))^2}{\sum_{\mathbf{p}' \in A} D(\phi(\mathbf{p}'))^2} \cdot \sum_{\mathbf{p}' \in A} \min\{D(\phi(\mathbf{p}')), \|\phi(\mathbf{p}') - \phi(\mathbf{p})\|^2\} \\ &\quad + \frac{2}{|A|} \cdot \sum_{\mathbf{p} \in A} \frac{\sum_{\mathbf{p}' \in A} \|\phi(\mathbf{p}') - \phi(\mathbf{p})\|^2}{\sum_{\mathbf{p}' \in A} D(\phi(\mathbf{p}'))^2} \\ &\quad \cdot \sum_{\mathbf{p}' \in A} \min\{D(\phi(\mathbf{p}')), \|\phi(\mathbf{p}') - \phi(\mathbf{p})\|^2\} \end{aligned}$$

Im ersten Summanden schätzen wir

$$\min\{D(\phi(\mathbf{p}')), \|\phi(\mathbf{p}') - \phi(\mathbf{p})\|^2\} \leq \|\phi(\mathbf{p}') - \phi(\mathbf{p})\|^2$$

nach oben ab, und im zweiten Summanden

$$\min\{D(\phi(\mathbf{p}')), \|\phi(\mathbf{p}') - \phi(\mathbf{p})\|^2\} \leq D(\phi(\mathbf{p}'))$$

Damit erhalten wir vereinfacht:

$$E[\theta(A)] \leq \frac{4}{|A|} \sum_{\mathbf{p} \in A} \sum_{\mathbf{p}' \in A} \|\phi(\mathbf{p}') - \phi(\mathbf{p})\|^2$$

Nun können wir Lemma 3.2.2 anwenden und erhalten:

$$E[\theta(A)] \leq 8\theta_{\text{OPT}}(A)$$

□

Mit den Lemmas 3.2.2 und 3.2.3 folgt, dass die Wahl der Zentren von Kernel- k -means++ der Approximationseigenschaft genügt, solange der Algorithmus Zentren aus jedem Cluster des optimalen Clusterings C_{OPT} auswählt. Im dritten und letzten Schritt unserer Analyse zeigen wir nun, dass der entstehende Fehler im anderen Fall nach oben insgesamt durch $\mathcal{O}(\log k)$ beschränkt ist.

Lemma 3.2.4. Sei C ein beliebiges Clustering. Gegeben seien $u > 0$ viele freie Cluster (die übrigen Cluster nennen wir abgedeckt) von C_{OPT} . P_u seien die Punkte in diesen Clustern. Weiterhin sei $P_c = P - P_u$. Wenn $t \leq u$ zufällig gemäß der D^2 -Gewichtung gewählte Zentren zu C hinzugefügt werden, dann gilt für das erwartete neue Potenzial $E[\theta'(C)]$ von C :

$$E[\theta'(C)] \leq (\theta(P_c) + 8\theta_{\text{OPT}}(P_u)) \cdot (1 + H_t) + \frac{u-t}{u} \cdot \theta(P_u)$$

Dabei ist H_t die t -te harmonische Zahl $H_t = 1 + \frac{1}{2} + \dots + \frac{1}{t}$.

Beweis. Wir beweisen die Aussage durch vollständige Induktion für $(t-1, u)$ und $(t-1, u-1)$.

Induktionsanfang. Falls $t = 0$ und $u > 0$, gilt:

$$1 + H_t = \frac{u-t}{u} = 1$$

Wir betrachten nun den Fall, dass $t = u = 1$. Wir wählen genau ein neues Zentrum aus der Menge von genau einem freien Cluster mit der Wahrscheinlichkeit $\frac{\theta(P_u)}{\theta(C)}$. In diesem Fall folgt aus Lemma 3.2.3, dass

$$E[\theta'(C)] \leq \theta(P_c) + 8\theta_{\text{OPT}}(P_u)$$

Auch wenn wir ein Zentrum aus einem abgedeckten Cluster wählen, gilt $\theta'(C) \leq \theta(C)$. Somit können wir das erwartete neue Potenzial in diesem Fall nach oben abschätzen:

$$\begin{aligned} E[\theta'(C)] &\leq \frac{\theta(P_u)}{\theta(C)} \cdot (\theta(P_c) + 8\theta_{\text{OPT}}(P_u)) + \frac{\theta(P_u)}{\theta(C)} \cdot \theta(C) \\ &\leq 2\theta(P_c) + 8\theta_{\text{OPT}}(P_u) \end{aligned}$$

Bezüglich der Behauptung gilt in diesem Fall, dass $1 + H_t = 1 + 1 = 2$ (und $\frac{u-t}{u} = \frac{1-1}{1} = 0$). Daher gilt die Behauptung auch in diesem Fall.

Induktionsschritt. Wir nehmen eine Fallunterscheidung bezüglich des Clusters vor, aus dem wir unser erstes Zentrum wählen.

