

Contents

1	Spectral Clustering for Data Streams	1
1.1	Introduction	1
1.2	Literature Review	2
1.3	Spectral Clustering	8
1.4	Micro-cluster based spectral clustering in data streams	14
1.4.1	The challenges of clustering datastreams	14
1.4.2	Clustream	15
1.4.3	Weighting the Micro-Clusters	21
1.5	Experimentation	25
1.5.1	Methodology	26
1.5.2	Simulated Results	28
1.5.3	Real Data Results	30
1.5.4	Optional Extras	35
1.6	Conclusion	37
	Bibliography	38

¹ Chapter 1

² Spectral Clustering for Data Streams

³ 1.1 Introduction

^{starts off abruptly}

- ⁴ Spectral clustering has been empirically shown to perform well under difficult cluster scenarios, but is computationally intensive to run for large data sets. KASP offers an approach to deal with this burden, by summarising the data using k-means, and feeding the centroids as input to spectral clustering. We show how spectral clustering can be applied to online data by combining the Clustream framework with the concepts developed in KASP Yan et al. (2009). Clustream is a framework for dealing with a constant stream of data, and splits the clustering process into two stages: micro-clustering and macro-clustering. We analyse the performance of spectral clustering in an online setting on simulated and real texture-based data sets, both for static and evolving streams. We discuss how to quantify the information loss from summarising the data stream and investigate empirically and algebraically how this perturbation may effect the overall performance. **RD: Rewrite, excluding any references to KASP, focus more on unweighted / weighted.** (I am not sure this is the best strategy given that the unweighted is not working well. Perhaps postpone writing this until you have conclusions from experiments.)
- ¹⁵

- ¹ We discuss the literature in Section 1.2. Spectral clustering is introduced in Section 1.3.
- ² Online clustering and the Clustream algorithm are discussed in Section 1.4 and experimental
- ³ results shown in Section 1.5.

⁴ 1.2 Literature Review

- ⁵ This review introduces the historical origins of Spectral Clustering, discusses methods to
- ⁶ speed up the computation time of offline Spectral Clustering algorithms and comments on
- ⁷ existing incremental spectral clustering methods.

⁸ Spectral Clustering is a technique for separating data into groups or clusters. The elements
general cluster definition

⁹ in each group should be similar to each other, and elements in different groups should be

¹⁰ dissimilar. In Spectral Clustering a weighted distance matrix of the data points called the
give more detail, formula?

¹¹ affinity matrix is formed and by evaluating the spectral properties of this affinity matrix

¹² clusters can be estimated. The idea of splitting a data set into groups or clusters with high

¹³ within-cluster similarity and between-cluster dissimilarity, can also be thought of as a graph

¹⁴ partitioning problem with the data points as graph vertexes and the elements in the affinity

¹⁵ matrix representing weights on the graph edges. Imagine that all of the data points in a data

¹⁶ set are joined together by string. If you were asked break this graph into two distinct pieces,

¹⁷ which strings would you cut?

¹⁸ Graph partitioning problems have long been of interest to mathematicians, ever since

¹⁹ Euler solved the popular “Seven Bridges of Königsberg” problem in 1736. However it was

²⁰ not until the 1970’s that the link between graph partitioning and matrix spectra was ob-

²¹ served when Donath and Hoffman (1973) established a lower bound for the number of edges

1 which must be cut in a graph partitioning problem, based on a factor of the eigenvectors of
 2 the affinity matrix. Around the same time Fieldler (1973) published work highlighting the
 3 importance of the value of the second smallest eigenvector of the affinity matrix when produc-
 4 ing a two way graph partition, also known as a bi-partition. As bi-partitioning was studied
k? as we later use n for # of observations
 5 the question arose, how to step from bi-partitioning problems to n-partitioning problems?
 6 The obvious two approaches are to either recursively apply bi-partitioning on the graph or
 7 calculate an n-partition directly from the n smallest eigenvectors. This is discussed further
 8 in Alpert (1995) where they claim that recursive bi-partitioning, which was more popular at
 9 the time, is inferior compared with creating one partition using many eigenvectors. For our
 10 purposes we will consider only the direct n-partitioning problem. For a full review of the
 11 historical origins of spectral clustering, see the tutorial (von Luxburg et al. (2008)).
not appropriate I think for a PhD chapter

