

October 22, 2015
DRAFT

An Experimental Study into Spectral and Geometric Approaches to Data Clustering

Prashant Sridhar

October 2015

School of Computer Science
Carnegie Mellon University
Pittsburgh, PA 15289

Thesis Committee:

Dr. Gary Miller, Chair

Dr. Alex Smola

*Submitted in partial fulfillment of the requirements
for the degree of Masters in Computer Science.*

October 22, 2015
DRAFT

Keywords: Geometry, Nonparametric, Clustering

October 22, 2015
DRAFT

For my parents

October 22, 2015
DRAFT

October 22, 2015
DRAFT

Abstract

October 22, 2015
DRAFT

Acknowledgments

My advisor is cool.

October 22, 2015
DRAFT

Contents

1	Introduction	1
2	Roadmap	5
3	Traditional Non-parametric clustering	7
4	Density based distance metrics and clustering	9
4.1	Approximation of density based distance metrics	9
4.2	Gabriel graphs and approximate Gabriel graphs	11
4.3	Weakly Gabriel graph and fast linear spanner	17
4.4	Results and Discussion	21
5	Geometric non-parametric clustering	25
	Bibliography	27

October 22, 2015
DRAFT

List of Figures

4.1	Quadratic size Gabriel graph example	13
4.2	Gabriel edges with small angle	14
4.3	Candidate hyperplane elimination	19

October 22, 2015
DRAFT

List of Tables

October 22, 2015
DRAFT

Chapter 1

Introduction

Clustering of data points is a problem that has existed in machine learning literature for quite some time now. The general definition is to partition the data set into "clusters" based on their similarity. These techniques have several applications to fields like computational biology, computer vision, robotics, and others. The defining characteristic of clustering that also makes it so powerful is that it is traditionally defined in the unsupervised sense. This means that data is presented to the algorithms without any prior labels designating clusterings. This unsupervised nature is part of the appeal since unlabeled data is often easy to collect and quite plentiful, while labeling data is often laborious, requiring a human to hand label the data set. However, while unsupervised learning is often easy to collect data for, it is notoriously difficult to judge when the algorithm has produced a good clustering. Parametric clustering methods often have a well defined statistical interpretation, but perform poorly on data sets that do not satisfy the model assumptions. Non-parametric methods, on the other hand, are either very slow to implement in practice or not very well defined statistically. This thesis provides a statistical basis for a non-parametric clustering method that is efficient in practice.

One standard approach to nonparametric clustering is mean shift clustering [1]. The intuition behind this approach is that the underlying probability distribution has several maxima. Each point in the space can be attributed to one of these maxima or modes by following its gradient ascent curve to a maxima. In order to compute the ascent curve, most approaches use multiple start gradient ascent. The function maximized is the estimated probability density function. This estimate is obtained using some kernel, often a gaussian. While the idea is quite simple, it runs into several issues that prevent it from being scaled up to large data sets. Firstly, the algorithm requires a quadratic number of operations in general. Let each point require m rounds to converge and the size of the data set be N . Mean-shift clustering requires $O(mN^2)$ steps to converge. Contrast this with k-means [11] which only requires $O(mN)$ steps to converge. This quadratic run time often proves infeasible for large data sets. Further, kernel density estimation itself carries some run time concerns. If infinite support kernels like the gaussian kernel are used, then estimating the density at a given point is itself a linear time operation. This causes the final run time of the mean-shift algorithm to become cubic. While finite support kernels can be used, the question of kernel bandwidth still remains. The kernel bandwidth is a hyperparameter to this model and can often only be chosen by cross-validation and grid search. Finally, there are also some concerns of the quality of the clusters generated. Mean-shift requires high data density in

order to have a consistent kernel density estimate. In particular, the approach often fails to cluster outlier points that lie between two natural clusters.

Spectral methods are another common non-parametric way to cluster data sets [15]. The data points are connected together to form a graph with heavy weight edges between similar points and light edges between dissimilar points. This graph can then be cut in order to generate the appropriate clustering. Strong edges will rarely be cut because they would greatly decrease the cut quality. Conversely, weak edges are encouraged to be cut in order to keep the cut value low. This leaves similar nodes in the same cluster while separating dissimilar nodes. As a non-parametric method it is resistant to non-linearities in the data set, however, it is often quite slow to run in practice. If a complete graph of the data is constructed, storing the graph itself will require $O(N^2)$ space. While quadratic time algorithms are slow to run in practice on larger data sets, quadratic space algorithms are completely infeasible even on smaller samples. Storing the graph itself will require several gigabytes of memory in data sets larger than a few ten thousands of points preventing most computers from even loading the graph into RAM. Further, if the graph is to be partitioned into k clusters, algorithm will require finding the k largest eigenvalues of graph laplacian. If the graph is very near complete, and consequently the laplacian is very dense, this operation will be very expensive making it infeasible to run on anything but the smallest data sets. Finally, deciding exactly how to weight the edges of the graph is still a matter of trial and error. The most common method is to use a gaussian kernel to decide the edge weights, but there is little statistical basis to this. It is merely a way to handle non-linearity of the input data. Choosing the bandwidth of this gaussian kernel, again, must be done with grid search and cross validation.

The method of clustering proposed in this thesis solves all the problems mentioned above. Our method also relies on partitioning a graph, however, we provide some statistical explanation for the edge weights. Additionally, our method constructs a graph with only a linear number of edges in $O(n \log n)$ time. This linear size graph is feasible to store and results in a sparse laplacian, which can also be solved easily. We view clustering of data points as a partitioning of the probability space. Our interpretation of a good partition is one separating two high density regions while having a low cross-sectional cut volume. This encourages splitting up maxima across lines drawn through density troughs. This idea is also very similar to sparsest cuts on graphs and lends itself very easily to graphical interpretations. The goal is to construct a sparse graph wherein cut quality is similar to that of the distribution generating the data. The probability distribution can be viewed as a heat distribution over a metal plate. Such heat distributions can be modelled by triangular meshes over the data with added Steiner points. Since the mesh is a triangulation, it will only have a linear number of edges. The triangular mesh can then be cut with standard spectral methods in order to retrieve the clustering. More specifically, the mesh approximates second nearest neighbour distances for the distribution, but these distances can easily be used as an estimate of the density function. Since we are only concerned with clusters and not the point-wise estimation of density, we do not face the traditional problems faced by other knn approaches of having non-convergent variance [8]. Integrated over finite measure sets, variance of the estimator drops off with the size of the sample. Further, clusters of at least $O(\log n)$ elements should be identified by our approach since by Chernoff bounds, enough points will be sampled from these clusters to determine the cut. This gives provides a guarantee that small clusters will be identified by our approach. Finally, as a spectral, non-parametric method,

this algorithm is insulated against non-linearities in the data allowing it to identify clusters with non elliptical shapes. Thus, the approach overcomes the shortcomings of previous clustering algorithms.