- (a) Wir nehmen an, das erste gewählte Zentrum stammt aus einem abgedeckten Cluster. Analog zum Induktionsanfang ist die Wahrscheinlichkeit dafür $\frac{\theta(P_c)}{\theta(C)}$. Da das Zentrum aus einem abgedeckten Cluster stammt, kann es $\theta(C)$ nicht vergrößern. Wir verwenden die Induktionsannahme mit $t - 1$. Folglich tragen wir erwartet höchstens den folgenden Betrag zu $E[\theta'(C)]$ bei:

$$\frac{\theta(P_c)}{\theta(C)} \cdot \left((\theta(P_c) + 8\theta_{\text{OPT}}(P_u)) \cdot (1 + H_{t-1}) + \frac{u - (t - 1)}{u} \cdot \theta(P_u) \right)$$

- (b) Wir nehmen an, das erste gewählte Zentrum stammt aus einem freien Cluster A . Auch hier ist die Wahrscheinlichkeit für dieses Ereignis $\frac{\theta(A)}{\theta(C)}$. Für alle $a \in A$ sei p_a die Wahrscheinlichkeit, dass wir a als Zentrum wählen. Weiterhin sei θ_a das Potenzial $\theta(A)$, nachdem a als Zentrum gewählt wurde. Wir wenden wieder die Induktionsannahme an, setzen jedoch $t = t - 1$, $u = u - 1$ und fügen A zur Menge der abgedeckten Cluster hinzu. Dann tragen wir erwartet höchstens den folgenden Betrag zu $E[\theta'(C)]$ bei:

$$\begin{aligned} & \frac{\theta(A)}{\theta(C)} \cdot \sum_{a \in A} p_a ((\theta(P_c) + \theta_a + 8\theta_{\text{OPT}}(P_u) - 8\theta_{\text{OPT}}(A))) \\ & \cdot (1 + H_{t-1}) + \frac{u - (t - 1)}{u - 1} \cdot (\theta(P_u) - \theta(A)) \\ & \leq \frac{\theta(A)}{\theta(C)} \cdot \left((\theta(P_c) + 8\theta_{\text{OPT}}(P_u)) \cdot (1 + H_{t-1}) \right. \\ & \quad \left. + \frac{u - t}{u - 1} \cdot (\theta(P_u) - \theta(A)) \right) \end{aligned}$$

Die zweite Abschätzung gilt, da $\sum_{a \in A} p_a \leq 8\theta_{\text{OPT}}(A)$ direkt aus Lemma 3.2.3 folgt. An dieser Stelle müssten wir über alle freien Cluster summieren, um den gesamten Potenzial-Beitrag nach oben abzuschätzen. Dazu setzen wir die Power-Mean-Ungleichung ein. Mit dieser gilt: $\sum_{A \subset P_u} \theta(A^2) \geq \frac{1}{u} \cdot \theta(P_u)^2$. Damit ist es uns nun möglich, den über alle freien Cluster aufsummierten Beitrag zum Potenzial nach oben zu beschränken:

$$\begin{aligned} & \frac{\theta(P_u)}{\theta(C)} \cdot (\theta(P_c) + 8\theta_{\text{OPT}}(P_u)) \cdot (1 + H_{t-1}) \\ & + \frac{1}{\theta(C)} \cdot \frac{u - t}{u - 1} \cdot \left(\theta(P_u)^2 - \frac{1}{u} \cdot \theta(P_u)^2 \right) \end{aligned}$$

Dieser Term lässt sich etwas vereinfachen zu:

$$\frac{\theta(P_u)}{\theta(C)} \cdot (\theta(P_c) + 8\theta_{\text{OPT}}(P_u)) \cdot (1 + H_{t-1}) + \frac{u-t}{u} \cdot \theta(P_u)$$

Wir kombinieren nun den erwarteten Potenzial-Beitrag zu $\theta'(C)$ aus den Fällen (a) und (b) und erhalten die folgende obere Schranke:

$$\begin{aligned} E[\theta'(C)] &\leq (\theta(P_c) + 8\theta_{\text{OPT}}(P_u)) \cdot (1 + H_{t-1}) + \frac{u-t}{u} \cdot \theta(P_u) + \frac{\theta(P_c)}{\theta(C)} \cdot \frac{\theta(P_u)}{u} \\ &\leq (\theta(P_c) + 8\theta_{\text{OPT}}(P_u)) \cdot \left(1 + H_{t-1} + \frac{1}{u}\right) + \frac{u-t}{u} \cdot \theta(P_u) \end{aligned}$$

In der Aussage von Lemma 3.2.4 ist $t \leq u$ vorausgesetzt. Daher gilt $\frac{1}{u} \leq \frac{1}{t}$. Damit folgt der Induktionsschluss. \square

Wir können Satz 3.2.1 als Spezialfall von Lemma 3.2.4 beweisen. Dazu betrachten wir das Clustering C , welches wir erhalten, nachdem die Zentren-Wahl von Kernel- k -means++ abgeschlossen ist und bevor der Kernel- k -means Algorithmus ausgeführt wird. Sei A das Cluster aus C_{OPT} , aus dem wir das erste Zentrum ausgewählt haben. Wir wenden Lemma 3.2.4 mit $t = u = k - 1$ an und legen dazu A als einziges abgedecktes Cluster fest. Dann gilt:

$$E[\theta_{\text{OPT}}] \leq (\theta(A) + 8\theta_{\text{OPT}} - 8\theta_{\text{OPT}}(A)) \cdot (1 + H_{k-1})$$

Es ist allgemein bekannt, dass für $t \geq 2$ die t -te harmonische Zahl nach oben beschränkt ist durch $H_t < \ln t + 1$. Damit gilt insbesondere $H_{k-1} \leq \ln k + 1$. Somit gilt:

$$E[\theta_{\text{OPT}}] \leq (\theta(A) + 8\theta_{\text{OPT}} - 8\theta_{\text{OPT}}(A)) \cdot (\ln k + 2)$$

Satz 3.2.1 folgt mit Lemma 3.2.2.