12 The first paper to bring Spectral Clustering to the Machine Learning community was Shi
 13 and Malik (2000) which treated image segmentation as a graph-cut problem. Mathematically,
 14 a cut of a graph $G = (V, E)$ is the set of edges S , such that the removal of S from G disconnects
? G need not be connected to begin with
 15 G . A graph is connected if there is a path from every vertex to every other vertex in G . The
 16 weight of the cut can be calculated by summing the weights of the edges which will be broken
 17 when a cut is made, as given in equation (1.2.1),
is it necessary that $A \cap B = \emptyset$?

$$\text{cut}(A, B) = \sum_{p \in A, q \in B} w(p, q) \quad (1.2.1)$$

18 where A and B are the sub-graphs formed by cutting graph G , and p and q are the edges
 19 belonging to sub-graphs A and B respectively. The Minimum cut (Wu and Leahy, 1993) is
 20 simply the cut which generates the minimum cut weight. However the minimum cut does

with respect to which objective?

- 1 not always provide the best solution; it tends to favour removing small sections and isolated
why is this a problem for clustering?
- 2 nodes of the graph. The main proposal from Shi and Malik was to consider the cut function
- 3 as a fraction of the total edge connections to all nodes in the graph, instead of using the total
you could explain in words what the difference is
- 4 edge weight connecting the two partitions. This subtle yet important distinction provides
- 5 us with the highly acclaimed normalised cut (NCut), defined in equation (1.2.2), which is
- 6 an integral part of spectral clustering algorithms. The main idea is that big clusters will
- 7 increase $\text{assoc}(A, V)$, thus decrease $\text{Ncut}(A, B)$. This will encourage splitting the data into
- 8 fairly evenly sized clusters, and avoid the minimum cut issue of segmented isolated points.

reference?

- 9 Minimising the normalised cut is an NP-complete problem but an approximate solution can
you could also comment on accuracy of this solution
- 10 be found efficiently by finding the eigenvector with the second smallest eigenvalue. Figure
- 11 1.2.1 shows both the minimum cut and normalised cut solutions for a particular graph. The
- 12 shaded/non-shaded regions represent the partitioning. The minimum cut isolates one node
- 13 from the rest of the graph, whilst the normalised cut provides a more balanced and sensible
- 14 partition.

$$\text{Ncut}(A, B) = \frac{\text{cut}(A, B)}{\text{assoc}(A, V)} + \frac{\text{cut}(A, B)}{\text{assoc}(B, V)} \quad (1.2.2)$$

$$\text{assoc}(A, V) = \sum_{u \in A, t \in V} w(u, t)$$

- 15 On Spectral Clustering Ng et al. (2001) was one of the first papers to provide theoretical
- 16 guarantees on performance for spectral clustering algorithms. Unlike previous authors who
- 17 had only shown empirical results to justify spectral clustering abilities, they prove that their
- 18 spectral clustering algorithm will produce a reasonable clustering, given certain assumptions

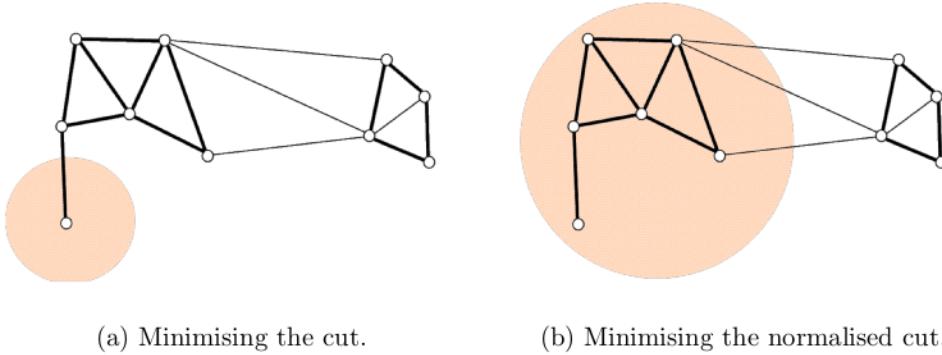


Figure 1.2.1: Two solutions to the bi-partition problem. The partitioning is indicated by shading/non shading of nodes.