This thesis serves mostly as an experimental exploration into the idea of sparse graphs for data clustering since much of the theory is still being developed. Primarily, this paper focuses on data with "rare" cluster populations and highly non-elliptical clusters that are difficult to cluster with standard methods like gaussian mixture modeling [6] or k-means. Specifically, we will deal with flow cytometry data. Flow cytometry is a laser based approach to differentiating between populations of cells based on certain bio markers on the cells surface or inside its body. Cytometers can take measurements from several thousands of cells per second leading to large datasets. Flow cytometry is often used in the diagnosis of several health disorders including various types of cancers. It is also used in clinical trials including research to autonomously compute prognoses for HIV patients. Many of these applications require clustering cell types as a primitive before prediction can take place, and at present, such clustering is often done by hand. Several model based clustering algorithms have been developed for this problem such as, FLAME, flowClust, flowMerge. These methods are all mixture models over t-distributions or skew t-distributions. Skew t-distributions, in theory, should be able to account for eccentricities in the data, but the time complexity scales to the fourth power with dimension rendering it useless in practice beyond 5 dimensions. Further, these model based methods often fail to find rare cell populations which are of particular interest to biological applications. Non-parametric methods have also been tried in this space and mostly rely on spectral clustering. SAMSppectral solves the issue of quadratic space by heuristically sampling data, which the authors call faithful sampling, to form a representative population. As with all heuristics, there is no theoretical basis to this approach and it provides no provable guarantees. Our geometric spectral clustering method outperforms SAMSppectral while being able to produce guarantess on quality and run time.

October 22, 2015
DRAFT

Chapter 2

Roadmap

October 22, 2015
DRAFT

Chapter 3

Traditional Non-parametric clustering

More formally, consider a sample space over d dimensions, \mathbb{R}^d with an associated probability density function f . For each point, we can define a unique integral path or gradient ascent path as follows,

$$\pi'_x(t) = \nabla f(\pi_x(t))$$

October 22, 2015
DRAFT

Chapter 4

Density based distance metrics and clustering

Our initial approach to combating the non-linearity of cytometry data, was to try non-euclidean distance metrics. Under a density based metric, two points are closer together if joined by a path through high density space, and farther apart if only connected through low density space. This should allow cluster detection even for clusters with highly non-elliptical shape. However, density based metrics are often intractable to compute exactly leading most such approaches to use reliable approximations. A recent approximation method developed by Choen et. al., [2], provides a graph based method to approximate a density based metric by computing shortest paths on a Gabriel graph over the points. The final graph has a linear(in the points) number of edges leading to quick single source shortest path queries. However, naive implementations of k-means replacing Euclidean distance with density based distance will still run into the problem of having to store all pairs shortest paths. This will still require $O(n^2)$ space rendering the approach infeasible. While, straightforward applications of the distance metric might not be feasible, spectrally cutting the Gabriel graph itself intuitively leads to a good clustering. This section of the thesis does not contain fully reasoned theory. Instead, it is an intuitive exploration of the usefulness of the Gabriel graph. While it is hard to arrive at provable bounds for cluster quality, Gabriel graphs are not as prone to the curse of dimension. While the number of edges grows exponentially with dimension, it is still upper bounded by $O(n^2)$, or the complete graph. On the other hand, the main method presented in this paper is based on Delauny triangulations and require time exponential in the dimension to construct. Therefore, the Gabriel graph may provide a clue on how to deal with high dimensional data sets. Before explaining how to cluster using the Gabriel graph, this thesis will first go over the work of Cohen et. al. to provide some background for the discussion and define relevant terms.

4.1 Approximation of density based distance metrics

If \mathbb{R}^d space is Endowed with a distance metric, we can define the length of an arbitrary path γ through this space as follows. Let the length function of a path under Euclidean distance be ℓ_e , and let γ begin at the point x at time 0 and end at the point y at time 1.

$$\begin{aligned} \ell_e(\gamma) &= \int_{\gamma} ds \\ &= \int_0^1 \left| \frac{d\gamma(t)}{dt} \right| dt \end{aligned}$$

Which is simply the path integral of γ in \mathbb{R}^d . Here, $\left| \frac{d\gamma(t)}{dt} \right|$ is the velocity of trajectory at time t . Note that this integral is minimized by using the straight line path between the two points. Using this formulation for the length of a path, we can define the Euclidean distance, \mathbf{d}_e between the points x and y to be,

$$\mathbf{d}_e(x, y) = \inf_{\gamma} \ell_e(\gamma)$$

This is merely a mathematical precise way of saying that the Euclidean distance between two points is the length of straight line shortest path between the points.

With Euclidean distance formally defined, this can now be extended to more general distance metrics. By adding a cost function based on location in the space, segments of a path through different regions are priced differently thereby contorting the shortest path away from the straight line path between the points. Let the cost function for a certain metric k be $c(x)$. We can now define the length of a path γ for this metric, ℓ_k , as follows.

$$\begin{aligned} \ell_k(\gamma) &= \int_{\gamma} c(s) ds \\ &= \int_0^1 c(s) \left| \frac{d\gamma(t)}{dt} \right| dt \end{aligned}$$

There is also an analogous definition for distance, \mathbf{d}_k , for this metric.

$$\mathbf{d}_k(x, y) = \inf_{\gamma} \ell_k(\gamma)$$

Density based distance metrics penalize paths through low density regions and reward paths through high density regions. Therefore, when used as a similarity metric, it will cluster together points belonging to the same high density regions. Let P be the probability distribution over \mathbb{R}^d , and $f(x)$ be its density function at point x . Define $c_p(x) = f(x)^{\frac{2(1-p)}{d}}$ to be the cost function used for the density based distance metric. Also define $\ell_p(\gamma)$ and $(d)_p(x, y)$ to be the associated path length function and distance metric. Note that $p = 1$ returns the Euclidean metric.

The problem with using density based metrics is that for $p \neq 1$, computing $\ell_p(\gamma)$ is near intractable for many distributions and paths. Often the true distribution is not even known and must therefore be estimated. Worse, computing $\mathbf{d}_p(x, y)$ involves finding an infimum over an infinite number of paths. Hence, the distance can never be exactly computed and must be approximated.