4 Experimentelle Evaluation

In diesem Kapitel wollen wir empirisch evaluieren, wie performant unsere Algorithmen sind. In Kapitel 1 haben wir bereits erwähnt, dass die Clusteranalyse zahlreiche Anwendungen hat. Zudem müssen gegenwärtig immer häufiger Clusteringprobleme mit großen bis riesigen Eingaben gelöst werden, sodass effiziente Algorithmen zunehmend praktisch erforderlich werden.

Wir gehen folgendermaßen vor: wir beschreiben zunächst in Abschnitt 4.1 die Umgebung, unter der wir unsere Experimente durchgeführt haben. Damit wollen wir die Reproduzierbarkeit unserer Ergebnisse garantieren und eine grobe Vorstellung von der Leistungsfähigkeit der eingesetzten Hardware ermöglichen. Wir stellen in diesem Abschnitt auch die eingesetzten Eingabedaten vor.

In Abschnitt 4.2 untersuchen wir einleitend zunächst, inwieweit der Einsatz der Kernel Methode tatsächlich die Qualität von Clusterings verbessern kann. Außerdem untersuchen wir insbesondere, ob sich der Overhead, der durch die Verwendung eines Kernels entsteht, in der Praxis tatsächlich so gering ist, wie wir zuvor bei der Einführung des Kernel-Tricks argumentiert haben.

In Abschnitt 4.3 diskutieren wir, ob wir mit dem Kernel- k -means++-Algorithmus den gewünschten Trade-off im Bereich der spektralen Clusteranalyse von Graphen erzielen konnten. Dazu betrachten wir einerseits Experimente mit den klassischen Graphclustering-Zielfunktionen und führen zusätzlich die vielleicht wichtigste Anwendung in der Bildsegmentierung durch.

4.1 Experimentelle Rahmenbedingungen

Um Vergleichbarkeit der Testläufe sicherzustellen, wurden sämtliche hier vorgestellten Ausführungen der Algorithmen auf derselben Hardware durchgeführt. Der Rechner wurde unter 64 Bit Linux (Arch Linux) betrieben und verfügte über einen Intel i7 3,5 GHz Vierkern-Prozessor sowie 16 GB Arbeitsspeicher. Die Laufzeiten der Algorithmen wurden jeweils über die reine Ausführung des jeweiligen Algorithmus ermittelt, das heißt Vorverarbeitung wie beispielsweise das Einlesen von Daten aus einer Datei wird in den angegebenen Laufzeiten nicht berücksichtigt.

Die Algorithmen wurden mit einer Ausnahme in C++ implementiert. Als Compiler wurde die GCC in der Version 4.9.2 verwendet. Der Code wurde jeweils für 64-Bit Architekturen übersetzt, das heißt es wurde das GCC-Flag `-m64` gesetzt. Außerdem wurden für die Optimierung die Flags `-O3` und `-fPIC` gesetzt. Lediglich für die Bildsegmentierung wurde MATLAB verwendet.

Für die Graphclustering-Experimente in Abschnitt 4.3 wurden Graphen aus dem „Graph Partitioning Archive“ [SWC04] eingesetzt. Die Daten stammen aus Anwendungsfällen, in denen Graphpartitionierungen berechnet werden müssen und wurden von diversen Forschern im Laufe der Zeit bereitgestellt. Das Archiv ist unter der folgenden URL abrufbar: <http://staffweb.cms.gre.ac.uk/~wc06/partition/> (Stand: 02.02.2015).

Für die Experimente in den Abschnitten 4.2 und 4.4 wurden Clusteringeingaben aus dem „UC Irvine Machine Learning Repository“ [Lic13] verwendet. Das UCI Machine Learning Repository ist unter der folgenden URL erreichbar: <http://archive.ics.uci.edu/ml/> (Stand: 02.02.2015). Es enthält über 300 Datensätze aus Anwendungen der Community für Maschinelles Lernen, von denen einige Clusteringprobleme sind. Das jeweilige zu lösende Problem kann in der durchsuchbaren Übersicht unter <http://archive.ics.uci.edu/ml/datasets.html> dem Attribut „Default Task“ entnommen werden. Die meisten Datensätze sind Klassifizierungsprobleme. Wir haben uns auf Datensätze beschränkt, die als Clustering-Probleme ausgezeichnet sind, da nicht bei allen Datensätzen klar wird, aus welchem Kontext, beziehungsweise welcher Anwendung die Daten stammen.

4.2 Die Kernel Methode

Wir wollen einleitend zunächst untersuchen, ob die Kernel Methode tatsächlich signifikante Verbesserungen der Clusteringqualität erbringen kann. Wir betrachten zunächst die Güte der berechneten Clusterings der Algorithmen k -means (Lloyds Algorithmus) und Kernel- k -means++. Wir bereiten dazu Ergebnisse des Papiers [KLLL05] auf, welches sich mit exakt dieser Fragestellung beschäftigt. Die Datensätze, mit denen diese Experimente nachgestellt wurden, stammen aus diversen Papieren, deren Ergebnisse in [BPKK99] zusammengefasst sind. Tabelle 2 zeigt, dass es sich um kleine Datensätze handelt. Für diese Datensätze sind jedoch die k -means-Zielfunktionswerte der jeweils optimalen Lösungen bekannt [BPKK99]. Diese konnten vermutlich gerade aus dem Grund, dass die Datensätze sehr klein sind, mit überschaubarem Rechenaufwand ermittelt werden. Zudem decken die Datensätze ein großes Spektrum an geometrischen Clustern ab, insbesondere sind hyperspherische und hyperellipsoide Cluster enthalten. Graphische Plots der Cluster sind in [KLLL05] enthalten.