This is interesting but you need to make more of it.

- 1 that the clusters are well-spaced. This version of the spectral clustering algorithm has been
- 2 popularly cited throughout the literature and is the one that we shall reference in Section
- 3 1.3.

4 Although spectral clustering has been shown to perform well empirically on simple data
5 sets, computational problems arise as the data set size increases. There have been a num-
6 ber of ways to deal with speed up in the static case, one of the most popular methods is
7 to use Nyström efficiency methods (Williams and Seeger (2001); Fowlkes et al. (2004)).
8 The Nyström method samples the columns of the affinity matrix and approximates the full
9 matrix by using correlations between the sampled columns and the remaining columns. Ef-
10 fectively we can think of this sampling as a dial which the user has control over. Sampling
11 more columns will provide a better results but at a cost. Although Nyström methods are
12 approximation techniques to speed up the computation of spectral clustering, the working
13 memory can still be high. Another drawback with Nyström is that due to random sam-
14 pling it is possible to under represent or entirely miss smaller clusters. One alternative to
15 Nyström methods is to perform a permutation on the data to act as a pre-processing step

why is this? Mention computational complexity of eigen-vector
problem

I don't understand the permutation idea

- 1 to generate a smaller summary data set. By feeding a smaller set of representative points
- 2 into the spectral clustering algorithm instead of the whole data set, we can lessen the effect
- 3 of the computational bottleneck that comes with eigen analysis. Fast Approximate Spectral
- 4 Clustering Yan et al. (2009) explores the theoretical guarantees on misclustering the data set,
Do you mean perturbation?
- 5 given that some permutation has been performed on the data before the spectral clustering.
- 6 Specifically they exactly quantify the misclustering of data sets, given that the original data
- 7 is summarised using K-means (KASP) or a Random Project tree (RASP) as preprocessors.
- 8 This is explored further in Section 1.3.

- 9 Another method to deal with the computational challenge in spectral clustering is Local
- 10 Information-based Fast Approximate Spectral Clustering (Li-ASP) introduced in Cao et al.
- 11 (2014). Li-ASP consists of two upgrades; a sparse affinity graph to speed up computation
- 12 and local interpolation to improve clustering performance. The sparse affinity graph uses
- 13 a k nearest neighbour or an ϵ neighbourhood to set many elements in the affinity graph to
- 14 zero. Local interpolation is suggested based on an issue identified about KASP Yan et al.
- 15 (2009) that if a representative point is miss-assigned to the wrong cluster, then all data
- 16 points represented by that point will also be miss-assigned. Local interpolation can avoid
- 17 these situations by using labelling data points with a weighted version of their p closest
- 18 representative points labels, rather than labelling based just on the label of the single closest
Can you use equations here? I don't quite get what the method is supposed to be doing
- 19 representative point.

- 20 One constant challenge in clustering is selecting the number of clusters to search for.
- 21 Zelnik-manor et al. (2004) present a method for automatically choosing the true number
- 22 of clusters using the eigenvectors to inform their choice. More commonly the eigenvalues
- 23 are used to estimate the number of clusters, but if the clusters are not clearly separated

- 1 identifying the number of clusters from eigenvalues alone is not trivial. We shall assume that
 2 the true number of clusters is known.

3 So far, we have only discussed spectral clustering of static data sets. We are interested
 4 in data streams, which are of great interest in today's world of communication graphs such
You can make more of this also
references ... this should go in the beginning
 5 as the Internet and social networks. Applications also include health care such as modelling
 6 epidemics and understanding sensor networks. Re-clustering the graph at each time step
 7 whenever new information arrives is not feasible especially if data is arriving rapidly.