Assume that n points are uniformly randomly sampled from P to form the training data set D . We can now define an undirected graph G over these points such the edge length, e_{xy} ; $x, y \in D$, is equal to the p^{th} power of the Euclidean distance between the points, $e_{xy} = (d)_e(x, y)^p$. Define $\mathbf{d}_G(x, y)$ to be the shortest path distance from x to y in the complete graph G . The following theorem from Hwang et. al., [9] proves a strong connection between $\mathbf{d}_G(x, y)$ and $\mathbf{d}_p(x, y)$ at the price of a few assumptions on the sample space and the density function. Assume that P is instead defined over a manifold M , the same manifold that D is drawn from.

Theorem 4.1.1. *Assume M is compact, and that f is continuous and supported everywhere over M . There exists a constant $C(d, p) > 0$, which only depends on d and p , satisfying the following. Let $\epsilon > 0$ and $b > 0$, then there exists $\theta > 0$ such that*

$$P \left(\sup_{x, y} \left| \frac{\mathbf{d}_G(x, y)}{n^{\frac{(1-p)}{d}} \mathbf{d}_p(x, y)} - C(d, p) \right| > \epsilon \right) < e^{-\theta n^{\frac{1}{(d+2p)}}}$$

for all sufficiently large n where $x, y \in M$ and $\mathbf{d}_e(x, y) \geq b$.

In the case where M is closed instead of compact, the paper presents a second theorem giving an almost surely limit relating the two metrics.

Theorem 4.1.2. *Assume M is closed and that f is continuous and supported everywhere over M . There exists a constant $C(d, p) > 0$, which only depends on d , and p , satisfying the following.*

$$\lim_{n \rightarrow \infty} n^{\frac{p-1}{d}} \mathbf{d}_G(x, y) = C(d, p) \mathbf{d}_p(x, y)$$

for fixed $x, y \in M$.

The assumptions of note here are that M is either compact or closed. This means that the theorems do not work in the general case where data can be drawn from anywhere in \mathbb{R}^d . In practice, however, this assumption is not a problem because the set of valid data to be clustered can often be bounded between some finite intervals. The density function is also assumed to be continuous and supported everywhere over its domain. This can be achieved by simply shrinking its domain to only contain space in which f is supported. Since points cannot be drawn from 0 probability regions, this assumption is easily met in practice. What these assumptions buy though, are strong bounds relating \mathbf{d}_p and \mathbf{d}_g thereby giving a way to approximate density based distance metrics. However, while computing shortest paths on G is tractable, computing G itself requires quadratic space rendering this approach useless on larger data sets. This problem can be avoided by considering the special case of $p = 2$, as is done by Cohen et. al.

4.2 Gabriel graphs and approximate Gabriel graphs

Once we restrict $p = 2$, we can immediately make a few observations about the graph G . Since edge weights are now squared Euclidean distance, the shortest path between two points will not a long edge if another point exists inside of the circle inscribed by the edge as a diameter. Let the graph constructed for $p = 2$ be G_2 .

Definition 4.2.1. *An edge e_{xy} between points x, y is said to be Gabriel if there exists no z such that z is contained within the circle inscribed on e_{xy} as diameter.*

Lemma 4.2.1. *For all $x, y \in G_2$, the shortest path from x to y will only contain edges that are Gabriel.*

Proof. Assume for the sake of contradiction that the shortest path from x to y contains an edge e_{uv} that is not Gabriel. By 4.2.1, there exists a point z inside the circle formed using e_{uv} as diameter. Consider the triangle formed by u, v, z . By the properties of triangles, the angle formed at z will be at least 90° . Therefore, $e_{uv} \geq e_{uz} + e_{vz}$ by the Pythagorean theorem. Therefore we have found a shorter path and have arrived at a contradiction. \square

Definition 4.2.2. *The graph GG_2 , with $V = D$ and $E = \cup_{e_{xy}, e_{yx}}$ is said to be the Gabriel graph over the points in D .*

Note that the set of Gabriel edges is a subset of Delaunay edges since the Gabriel requirement is a stricter form of the Delaunay requirement over edges.

The Gabriel graph serves as a possibly sparse graph that can be used to compute the density metric. Because the Gabriel graph is a subset of the Delaunay triangulation, it serves as a linear spanner for $d = 2$. This is because the Delaunay triangulation in the 2D case forms a planar graph which necessarily has a linear number of edges. However, even in 3 dimensions, examples can be constructed where the Gabriel graph has $O(n^2)$ edges. Consider the unit sphere in 3 dimensions and two small arcs on along this sphere on the XY plane and the YZ plane with both arcs centered around $Y = 0$. Arrange $n/2$ points equally spaced on each arc. By construction, every edge between a pair of points on different arcs will be Gabriel since each of those edges is approximately a diametric chord for the sphere. Therefore, this example leads to at least $n^2/4$ edges. This example can be seen in 4.1.

This prompts the need for spanners on the Gabriel graph. It is worth noting that even in higher dimensional cases where the Gabriel graph is linear, there is no fast algorithm to compute it. The naive approach of checking if each edge is Gabriel requires $O(n^3)$ time even when only a linear number of edges are Gabriel.

The simplest solution is to use a Euclidean spanner over the points. It can be shown via induction that the Euclidean spanner is also a spanner for the Gabriel graph

Theorem 4.2.2. TODO: *Prove that euclidean spanner is spanner for gabriel graph.*

While very efficient algorithms [4], [3], working in $O(n \log n)$, time exist for computing Euclidean spanners, these algorithms are still exponential in dimension because they rely on space partitioning trees. In practice, these algorithms do not work well for $d > 10$. For this reason, other approaches that do not scale so poorly with dimension are desirable.