Um den klassischen k -means-Algorithmus mit dem Kernel- k -means-Algorithmus zu vergleichen, haben wir die beiden Implementierungen der Algorithmen in der Dlib-Bibliothek¹ auf die Datensätze angewandt. Für den Kernel- k -means-Algorithmus

¹Version 18.12, <http://dlib.net/> (Stand: 01.02.2015)

Datensatz	Anzahl Punkte	Clusteranzahl k
BENSAID	49	3
DUNN	90	2
IRIS	150	3
ECOLI	336	7
CIRCLE	108	2
BLE-3	320	3
BLE-2	312	2
UE-4	262	4
UE-3	377	3
ULE-4	298	4

Tabelle 2: Eigenschaften der Datensätze aus [BPKK99].

wurde wie in [KLLL05] der Gauss-/RBF-Kernel verwendet. Wir haben diesen für unsere Experimente mit $\sigma = 0.1$ parametrisiert. In [KLLL05] finden sich keine Angaben über die Wahl des Parameters, die geringfügigen Abweichungen unserer Ergebnisse stammen daher vermutlich von der Wahl von σ . Zudem wurde die maximale Anzahl an Iterationen auf 100 beschränkt.

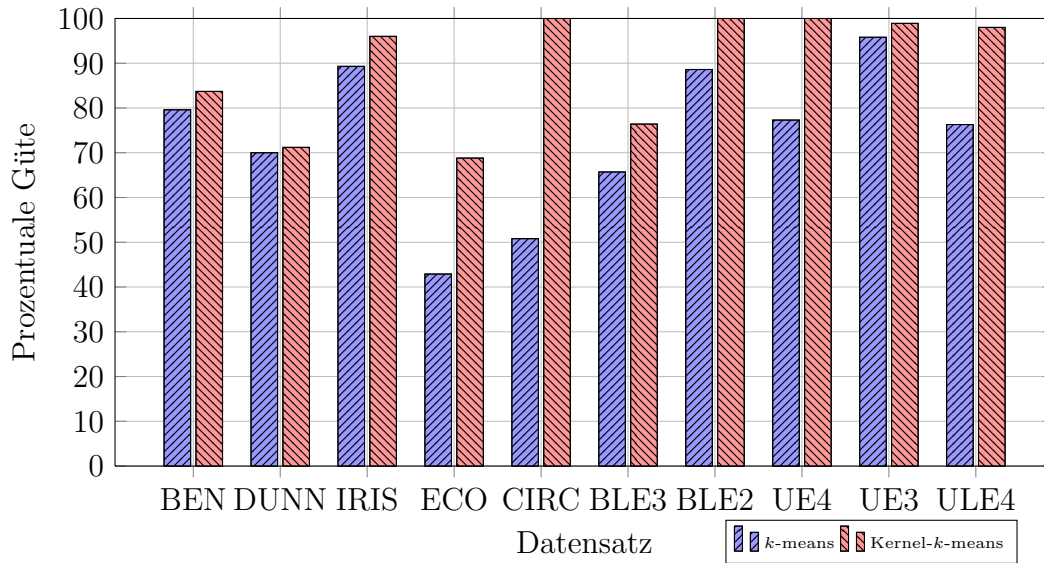


Abbildung 1: Die Clusteringqualität von k -means und Kernel- k -means.

Abbildung 1 und Tabelle 3 zeigen, dass die Kernel Methode bei geometrischem Clustering auf allen Instanzen zu einem besseren Clustering führt. Auf den schwer zu clusternden Instanzen BENSAID und DUNN beträgt die Verbesserung schon einige Prozentpunkte, während das Potenzial der Kernel Methode auf den Instanzen ECOLI und CIRCLE besonders gut zu erkennen ist. Dabei handelt es sich um

Datensatz	k-means	Kernel-k-means
BENSAID	79.6	83.7
DUNN	70	71.2
IRIS	89.3	96
ECOLI	42.9	68.8
CIRCLE	50.8	100
BLE-3	65.7	76.4
BLE-2	88.6	100
UE-4	77.3	100
UE-3	95.8	98.9
ULE-4	76.3	98

Tabelle 3: Die Clusteringqualität von k -means und Kernel- k -means.

Instanzen, auf denen die Cluster hyperellipsoide, beziehungsweise konzentrische Form haben. Auf diesen Instanzen clustert der k -means-Algorithmus wegen der Hyperebenen-Beschränkung nur sehr schlecht, während mit der Kernel Methode auch nicht linear separierte Cluster möglich sind. Bemerkenswert ist auch, dass der Kernel- k -means-Algorithmus sogar auf Instanzen mit mehreren hundert Punkten noch die optimale Lösung berechnet, wie man an den Instanzen CIRCLE, BLE-2 und UE-4 erkennen kann.

Die Datensätze sind nicht geeignet, um Laufzeiteinbußen bei der Kernel Methode zu ermitteln, da beide Algorithmen auf Instanzen dieser Größe nahezu unmittelbar terminieren. Wir verwenden daher größere Datensätze aus dem UCI Machine Learning Repository. Tabelle 4 gibt einen Überblick über die Beschaffenheit der Datensätze.