8 To my knowledge, there is not currently a fully online method for Spectral Clustering. The
 9 problem of performing spectral clustering in data streams has been considered, but framed
 10 as an evolving network rather than a stream where data appears rapidly. Time evolving
 11 graphs are still an interesting problem, and are often found in social networks or biological
 12 applications.

13 The first incremental spectral clustering algorithm concerns topological mapping (Valgren
 14 et al. (2007)). Their algorithm updates the cluster estimates whenever ever a new data
 15 point (or batch of data points) arrives. The cluster membership is updated directly. The
how?
 16 affinity matrix is periodically compressed to deal with large data sets. If a new data point
 17 is sufficiently far from its closest representative points, it is considered the start of a new
 18 cluster, this means that the number of overall clusters must always increase.

19 An incremental update algorithm is proposed in Ning et al. (2007) and Ning et al. (2010)
 20 which can deal with both additional data points joining the network, and similarity weights
 21 changing between existing data points. The algorithm updates the eigenvectors and eigenval-
 22 ues directly without performing a full eigen-decomposition. The addition of a new data point
 23 is treated as a series of n weight changes, where n is the number of currently observed data

1 points. However the authors recommend a full re-clustering in batch to minimise cumulative
 2 errors. There are some issues with update method, mainly that the updating of eigenvectors
 3 means that the orthogonality property may be lost. Also if the spatial neighbourhoods of
 4 often changing vertices are large it can still be computationally difficult as the eigenvector
 5 update step involves the inversion of a matrix. Finally the authors recommend a full spectral
 6 re-clustering occasionally to prevent the accumulation of errors in the eigenvectors, this is not
 7 feasible in the streaming setting. Generally this method is not suitable for data streaming,
 8 as the size of the Laplacian can grow unbounded for an infinite data stream. However since
 9 it is the most well-known relevant “online” algorithm existing, we will compare performance
 10 against it in Section 1.5. [RD: Might move Ning from experiments to 1.4.](#)

11 Other incremental spectral clustering algorithms include Kong et al. (2011), Langone
 12 et al. (2014) and Dhanjal et al. (2011) which approximates the eigen decomposition of the
 13 Laplacian incrementally but still requires regular full re-clustering. insufficient description

14 We consider an online spectral clustering algorithm based on the Clustream model of
 15 Aggarwal et al. (2003). The Clustream algorithm is introduced in full in Section 1.4, but
 16 first we discuss spectral clustering in more detail.

17 1.3 Spectral Clustering

should this be after or before the previous discussion? In my view it should be prior

18 In this Section we provide a brief introduction to Spectral Clustering, and discuss the choice
 19 of affinity matrices. A fast offline spectral clustering algorithm Yan et al. (2009) algorithm
 20 will be introduced, which inspired our implementation in Section 1.4.

21 The goal of clustering algorithms is to partition data ($\mathbf{X} = \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n, \mathbf{x}_i \in \mathbb{R}^d$) into

- 1 k disjoint classes such that each \mathbf{x}_i belongs to one and only one class. Clusters come in all
 2 shapes and sizes, they can be spherical and compact as in Figure 1.3.1a or connected but
I am a little confused by the use of the term "compact" should this be convex? The clusters in 1.3.1b are compact (in the sense that they are closed and bounded)
 3 not visually compact as in Figure 1.3.1b. Data which is compact may be simple to cluster as
 4 the gaps between clusters are easy for simple clustering algorithms like k-means to identify.
 5 Connected but non-compact data sets are can be much more challenging than compact data
 6 sets, and can cause some simple clustering algorithms to fail. Spectral clustering can provide
 7 good quality segmentation on even these difficult cases, however its performance comes at the
 8 cost of computational complexity. We shall use the Jordan-Weiss (NJW) spectral clustering
 9 framework Ng et al. (2001) which is described in Algorithm 1.