While the original paper did not present a way to compute the Gabriel graph quickly, they do present an algorithm to produce a linear, $(1 + \epsilon)$ spanner given the Gabriel graph. The simple construction involves constructing a conflict graph for each node p based on the angles formed by its neighbours. For any node p , construct a graph, H_p by connecting p to all its neighbours, and connect any pair of neighbours a, b if the angle formed by e_{pa} and e_{pb} is sufficiently small, say, smaller than θ . Define I_p to be the set of edges in the maximal independent set H_p . The linear spanner of the Gabriel graph is the graph H_θ such that

$$E_{H_\theta} = \{(p, q) | e_{pq} \in (I_p \cup I_q)\}$$

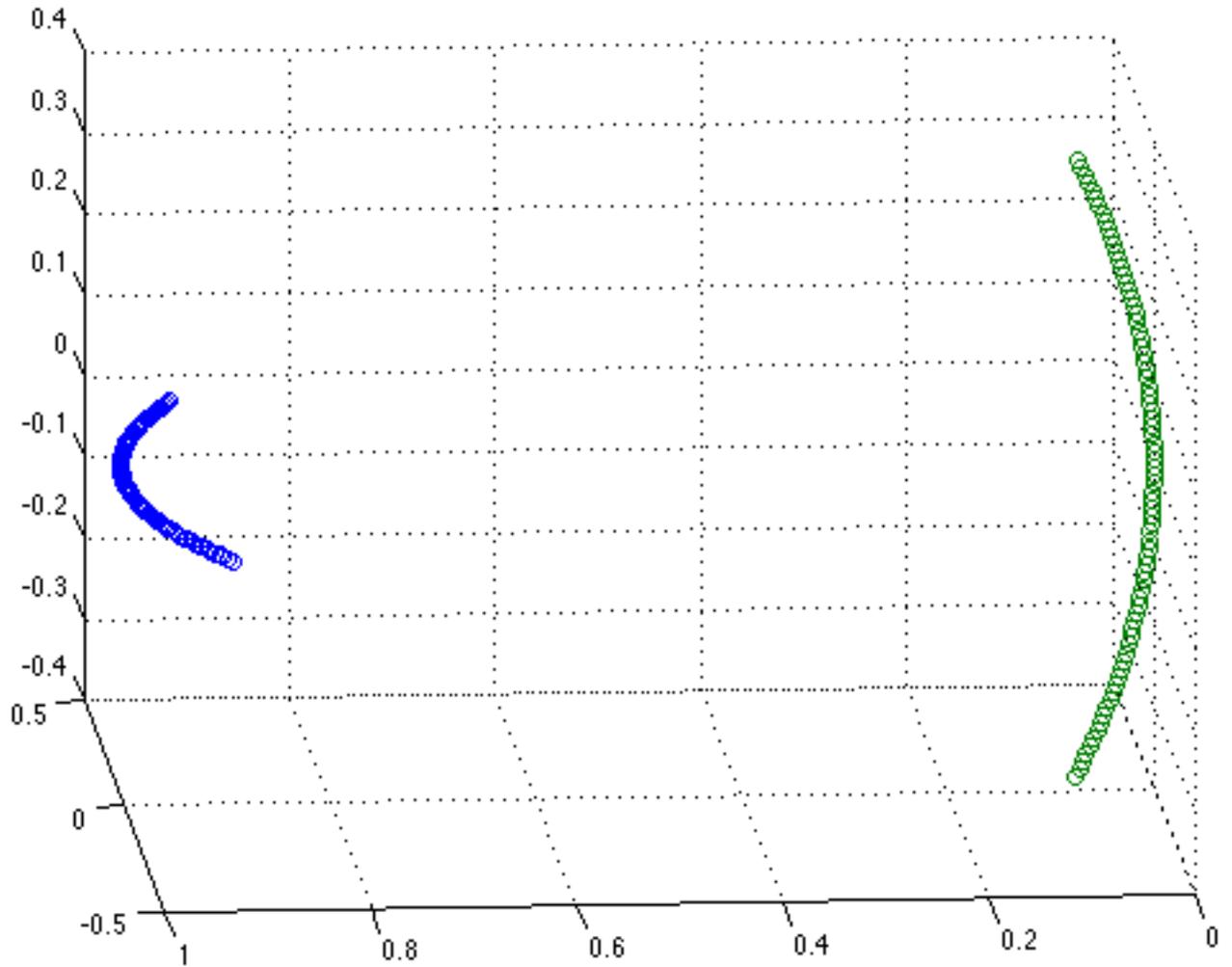


Figure 4.1: Quadratic size Gabriel graph example

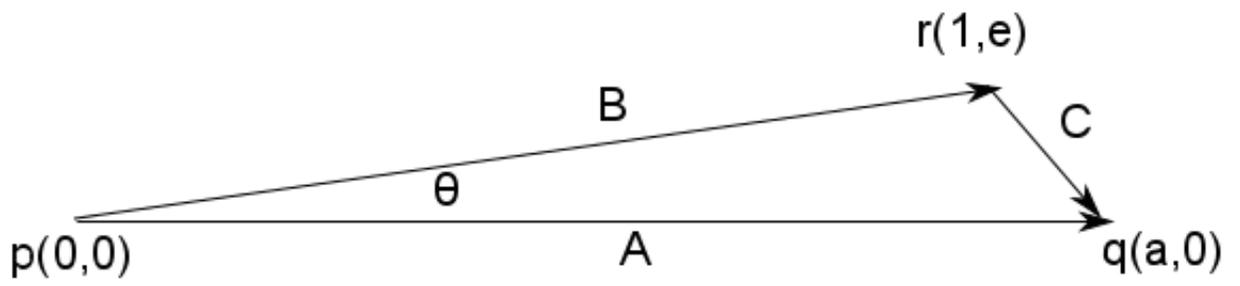


Figure 4.2: Gabriel edges with small angle

Before presenting the proof of quality for the spanner H_θ , we must first prove two supporting lemmas. We show that if two edges make a sufficiently small angle, then the length of the path does not change by much depending on which we pick.

Theorem 4.2.3. *Consider a point p with neighbours q, r such that e_{pq}, e_{pr} are both Gabriel and make an angle of at most θ . Let A be the Euclidean length of e_{pq} and B be the Euclidean length of e_{pr} . Then it is the case that for any path through the graph containing e_{pq} , its length would not change by more than $(1 + 2\tan^2\theta)|A|^2$ if e_{pr} were traversed instead.*

Proof. Begin by noticing that we can assume e_{qr} exists in the graph. If e_{qr} is not in the graph, then it is not Gabriel, implying that there exists a point s within the circle with qr as diameter. Then, traversing e_{qs} followed by e_{sr} would be shorter than directly traversing e_{qr} . Therefore, e_{qr} 's existence serves as an upper bound assumption for the stretch generated by not traversing e_{pq} and does not weaken the proof.

Let C be the Euclidean length of e_{qr} . We want to show that $|B|^2 + |C|^2 \leq (1 + 2\tan^2\theta)|A|^2$.