Datensatz	Anzahl Punkte	Dimension d
NYSK	10241	7
CENSUS	2458285	68
BAG-OF-WORDS	8000000	100000
PLANTS	22632	70
130-US	100000	55
GFE	27965	100
YOUTUBE	120000	1000000

Tabelle 4: Eigenschaften der Datensätze aus [Lic13].

In Abbildung 2 und Tabelle 5 ist gut zu erkennen, dass der Overhead, der durch die Kernel Methode entsteht bei den kleinen Instanzen (bis zu 100.000 Punkten)

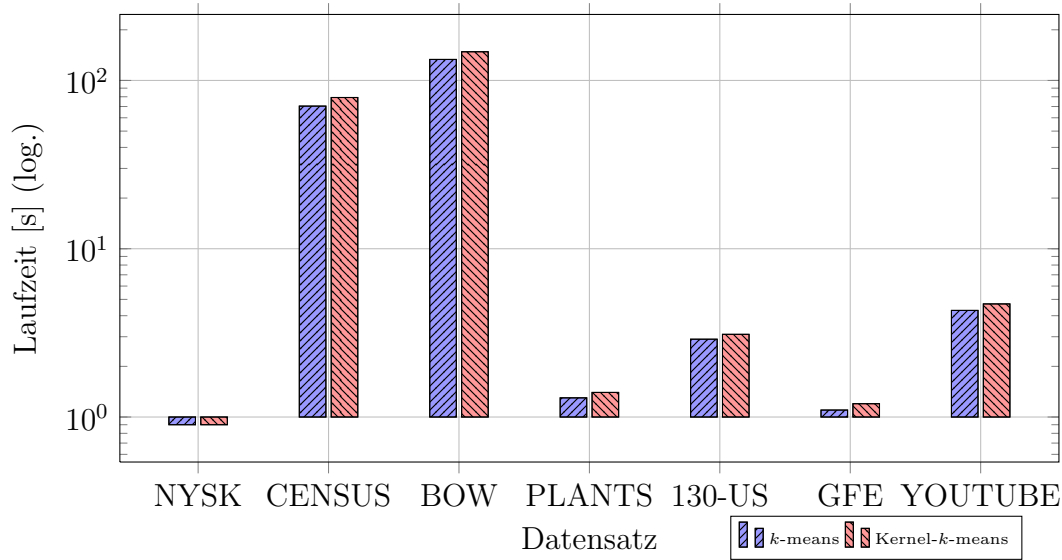


Abbildung 2: Die Laufzeiten von k -means und Kernel- k -means.

vernachlässigbar gering ist. Auf den größeren Instanzen sind Laufzeiteinbußen in der Größenordnung von rund 10% zu verzeichnen. Diese sind nicht zu vernachlässigen. In Relation zur potenziellen Qualitätssteigerung aus den vorhergegangenen Experimenten betrachtet, handelt es sich jedoch um einen Trade-off, der sich in den allermeisten Fällen lohnt.

Datensatz	k -means	Kernel- k -means
NYSK	0.9	0.9
CENSUS	70.4	79.2
BOW	133.3	148.1
PLANTS	1.3	1.4
130-US	2.9	3.1
GFE	1.1	1.2
YOUTUBE	4.3	4.7

Tabelle 5: Die Laufzeiten von k -means und Kernel- k -means in Sekunden.

Wir diskutieren im Folgenden die experimentellen Ergebnisse im Zusammenhang mit dem Kernel- k -means++-Algorithmus.

4.3 Der Kernel- k -means++-Algorithmus

Für den Kernel- k -means++-Algorithmus haben wir das Framework „Graclus“² von Kulis und Guan in der Version 1.2 erweitert. Dabei handelt es sich um eine

²<http://www.cs.utexas.edu/users/dml/Software/graculus.html> (Stand: 27.01.2015)

Zielfunktion	Knotengewichte	Kernelmatrix
Ratio Association	1 für alle Knoten	$K = \sigma I + A$
Ratio Cut	1 für alle Knoten	$K = \sigma I - L$
Kernighan-Lin	1 für alle Knoten	$K = \sigma I - K$
Normalized Cut	Grad des Knoten	$K = \sigma D^{-1} + D^{-1} A D^{-1}$

Tabelle 6: Die σ -verschobenen Kernelmatrizen für die Graphschnitt-Zielfunktionen [DGK04].

Implementierung des gewichteten Kernel- k -means Algorithmus 7 für die spektrale Clusteranalyse beziehungsweise Graphpartitionierung. Eine detaillierte Beschreibung, die recht nah an der Implementierung orientiert ist, haben die Autoren in [DGK07] gegeben. Graculus ist im Bereich quelloffener Graphpartitionierungs-Software gegenwärtig das vielleicht schnellste Framework, wie ein Vergleich im „Berkeley Segmentation Dataset and Benchmark“³ vermuten lässt. Graculus wird in einer Reihe von quelloffener Bildverarbeitungs-Software verwendet. Das prominenteste Beispiel, das Graculus verwendet, ist vermutlich „Clustering Views for Multi-view Stereo (CMVS)“⁴.