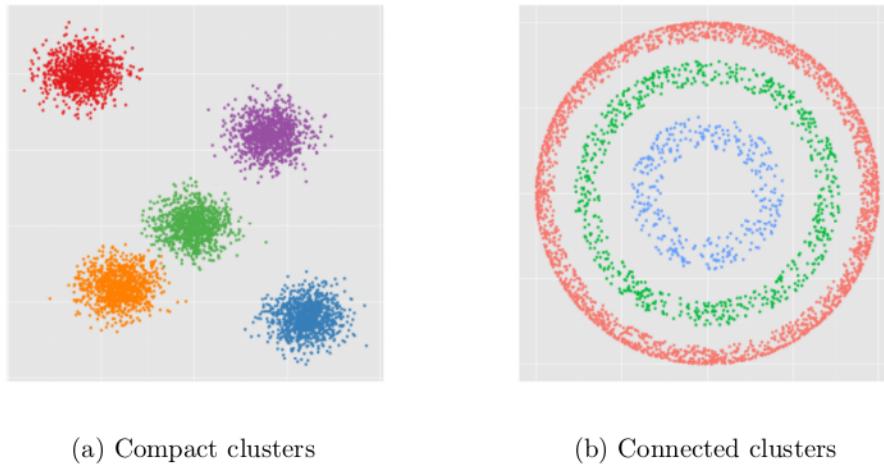


Figure 1.3.1: Examples of different types of clusters

- 10 The affinity matrix $A = (a_{ij})_{i,j=1}^n$ represents the pairwise similarities or distances between
 11 all data points x_i and x_j . A popular choice is to define A to be the Gaussian kernel, as defined
 12 in equation (1.3.1), where the parameter σ controls the width of the local neighbourhoods
 13 which we want to model. What of the ability to detect non-convex clusters through the use of different kernels? this is not mentioned at all

$$a_{i,j} = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right), i, j = 1, \dots, n. \quad (1.3.1)$$

If x_i and x_j are very close, then $a_{i,j} \rightarrow 1$, and if they are far apart $a_{i,j} \rightarrow 0$. A gaussian kernel affinity matrix will have ones along the diagonal and is symmetric ($a_{ij} = a_{ji}$).
The value scaling parameter σ is usually chosen manually. Ng et al. (2001) automatically chose σ by running their clustering algorithm repeatedly for a number of values and selecting the one which provides least distorted clusters. Zelnik-manor et al. (2004) argue that for data which has a cluttered background, or multi-scale data, one global parameter choice for σ is not sufficient. They calculate a localised parameter σ_i for each data point x_i based on it's neighbourhood. In our experimentation we did not find it necessary to do this and use a global σ .

If we mainly wish to model the local relationships, using all of the possible pairwise data connections may not be necessary. It is possible to use a weighted k-nearest neighbour structure to build the affinity matrix once corrections have been made to ensure that this matrix is symmetrical. Another option is to choose some threshold ϵ and only consider connections between data points whose pairwise distances are smaller than ϵ . This is an ϵ -neighbourhood graph. Although it is possible to weight this graph by ϵ , if we choose ϵ to generate a small ϵ -neighbourhood, then the differences between the weights will be small that weighting may become eligible. Therefore a simple construction of the ϵ -neighbourhood graph is as shown in equation (1.3.2).

$$a_{ij}^* = \begin{cases} 1, & \text{if } a_{ij} < \epsilon \\ 0, & \text{otherwise} \end{cases} \quad (1.3.2)$$

1 Using this construction will give a sparse affinity matrix instead of a fully connected
 2 graph, which will help lower the computational complexity. In Section 1.5, we use a fully
 3 connected graph for all experiments.**RD: Possibly KNN now** The degree of each vertex is
 4 defined as, $d_i = \sum_{j=1}^N A_{ij}$, the sum of the rows of the affinity matrix. The degree matrix,
 5 D , is then defined as a diagonal matrix with i^{th} diagonal element equal to d_i . We use the
 6 Normalised Symmetric Laplacian (Chung, 1997) defined in equation (1.3.3), which relates to
 7 an approximation of minimising the Normalised Cut, discussed in Section 1.2. **RD: Should I**
 8 **introduce unnormalised and random walk normalised too?** YES

$$L = D^{-1/2} A D^{-1/2} \quad (1.3.3)$$

9 **RD: Would this be clearer if I explain unnormalised and state** $L_{\text{symm}} = I - D^{-1/2} A D^{-1/2}$