Without loss of generality, assume that p is at the origin, q is at $(a, 0)$ and r is at $(1, e)$. This is shown in 4.2

$$\begin{aligned} \frac{|C|^2 + |B|^2}{|A|^2} &= \frac{(1-a)^2 + e^2 + 1 + e^2}{a^2} \\ &= \frac{2 - 2a + a^2 + 2e^2}{a^2} \end{aligned}$$

This ratio is at its worst when $a = 1$, or when A and C are perpendicular. To see this, differentiate the above ratio with respect to a . The derivative is negative for $1 \leq a \leq 1 + e^2$. For a outside of this range, either q or r will encroach on the ball of the other. Therefore, if the angle between A and B is less than θ ,

$$\begin{aligned} \frac{|B|^2 + |C|^2}{|A|^2} &\leq \frac{|B|^2 + |B|^2 \sin^2\theta}{|B|^2 \cos^2\theta} \\ &= \frac{1 + \sin^2\theta}{\cos^2\theta} \\ &= 1 + 2\tan^2\theta \\ \implies |B|^2 + |C|^2 &\leq (1 + 2\tan^2\theta)|A|^2 \end{aligned}$$

□

This theorem justifies why a maximal independent set of the conflict graph of a point is sufficient to preserve path lengths. If an edge is not present in the maximal independent set, then it must have been the case that it conflicted with some other edge at a low degree. Therefore, the other edge can be traversed instead with small cost to total length. One other lemma is needed for the the proof of the spanner. We need to show that the $|C|$ is not too large as compared to A .

Theorem 4.2.4. *Consider a point p with neighbours q, r such that e_{pq}, e_{pr} are both Gabriel and make an angle of at most θ . Let A be the Euclidean length of e_{pq} , B be the Euclidean length of e_{pr} and C be the Euclidean length of e_{qr} . Then $|C|^2 \leq \tan^2\theta|A|^2$.*

Proof. Once again, without loss of generality, assume p is on the origin, q is at $(a, 0)$ and r is at $(1, e)$.

$$\begin{aligned} \frac{|C|^2}{|A|^2} &= \frac{(1-a)^2 + e^2}{a^2} \\ &= \frac{1 - 2a + a^2 + e^2}{a^2} \end{aligned}$$

This ratio is again maximized if e_{qr} and e_{pq} are perpendicular. Therefore, if the angle between the edges is θ , we have

$$\begin{aligned} \frac{|C|^2}{|A|^2} &\leq \frac{|B|^2 \sin^2\theta}{|B|^2 \cos^2\theta} \\ &= \tan^2\theta \\ \implies |C|^2 &\leq \tan^2\theta|A|^2 \end{aligned}$$

□

With this theorem in place, we are now ready to prove correctness for the spanner. The proof will be by induction.

Theorem 4.2.5. *For any pair of points p, q in GG_2 , the path from p to q through H_θ will larger by at most a fraction of $(1 + \epsilon)$ for any θ such that $\tan^2\theta \leq \frac{\epsilon}{2+\epsilon}$. Further, H_θ will have a linear number of edges in fixed dimension.*

Proof. This theorem will be proved using induction over the length of the longest edge in a path. For the base case, consider (p, q) the closest pair in the data set. Since they are the closest pair, p is q 's nearest neighbour. Nearest neighbour edges are always Gabriel because if they were not, a point would lie within the circumscribing circle thereby generating a nearer neighbour. e_{pq} is also in H_θ because if it were not, then it must conflict with another edge, say, e_{pr} . By 4.2.4, e_{qr} must be shorter still meaning p, q is not the closest pair.

For the inductive step, assume that for all paths between pairs (p, q) with longest edges shorter than $|A|$, H_θ is a $(1 + \epsilon)$ spanner. Now consider all path with a longest edge of length $|A|$. If this edge is in H_θ , then the above theorem holds. If the edge is not in H_θ , then it must conflict with some neighbour. Let the endpoints of the longest edge be p, q , and the end points of the edge it conflicts with be p, r . Let B be the length of e_{pr} be $|B|$ and the length of e_{qr} be $|C|$. Let P be the path from q to r in H_θ . Since $|C|^2 \leq \tan^2\theta|A|^2$, the shortest path from q to r in GG_2 cannot have any edges as long as $|A|$. Therefore by the inductive hypothesis,

$$|P| \leq (1 + \epsilon)|C|^2$$

Again, note that if e_{qr} is not in GG_2 , the length of the path would be even shorter. Hence the assumption that e_{qr} is in GG_2 is an overestimate.

We now want to bound $|B|^2 + |P|$ since that is the path from p to q in H_θ .

$$\begin{aligned} |B|^2 + |P| &\leq |B|^2 + |C|^2 + \epsilon|C|^2 \\ &\leq (1 + 2\tan^2\theta)|A|^2 + \epsilon\tan^2\theta|A|^2 \\ &\leq (1 + (2 + \epsilon)\tan^2\theta)|A|^2 \\ &\leq (1 + \epsilon)|A|^2 \end{aligned}$$

Finally, H_θ must have a linear number of edges in a fixed dimension because in any fixed dimension, each node may only have a constant out degree when constrained to only have neighbours separated by a minimum degree. \square

4.3 Weakly Gabriel graph and fast linear spanner

The last section went over how to construct Gabriel graphs to compute density based distance metrics and how to sparsify these graphs. This section weakens the conditions on Gabriel graphs to allow for faster algorithms to construct them. Further, we show that we can sparsify these more general Gabriel graphs to still retain linear sized graphs. To begin with, we define the notion of an edge being weakly Gabriel. The edge to a neighbour is said to be weakly Gabriel if there are no other neighbours within the dihedral ball of the edge. However, other, non-neighbours are still allowed to encroach the ball. Formally, we define weakly Gabriel as follows,

Definition 4.3.1. Consider a point p , with its neighbour set V_p . Let q be a point in the neighbour set, $q \in V_p$. Define the edge e_{pq} to be weakly Gabriel if, $\forall r \in V_p$, r is not in the dihedral ball of e_{pq} .

Definition 4.3.2. Consider a point p . The neighbour set V_p of p is said to be weakly Gabriel if $\forall q \in D$, e_{pq} is Gabriel, $q \in V_p$, and $\forall q \in V_p$, e_{pq} is weakly Gabriel.

Definition 4.3.3. For a data set D , the graph WG_2 is a weakly Gabriel graph over D if, $\forall p \in D$, V_p is weakly Gabriel.