Bevor wir unsere Experimente vorstellen, wollen wir kurz auf ein bislang offen gebliebenes Problem aus Abschnitt 3.1 eingehen. Wir hatten am Ende des Abschnitts angemerkt, dass wir für eine konkrete Implementierung des Algorithmus sicherstellen müssen, dass die verwendeten Kernelmatrizen positiv definit sind. In [DGK07] wird detailliert diskutiert, dass für alle Graphpartitionierungs-Zielfunktionen gezeigt werden kann, dass eine einfache diagonale Verschiebung der ursprünglichen Kernelmatrix K um einen Faktor σ die Analogie über die Formulierungen als Spurmaximierungsprobleme intakt bleibt. Die Wahl von σ muss dabei groß genug sein, damit der entsprechend verschobene Kernel K' eine positiv definite Matrix ist. Für manche Eingaben kann es vorkommen, dass die Gewichtsmatrix W so beschaffen ist, dass die Kernelmatrix auch ohne eine diagonale Verschiebung schon positiv definit ist. Die Experimente in [DGK04, DGK07] legen nahe, dass es in solchen Fällen sogar vorteilhaft sein kann, eine negative σ -Verschiebung vorzunehmen. Empirisch wird σ idealerweise so gewählt, dass die Spur der Kernelmatrix möglichst nach bei Null liegt. Tabelle 6 zeigt die σ -verschobenen Kernelmatrizen für die diversen Graphschnitt-Zielfunktionen. Wie bereits erwähnt sind für unsere Zwecke nur die Ratio Association und der Normalized Cut relevant.

³<http://www.eecs.berkeley.edu/Research/Projects/CS/vision/grouping/segbench/> (Stand: 29.01.2015)

⁴<http://www.di.ens.fr/cmvs/> (Stand: 29.01.2015)

Datensatz	Anzahl Knoten	Anzahl Kanten
BCSSTK31	35588	572914
T60K	60005	89440
598A	110971	741934
144	144649	1074393
M14B	214765	1679018
AUTO	448695	3314611

Tabelle 7: Eigenschaften der Datensätze aus dem Graph Partitioning Archive [SWC04].

Graphpartitionierung. Graculus selbst bietet drei Möglichkeiten der Initialisierung an: zufällige Initialisierung (analog zum klassischen Algorithmus von Lloyd), spektrale Initialisierung (entspricht der Initialisierung über Eigenwert-Berechnung, die wir zu Beginn von Abschnitt 3.2 diskutiert haben) und Initialisierung durch METIS. METIS [KK98] ist ein state-of-the-art Graphpartitionierungs-Framework, das in einem mehrschrittigen Verfahren Cluster gleicher Größe in Graphen berechnet. Es diente bereits in [DGK04, DGK07] als Benchmark-Vergleichsinstanz. Wir werden METIS auch in unseren Experimenten für Vergleiche heranziehen.

Wegen der Zufallskomponente im Kernel- k -means++-Algorithmus haben wir diesen sowie die vollkommen zufällige Initialisierung pro Datensatz und pro Clusterzahl k insgesamt 100 mal ausgeführt. In unseren Ergebnissen geben wir für die zufällige Initialisierung jeweils das arithmetische Mittel über alle 100 Läufe an. Für den Kernel- k -means++-Algorithmus geben wir zusätzlich die gemessenen Minima und Maxima über alle Durchläufe an, um eine Vorstellung der Größenordnung der Schwankungen zu ermöglichen. Wir betrachten insbesondere auch die Laufzeiten der Algorithmen, um den Trade-off zwischen Clusteringqualität und Ausführungszeit beurteilen zu können. Wir untersuchen die erzielten Werte für Ratio Association (je *größer* der Wert, desto besser ist die Qualität des berechneten Clusterings) und Normalized Cut (je *kleiner* der Wert, desto besser die Qualität des berechneten Clusterings). Für jeden Datensatz und jede Zielfunktion berechnen wir mit den Algorithmen jeweils Clusterings mit 128 und 512 Clustern. Damit erreichen wir eine höhere Anforderung an die spektralen Algorithmen. Wir haben bereits thematisiert, dass bei der Performanz für diese die Anzahl der Cluster kritisch ist, da sie die Anzahl der zu berechnenden Eigenwerte erhöht. Schließlich sei bereits vorab erwähnt, dass die von uns gemessenen Werte bezüglich der Laufzeit teilweise stark von denen aus [DGK07] (nach unten) abweichen. Dies ist mit großer Wahrscheinlichkeit auf die deutlich leistungsfähigere Hardware zurückzuführen:

in [DGK07] wurde ein Rechner mit einem 2.4 GHz Pentium IV-Prozessor und nur 1 GB Arbeitsspeicher eingesetzt. Es ist zudem naheliegend, dass der Prozessor über einen 32 Bit Chipsatz verfügte und Graculus dementsprechend auch für 32 Bit-Architekturen kompiliert wurde.