10

11 Spectral clustering can be challenging for very large data sets, constructing the affinity
 12 matrix A and computing the eigenvectors of L have computational complexity $\mathcal{O}(n^2)$ and
 13 $\mathcal{O}(n^3)$ respectively. A fast approximate spectral clustering algorithm is proposed (KASP) is
 14 proposed Yan et al. (2009) which uses a k-means pre processing step to lessen the compu-
 15 tational complexity whilst retaining good clustering performance. Firstly k-means is run on
 16 the whole data set where q is chosen to be large but such that $k \ll q$. The centres of the
 17 clusters are then used as representative data points for the whole data set. Spectral cluster-
 18 ing is performed on the representative set only, which is significantly faster than performing
 19 spectral clustering on the full data set. The resulting cluster labels for the representative
 20 data are linked back to the original data set such that every original data point acquires

Algorithm 1 NJW spectral clustering algorithm

Input: Data set $X = \mathbf{x}_1, \dots, \mathbf{x}_n$, number of clusters k **Output:** k -way partition of the input data

- 1: Construct the affinity matrix A by the following Gaussian kernel function:

$$a_{i,j} = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right), i, j = 1, \dots, n.$$

- 2: Compute the normalised affinity matrix $L = D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$, where D is the diagonal matrix

with $D_{ii} = \sum_{j=1}^n a_{ij}$.

- 3: Compute the k eigenvectors of L , v_1, v_2, \dots, v_k , associated with the k largest eigenvalues, and form the matrix $X = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k]$.

- 4: Renormalize each row of X to form a new matrix Y .

- 5: Partition the n rows into k clusters via a general cluster algorithm, such as the k-means algorithm.

- 6: Assign the original point \mathbf{x}_i to the cluster $k \iff$ the corresponding row i of the matrix Y is assigned to the cluster k .
-

- ¹ the same label as its associated k -means cluster centre. The KASP algorithm is repeated in
- ² Algorithm 2 with unified notation.

Algorithm 2 KASP

Input: Data set $X = \mathbf{x}_1, \dots, \mathbf{x}_n$, number of clusters k , number of representative points q

Output: q -way partition of the input data

- 1: Perform k-means with q clusters on $\mathbf{x}_1, \dots, \mathbf{x}_n$
- 2: Compute the cluster centroids y_1, \dots, y_q as the q representative points.
- 3: Build a correspondence table to associate each x_i with the nearest cluster centroids y_j .
- 4: Run a spectral clustering algorithm on y_1, \dots, y_q to obtain an k -way cluster membership for each of y_q .
- 5: Recover the cluster membership for each x_i by looking up the cluster membership of the corresponding centroid y_j in the correspondence table.

- ³ The KASP algorithm has been shown to perform well empirically, and using representative
- ⁴ points is a sensible way to lessen the computational burden. As the number of representative
- ⁵ points increases, the better performance should be, but the greater the computational
- ⁶ expense.

- ⁷ In online spectral clustering algorithm, we will use representative points as an input for
- ⁸ spectral clustering. Rather than generating representative points from a k-means step, we
- ⁹ constantly update the representative points by using the streaming algorithm, Clustream
- ¹⁰ which is introduced in Section 1.4.

1.4 Micro-cluster based spectral clustering in data streams

2 In this Section we discuss online streaming methods, introduce the streaming algorithm
3 Clustream, and consider how to use incorporate the streaming data into a spectral clustering
4 algorithm.