To begin with, notice that each data set has a unique Gabriel graph, but need not have a unique weakly Gabriel graph. Secondly, notice that a weakly Gabriel graph still contains all Gabriel edges and therefore preserves all shortest paths. Hence, the shortest path between any two nodes, and, by extension, the approximate density based distance, is not changed if a weakly Gabriel graph is used instead of a Gabriel graph. The trade-off is that a weakly Gabriel graph will always be at least as dense as a Gabriel graph and could potentially have a quadratic number of edges even when the Gabriel graph is linear. However, the main gain of using the weakly Gabriel graph is that its definition lends itself to a fast iterative, input sensitive algorithm where Gabriel graphs must be computed in $O(n^3)$ time.

The first observation needed for this algorithm is that for a given point, the edge to its nearest neighbour must necessarily be Gabriel. Necessarily, there may be no points in the shared edges

dihedral circle. This serves as a base case for the algorithm. Once we have identified a Gabriel edge, several candidate neighbours can be immediately eliminated as seen in the 4.3.

The points above the hyperplane can all be eliminated because the nearest neighbour would encroach on their ball. The hyperplane check can be performed simply by checking if the angle θ in the figure is obtuse. This is equivalent to checking if the cosine is negative, or simply checking if the dot product is negative. If only the Gabriel edges were desired, every single point would have to draw its hyperplane and eliminate potentially offending points. However, this will again require $O(n^3)$ time. Instead, we need only eliminate points with successive nearest neighbours in order to obtain a weakly Gabriel neighbour set. This is formalized in 1

Algorithm 1: Compute Weakly Gabriel Graph

```

Data: Data set  $D$ 
Result: Weakly Gabriel Graph,  $WGG = \{V, \{V_p\}\}$  over  $D$ 
 $WGG \leftarrow \{\}$ 
for  $p \in D$  do
   $V_p \leftarrow \{\}$ 
for  $p \in D$  do
   $d_p \leftarrow []$ 
  for  $q \in D - \{p\}$  do
     $append(d_p, (q, \mathbf{d}_e^2(p, q)))$ 
   $s_p \leftarrow sort(d_p)$  ascending by distance
  while  $s_p$  is not empty do
     $r \leftarrow pop(s_p)$  for  $s \in s_p$  do
      if  $cos(rp, rs) > 0$  then
         $remove(s_p, s)$ 
     $add(V_p, r)$ 
   $add(WGG, \{p, V_p\})$ 
  for  $q \in V_p$  do
    if  $p \notin V_q$  then
       $add(V_q, p)$ 

```

The run time of this algorithm is $O(n^2 \log n)$ for the sorting step, and $O(n^2 e)$ for the elimination steps, where e is the maximum out degree of a point. In practice, this algorithm works well at building a sparse graph on which to run shortest path queries because many real world data sets have linear size Gabriel graphs. For those applications, The weakened Gabriel graph algorithm produces only a small fraction of non-Gabriel edges. For example, in one flow cytometry data-set consisting of roughly 70,000 points in 16 dimensions, only 2% of the edges generated were non-Gabriel. This highlights the algorithms input sensitivity. In most practical use cases, the bottle-neck in runtime is the sorting step giving an overall complexity of $O(n^2 \log n)$ since e is often a constant. However, it is possible to construct edge cases with a linear number of Gabriel edges for which 1 produces a quadratic number of edges leading to a $O(n^3)$ time and $O(n^2)$ space complexity.