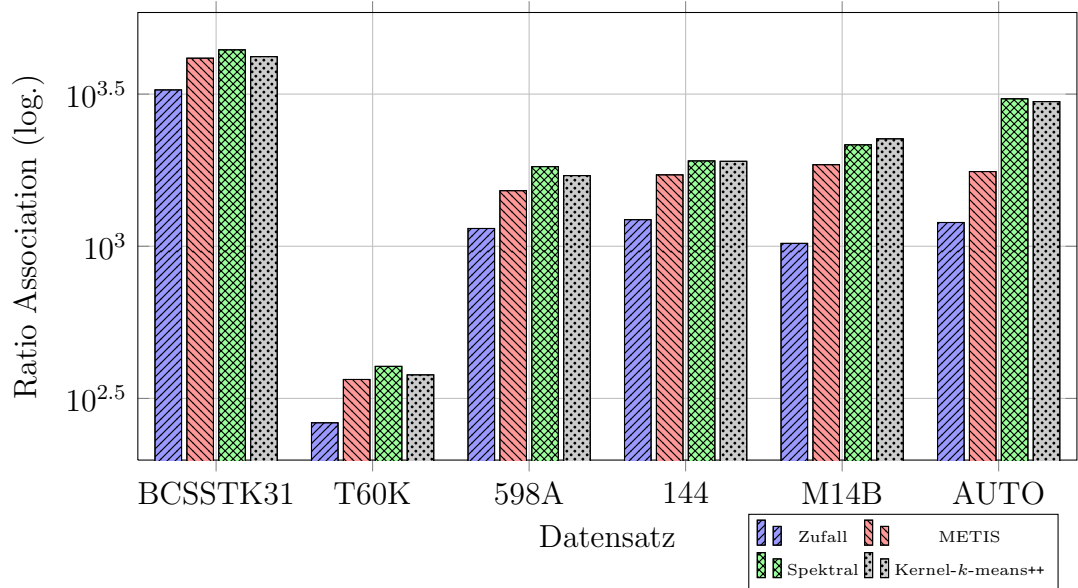


Abbildung 3: Ratio Association Werte der Algorithmen für $k = 128$ (je größer, desto besser).

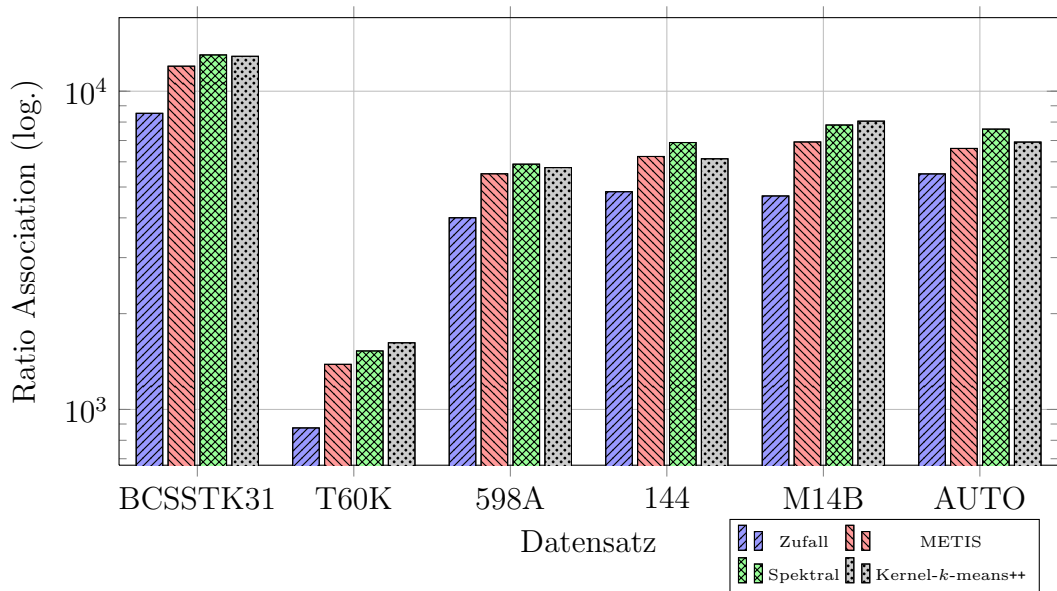


Abbildung 4: Ratio Association Werte der Algorithmen für $k = 512$ (je größer, desto besser).

Datensatz	Zufall	METIS	Spektral	Avg	Min	Max
BCSSTK31	3265	4150	4422	4199	3955	4322
T60K	263	365	403	378	345	385
598A	1144	1523	1827	1707	1622	1754
144	1223	1717	1908	1902	1521	2334
M14B	1022	1854	2155	2255	2133	2552
AUTO	1197	1760	3052	2988	2852	3128

Tabelle 8: Ratio Association Werte der Algorithmen für $k = 128$ (je größer, desto besser).

Datensatz	Zufall	METIS	Spektral	Avg	Min	Max
BCSSTK31	8522	11984	13002	12873	3955	4322
T60K	875	1386	1527	1620	345	385
598A	4003	5502	5899	5755	1622	1754
144	4833	6236	6899	6133	1521	2334
M14B	4687	6927	7832	8052	2133	2552
AUTO	5498	6610	7603	6922	2852	3128

Tabelle 9: Ratio Association Werte der Algorithmen für $k = 512$ (je größer, desto besser).