5 1.4.1 The challenges of clustering data streams

6 A relatively new challenge to clustering is working with data streams Gama (2010); Silva
7 et al. (2013). A data stream is data which arrives in an ordered sequence, continuously; for
8 example, sensor data or online shopping transactions. There is no control over the order in
9 which data objects should be processed. A data stream may be potentially unbounded, and
10 the data points are often discarded after processing. Much work has been done developing
11 offline clustering methods, such as spectral clustering, but it is not suitable to apply these
12 offline methods directly to the streaming scenario. Simply running an offline clustering
13 algorithm on all the data observed so far may not be feasible for three main reasons, storage
14 capacity, computational costs and ability to access to the data.

15 The first challenge is storing all of the data. As a data stream is a potentially an endless
16 sequence of observations obtained at a high frequency often, it may not be possible to store
17 all of the data in its entirety. Therefore the older data has to be thrown away to make room
18 for the new arrivals. This might not seem like a big issue as the data will naturally become
19 weighted temporally, but this is challenging if we wish to incorporate historical data into the
20 clustering.

21 Secondly, as we have discussed, clustering algorithms can be computationally expensive.

1 For example, computing the eigenvectors for spectral clustering has complexity $\mathcal{O}(n^3)$. As
2 data streams are potentially unbounded, standard clustering algorithms cannot be used.
3 Therefore we need to be able to update our idea of the data as new points arrive efficiently
4 and simply with little computational issues.

5 Finally, data streaming is often classed as a “one-pass-access” problem. Imagine a con-
6 stant stream of data flying past your window, you can view the data as it flies by the window,
7 but once it has passed by, it cannot be accessed again. Some traditional clustering algorithms
8 such DBScan (Ester et al. (1996)) require many passes or iterations of the data, therefore
9 these type of clustering methods are not directly suitable for the online data streaming case.

10 **1.4.2 Clustream**

11 Clustream Aggarwal et al. (2003) offers a framework which allows quick and easy updates and
12 the ability to perform sophisticated clustering algorithms. Clustream has proved popular,
13 since the paper was first published in 2003 it has been cited over 1400 times. The main idea
14 is to separate the clustering process in two stages, a micro-clustering stage and a macro-
15 clustering stage. The micro-clustering stage continuously updates statistical summaries of
16 the data stream, and the macro-clustering is more computationally intensive and run in batch
17 or on a user request.

18 **Micro-clustering**

19 The micro-clustering stage is a way of maintaining an active, evolving representative summary
20 of the data, without storing the absolute values of the data points. Micro-clusters are defined
21 as a temporal extension of the cluster feature vector first described in Zhang et al. (1996).

1 The data stream is summarised by many small clusters, which are initially generated by
 2 k -means. The online phase stores q micro-clusters in memory, where q is an input parameter.
 3 We take an initial training set, and perform k-means with q clusters but choose the value of
 4 q to be much larger than the expected number of true macro-clusters k , ($k \ll q$). The aim
 5 here is to create a fine scale summary of the data. The value of q should be chosen to be
 6 as large as computationally comfortable. The larger q is, the finer scale that the summaries
 7 will be. It is vital to ensure that the micro-cluster well represent the underlying data set or
 8 else the macro-clustering will under perform. These q clusters are our first micro-clusters.
 9 Over time, we will update these micro-clusters, adding new data points to them, merging
 10 them and removing old micro-clusters, although the number of micro-clusters should stay
 11 fixed throughout.

12 The micro-clusters can then be used on a user request to perform a macro-clustering using
 13 the summarised data rather than the full data set. If the micro-clusters represent the true
 14 underlying data stream well, then the difference between the clustering on the summarised
 15 data and the true full data should be small.

16 Assume that we have a data stream S which consists of d -dimensional data \mathbf{x}_i arriving
 17 in sequence. $S = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_i, \dots\}$, $\mathbf{x}_i \in \mathbb{R}^d$. Each micro-cluster M_j for ($j \in 1 \dots, q$)
 18 is stored as a $(2 \cdot d + 3)$ tuple $(\mathbf{CF1}_j^x, \mathbf{CF2}_j^x, n_j, CF1_j^t, CF2_j^t)$. The definitions are given
 19 in equation (1.4.1). $\mathbf{CF1}_j^x$ is the sum of all observed data in that micro-cluster, $\mathbf{CF2}_j^x$ is
 20 the sum of the squares of the data and n_j is the number of elements assigned to that micro-
 21 cluster. $CF1_j^t$ and $CF2_j^t$ refer to the sum of the time stamps, and the sum of squared time
 22 stamps respectively. Note that both $\mathbf{CF1}_j^x$ and $\mathbf{CF2}_j^x$ are d -dimensional vectors.