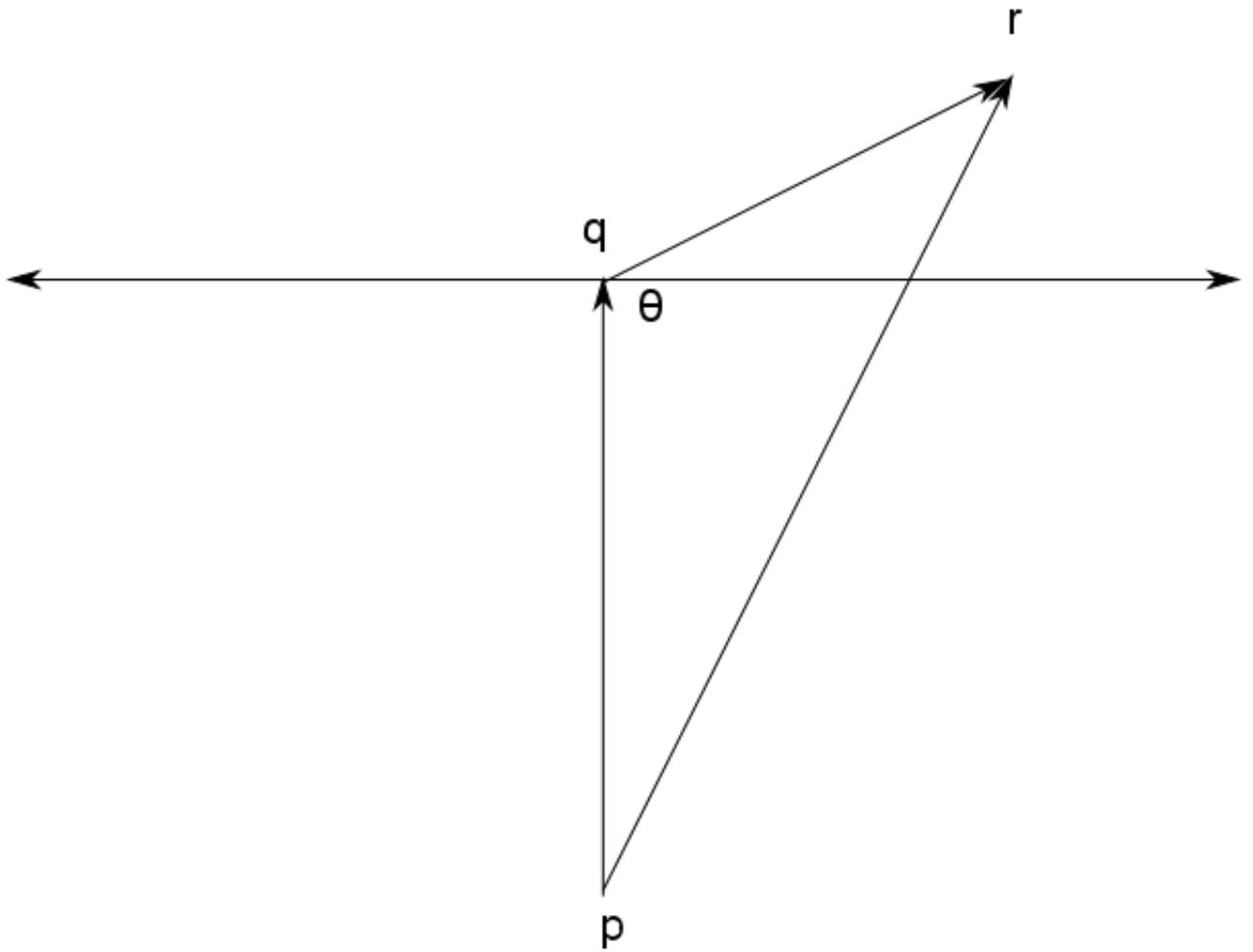


Figure 4.3: Candidate hyperplane elimination

The good news is that sparsifying based on angle will continue to work on weakly Gabriel graphs. If a non-Gabriel edge is ever picked in favour of a Gabriel edge, the stretch will still remain small because the edges must be comparable in length. If this were not the case, one would encroach on the dihedral ball of the other violating the requirement that both edges be at least weakly Gabriel. This preserves lemmas, 4.2.3 and 4.2.4. With these lemmas in place, the main proof carries through very similarly.

TODO: Decide if I should add the proof

Instead of first running 1 and then sparsifying the resultant graph, the two steps can be combined into one. When a point is added to the neighbour set and used to prune potential neighbours, the only check performed is the hyperplane check. Instead, the neighbour can prune points by its hyperplane and also by the angle made with p . Finally, because the resultant graph will have constant(in dimension) out degree from each node by virtue of being a linear spanner, the sorting step can be removed in favour of picking the closest neighbour each time. Since we will only have to pick a linear number of neighbours, this reduces the runtime from $O(n^2 \log n + n^2 e)$ to $O(n^2 e)$, which is the same as $O(n^2)$. This algorithm is formally represented in 2

Algorithm 2: Compute Sparse Weakly Gabriel Graph

Data: Data set D
Data: Error threshold ϵ
Result: Weakly Gabriel Graph, $WGG = \{V, \{V_p\}\}$ over D
 $WGG \leftarrow \{\}$
for $p \in D$ **do**
 $V_p \leftarrow \{\}$
for $p \in D$ **do**
 $d_p \leftarrow []$
 for $q \in D - \{p\}$ **do**
 $append(d_p, (q, \mathbf{d}_e^2(p, q)))$
 while d_p is not empty **do**
 $r \leftarrow findmin(d_p)$
 $remove(d_p, r)$
 for $s \in d_p$ **do**
 if $\cos(rp, rs) > 0 \parallel \tan^2(pr, ps) > \frac{\epsilon}{2+\epsilon}$ **then**
 $remove(s_p, s)$
 $add(V_p, r)$
 $add(WGG, \{p, V_p\})$
 for $q \in V_p$ **do**
 if $p \notin V_q$ **then**
 $add(V_q, p)$

4.4 Results and Discussion

While asymptotically we can achieve $O(n^2)$, in practice it is often faster to perform the sorting step and run at $O(n^2 \log n)$. It can be shown that for a point placed on the origin and surrounded by neighbours drawn uniformly from the unit ball, the expected number of Gabriel edges is $O(2^d)$. It can also be empirically observed that the average angle between edges is 90° . Therefore, each node often has a large out degree making repeatedly picking the nearest neighbour very slow if 2^d exceeds $\log n$. This also demonstrates that for data with no irregularities, sparsification is not necessary and constructing a weakly Gabriel graph is sufficient. Also, the number of non-Gabriel edges added by 2 is a very small percentage of the total edges. For the applications discussed in this paper, no more than 5% of the edges generated are non-Gabriel. Therefore, 2 is a significant improvement in computing approximations to density based distance.

Clustering with non-Euclidean distance metrics poses a challenge because computing the mean of a cluster is difficult thereby rendering k-means impossible to use. The solution is to assume that the cluster center will be a point in the data set. This gives rise to the k-medoids algorithm [10]. The most common realization of the idea is Partitioning Around Medoids (PAM) [16] which proceeds as follows.

- Select k random points without repetition to serve as cluster medoids.
- Assign each point to a cluster based on mediod distance.
- Compute total cost of clustering and repeat until cost converges
 - For each mediod m and non-mediod o
 - Swap m and o and recompute cost.
 - Keep swap if cost decreases.

Other mediod based algorithms also exist such as the following based on Voronoi iterations [12].

- Select k random points without repetition to serve as cluster medoids.
- Assign each point to a cluster based on mediod distance.
- Compute total cost of clustering and repeat until cost converges
 - For each cluster, select m as mediod if m minimizes the sum of distance to all other points.
 - Reassign points to the closest mediod.

If the number of iterations required to converge is assumed to be t and that pairwise distances are stored in a Gram Matrix G , then the runtime of PAM is $O(tkn^2)$. Each internal step require $O(n^2)$ time to swap a mediod with each non-mediod point and recompute the cost. The second method also has a similar runtime of $O(tkn^2)$ since it take $O(n^2)$ time to find the optimal mediod for each cluster. The advantage of the second approach is it more closely mirrors the more familiar k-means algorithm and is this easier to reason about. The main drawback to both of these approaches is that when the distance metric is density based distance, computing the distance between two points requires a shortest path computation using Dijkstra's algorithm, which will cost $O(en \log n)$ where e is the maximum out degree. This makes computing distance

on the fly infeasible. The other alternative is to precompute all pairs shortest paths and store the result in a Gram matrix. However, computing and storing the Gram matrix will have a $O(n^2)$ space complexity. While this approach is not further explored in this thesis, a solution to this problem might arise from the graph structure used to compute density based metrics. To begin, let us define the notion of closeness centrality [14] in a graph. The closeness centrality $C(x)$ of a graph is defined as follows, $C(x) = \frac{1}{\sum_y d_e(y,x)^2}$. The reciprocal of the closeness centrality is the distance from a node to all the other nodes in the graph. This is very similar to the k-medioids requirement of computing the distance from the mediod to every other node in the cluster. Several algorithms such as [17][13][5] deal with computing closeness centrality or a related metric called betweenness centrality [7] which is often considered harder to compute than closeness centrality. These algorithms rely on sampling to build approximations to closeness and the methods therein could significantly speed up mediod computation. Sampling techniques could also be useful in cluster assignment since the subsampled paths or nodes could build approximations to the distance to different medioids for each point.