4.4 Kernmengenkonstruktion

5 Zusammenfassung

Literatur

- [ADHP09] ALOISE, Daniel ; DESHPANDE, Amit ; HANSEN, Pierre ; POPAT, Preyas: NP-hardness of Euclidean sum-of-squares clustering. In: *Machine Learning* (2009), S. 245–248
- [AV07] ARTHUR, David ; VASSILVITSKII, Sergei: k-means++: The advantages of careful seeding. In: *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2007, New Orleans, Louisiana, USA, January 7-9, 2007*, 2007, S. 1027–1035
- [BGV92] BOSER, Bernhard E. ; GUYON, Isabelle ; VAPNIK, Vladimir: A Training Algorithm for Optimal Margin Classifiers. In: *Proceedings of the Fifth Annual ACM Conference on Computational Learning Theory, COLT 1992, Pittsburgh, PA, USA, July 27-29, 1992.*, 1992, S. 144–152
- [BPKK99] BEZDEK, James C. ; PAL, Mikhil R. ; KELLER, James ; KRISNAPURAM, Raghu: *Fuzzy Models and Algorithms for Pattern Recognition and Image Processing*. Norwell, MA, USA : Kluwer Academic Publishers, 1999
- [CKV13] CELEBI, M. E. ; KINGRAVI, Hassan A. ; VELA, Patricio A.: A comparative study of efficient initialization methods for the k-means clustering algorithm. In: *Expert Syst. Appl.* 40 (2013), Nr. 1, S. 200–210
- [Cov65] COVER, T.M.: Geometrical and Statistical Properties of Systems of Linear Inequalities with Applications in Pattern Recognition. In: *Electronic Computers, IEEE Transactions on* EC-14 (1965), June, Nr. 3, S. 326–334
- [DGK04] DHILLON, Inderjit ; GUAN, Yuqiang ; KULIS, Brian: A Unified View of Kernel k-means, Spectral Clustering and Graph Cuts / University of Texas at Austin. 2004 (TR-04-25). – Forschungsbericht
- [DGK07] DHILLON, Inderjit S. ; GUAN, Yuqiang ; KULIS, Brian: Weighted Graph Cuts without Eigenvectors A Multilevel Approach. In: *IEEE Trans. Pattern Anal. Mach. Intell.* 29 (2007), Nr. 11, S. 1944–1957
- [FSS13] FELDMAN, Dan ; SCHMIDT, Melanie ; SOHLER, Christian: Turning big data into tiny data: Constant-size coresets for k -means, PCA and projective clustering. In: *Proceedings of the Twenty-Fourth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2013, New Orleans, Louisiana, USA, January 6-8, 2013*, 2013, S. 1434–1453

- [GL96] GOLUB, Gene H. ; LOAN, Charles F.: *Matrix computations*. 3. Auflage. Johns Hopkins University Press, 1996
- [HM04] HAR-PELED, Sarel ; MAZUMDAR, Soham: On coresets for k-means and k-median clustering. In: *Proceedings of the 36th Annual ACM Symposium on Theory of Computing*, 2004, S. 291–300
- [HSS08] HOFMANN, Thomas ; SCHÖLKOPF, Bernhard ; SMOLA, Alexander J.: Kernel methods in machine learning. In: *Ann. Statist.* 36 (2008), 06, Nr. 3, S. 1171–1220
- [KK98] KARYPIS, George ; KUMAR, Vipin: A Fast and High Quality Multilevel Scheme for Partitioning Irregular Graphs. In: *SIAM J. Sci. Comput.* 20 (1998), Nr. 1, S. 359–392. – ISSN 1064–8275
- [KL70] KERNIGHAN, B. W. ; LIN, S.: An Efficient Heuristic Procedure for Partitioning Graphs. In: *The Bell System Technical Journal* 49 (1970), Nr. 1, S. 291–307
- [KLLL05] KIM, Dae-Won ; LEE, Ki Y. ; LEE, Doheon ; LEE, Kwang H.: Evaluation of the performance of clustering algorithms in kernel-induced feature space. In: *Pattern Recognition* 38 (2005), Nr. 4, S. 607–611
- [Lan50] LANCZOS, Cornelius: An iteration method for the solution of the eigenvalue problem of linear differential and integral operators. In: *Journal of Research of the National Bureau of Standards* 45 (1950), S. 255–282
- [Lic13] LICHMAN, M.: *UCI Machine Learning Repository*. <http://archive.ics.uci.edu/ml>. Version: 2013
- [Llo82] LLOYD, Stuart P.: Least squares quantization in PCM. In: *IEEE Transactions on Information Theory* 28 (1982), Nr. 2, S. 129–136
- [Lux07] LUXBURG, Ulrike von: A Tutorial on Spectral Clustering. In: *Statistics and Computing* 17 (2007), Nr. 4, S. 395–416
- [MS84] MEGIDDO, Nimrod ; SUPOWIT, Kenneth J.: On the Complexity of Some Common Geometric Location Problems. In: *SIAM J. Comput.* 13 (1984), Nr. 1, S. 182–196
- [NJW01] NG, Andrew Y. ; JORDAN, Michael I. ; WEISS, Yair: On Spectral Clustering: Analysis and an algorithm. In: *Advances In Neural Information Processing Systems*, 2001, S. 849–856

- [ORSS12] OSTROVSKY, Rafail ; RABANI, Yuval ; SCHULMAN, Leonard J. ; SWAMY, Chaitanya: The effectiveness of Lloyd-type methods for the k-means problem. In: *J. ACM* 59 (2012), Nr. 6, S. 28
- [SC04] SHAWE-TAYLOR, John ; CRISTIANINI, Nello: *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004
- [Sch14] SCHMIDT, Melanie: *Coresets and Streaming Algorithms for the k-means Problem and Related Clustering Objectives*, Technische Universität Dortmund, Diss., 2014
- [SM00] SHI, Jianbo ; MALIK, Jitendra: Normalized Cuts and Image Segmentation. In: *IEEE Trans. Pattern Anal. Mach. Intell.* 22 (2000), Nr. 8, S. 888–905
- [SWC04] SOPER, Alan J. ; WALSHAW, Chris ; CROSS, Mark: A Combined Evolutionary Search and Multilevel Optimisation Approach to Graph-Partitioning. In: *J. Global Optimization* 29 (2004), Nr. 2, S. 225–241

Hiermit versichere ich, dass ich die vorliegende Arbeit selbstständig verfasst habe und keine anderen als die angegebenen Quellen und Hilfsmittel verwendet sowie Zitate kenntlich gemacht habe.

Dortmund, den 14. Februar 2015

Lukas Pradel