The other alternative is to cluster directly based on the properties of the Gabriel graph. Instead of setting the edge weight e_{xy} to be the squared Euclidean distance, $e_{xy} = \mathbf{d}_e(x, y)^2$, set the edge weight to be the reciprocal, $e_{xy} = \frac{1}{\mathbf{d}_e(x, y)^2}$. Geometrically close points in this graph are usually connected by a single Gabriel edge or a series of short hops. In either of those cases, short edges will have very high weights leading the cut between the points to be very expensive. On the flip side, long edges on the graph are not very common since they are very likely to have some other point encroach on them. They mostly arise from the boundries of clusters. These long edges are also much cheaper to cut than short edges, thereby incentivising cuts between clusters. In non-sparsified Gabriel graphs, there is some chance that there many connecting edges between two clusters if the clusters are sufficiently compact and far away. This is similar to the scenario where $O(n^2)$ Gabriel edges exist in the graph. These dense inter cluster connections are problematic for clustering because they drive up the cost of cutting the two clusters with several nearly identical edges. Sparsified Gabriel garphs fix this issue by removing such "duplicate" edges. If several edges arise from and end with similar nodes, the edges will make small angles with each other and will hence be pruned. This makes sparsified Gabriel graphs very good to cluster over. Note that instead of the reciprocal of the squared Euclidean distance, a gaussian kernel could have been used for edge length. While results vary, they do not change significantly by altering the kernel. The weights that are mentioned in this section are used because they tie in to ideas of effective resistance and conductance of the edges. The experimental results section will go into more detail but for now consider the following figure of a clustering constructed by partitioning Gabriel graphs. **TODO:ADD GABRIEL GRAPH CLUSTERING PICTURE ON BARCODE AND TCELL.** The data set used here is called the Barcode data **TODO:CITE BARCODE** set and it captures the difficulty of clustering rare populations. The Barcode data set also provides an easily verified correct answer that is not available for other flow cytometry applications(where experts often also disagree on clusterings). If we contrast the quality of this clustering with Gaussian mixture models **TODO: ADD CLUSTERING OF BARCODE WITH GMMS**, we see that cutting the Gabriel graph performs much better. The primary win of this approach is that the Gabriel graph is still fast to construct in high dimensions. While the runtime will depend on the dimensionality of the data, the number of edges in the graph is

bounded by $O(n^2)$ and since 2 is an input sensitive algorithm, the runtime of that algorithm is also bounded. On the other hand, the approach mentioned later in this thesis requires computing the Delaunay triangulation, a procedure that scales exponentially with dimension. The primary downside of Gabriel graphs is that clustering on them is not based soundly in theory while this is not true for the approach to be discussed in the next chapter. However, our algorithm grounded in theory still calls upon the Delaunay triangulation, and since Gabriel edges are a subset of the Delaunay edges, the Gabriel graph might serve as an approximation to the more principled clustering approach in high dimension.

October 22, 2015
DRAFT

Chapter 5

Geometric non-parametric clustering

In the previous section we explored a clustering method with promising results that was not fully grounded in theory. This section will present a statistical interpretation to clustering and detail how this can be translated into an algorithm. In the process, we will provide some background for finite element methods.

October 22, 2015
DRAFT

Bibliography

- [1] Yizong Cheng. Mean shift, mode seeking, and clustering. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 17(8):790–799, 1995. 1
- [2] Michael B. Cohen, Brittany Terese Fasy, Gary L. Miller, Amir Nayyeri, Donald Sheehy, and Ameya Velingker. Approximating nearest neighbor distances. *CoRR*, abs/1502.08048, 2015. URL <http://arxiv.org/abs/1502.08048>. 4
- [3] Gautam Das and Giri Narasimhan. A fast algorithm for constructing sparse euclidean spanners. *International Journal of Computational Geometry & Applications*, 7(04):297–315, 1997. 4.2
- [4] Michael Elkin and Shay Solomon. Optimal euclidean spanners: really short, thin and lanky. *CoRR*, abs/1207.1831, 2012. URL <http://arxiv.org/abs/1207.1831>. 4.2
- [5] David Eppstein and Joseph Wang. Fast approximation of centrality. *J. Graph Algorithms Appl.*, 8:39–45, 2004. 4.4
- [6] Mario AT Figueiredo and Anil K Jain. Unsupervised learning of finite mixture models. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 24(3):381–396, 2002. 1
- [7] Linton C Freeman. A set of measures of centrality based on betweenness. *Sociometry*, pages 35–41, 1977. 4.4
- [8] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*, volume 1. Springer series in statistics Springer, Berlin, 2001. 1
- [9] Sung Jin Hwang, Steven B Damelin, and Alfred O Hero III. Shortest path through random points. *arXiv preprint arXiv:1202.0045*, 2012. 4.1
- [10] L. Kaufman and P.J. Rousseeuw. Clustering by means of medoids. in statistical data analysis based on the l_1 -norm and related methods. 1987. 4.4
- [11] Stuart P Lloyd. Least squares quantization in pcm. *Information Theory, IEEE Transactions on*, 28(2):129–137, 1982. 1
- [12] Hae-Sang Park and Chi-Hyuck Jun. A simple and fast algorithm for k-medoids clustering. *Expert Systems with Applications*, 36(2):3336–3341, 2009. 4.4
- [13] Matteo Riondato and Evgenios M Kornaropoulos. Fast approximation of betweenness centrality through sampling. In *Proceedings of the 7th ACM international conference on Web search and data mining*, pages 413–422. ACM, 2014. 4.4
- [14] Gert Sabidussi. The centrality index of a graph. *Psychometrika*, 31(4):581–603, 1966. 4.4

- [15] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 22(8):888–905, 2000. 1
- [16] Sergios Theodoridis and Konstantinos Koutroumbas. *Pattern Recognition, Third Edition*. Academic Press, Inc., Orlando, FL, USA, 2006. ISBN 0123695317. 4.4
- [17] Vladimir Ufimtsev and Sanjukta Bhowmick. An extremely fast algorithm for identifying high closeness centrality vertices in large-scale networks. In *Proceedings of the Fourth Workshop on Irregular Applications: Architectures and Algorithms*, IA3 '14, pages 53–56, Piscataway, NJ, USA, 2014. IEEE Press. ISBN 978-1-4799-7056-8. doi: 10.1109/IA3.2014.12. URL <http://dx.doi.org/10.1109/IA3.2014.12>. 4.4