Each micro-cluster M_j will have

$$\begin{aligned}
 \mathbf{CF1}_j^x &= \sum_{x_i \in M_j} \mathbf{x}_i , \\
 \mathbf{CF2}_j^x &= \sum_{x_i \in M_j} (\mathbf{x}_i)^2 , \\
 CF1_j^t &= \sum_{i|x_i \in M_j} t_i , \\
 CF2_j^t &= \sum_{i|x_i \in M_j} (t_i)^2 , \\
 n_j &= \sum_{x_i \in M_j} 1 .
 \end{aligned} \tag{1.4.1}$$

If a new data point x_{new} arrives at time t_{new} and is assigned to micro-cluster M_j , the update given in equation (1.4.2) is applied.

$$\begin{aligned}
 \mathbf{CF1}_j^x &\leftarrow \mathbf{CF1}_j^x + \mathbf{x}_{\text{new}} , \\
 \mathbf{CF2}_j^x &\leftarrow \mathbf{CF2}_j^x + (\mathbf{x}_{\text{new}})^2 , \\
 CF1_j^t &\leftarrow CF1_j^t + t_{\text{new}} , \\
 CF2_j^t &\leftarrow CF2_j^t + (t_{\text{new}})^2 , \\
 n_j &\leftarrow n_j + 1 .
 \end{aligned} \tag{1.4.2}$$

¹ Note that updating the micro-clusters requires only addition therefore updating is cost
² effective. Critically it is possible to use these summaries to calculate the centre of each
³ micro-cluster as in equation (1.4.3). It is these centres which are used as representative points
⁴ for input into the macro-clustering. As new points in the data stream arrive, they are either
⁵ allocated to a micro-cluster and the update procedure discussed above is carried out, or a
⁶ new micro-cluster is created. The decision for a new micro-cluster to be created is based on
⁷ whether the new data point is close enough to its nearest cluster centre.

$$\text{Centre of micro-cluster } j = \bar{\mathbf{M}}_j = \frac{\mathbf{C}\mathbf{F}^x_j}{n_j} \quad (1.4.3)$$

When a new data point arrives it's nearest micro-cluster M_* is identified using the Euclidean distance metric given in equation (1.4.4). If the data point falls within the Maximum Boundary Factor of it's nearest cluster centre, then it is absorbed as part of that cluster. If not, it is used to create a new micro-cluster. However we stated earlier than the number of micro-clusters must remain fixed throughout the process. Therefore if a new micro-cluster is formed, either an existing micro-cluster must be deleted, or two close micro-clusters should be merged. We follow the methodology in Clustream by first looking for an old micro-cluster to delete using the time-stamp references detailed in the original paper and otherwise combine the two nearest micro-clusters. In this way, the algorithm tracks the data stream as it evolves.

$$M^* = \arg \min_{M_j, j \in 1:q} \|\mathbf{x}_i - \bar{\mathbf{M}}_j\|^2 \quad (1.4.4)$$

In Clustream, the maximum boundary factor is defined as a factor of t of the RMS deviation of the data points in M_j from the centroid of M_j . The RMS deviation can only be defined for a cluster with more than one point. For a cluster with only one previous point, a heuristic is used to define the MBF. Details of the MBF can be found in Aggarwal et al. (2003)

With this online micro-cluster maintenance, the data stream should remain well represented over time. When a new data point arrives if it is the start of a new evolving cluster it will be allowed to grow however if it is an outlier no more points will be added to it and

¹ over time it may be deleted from the system all together.

If two micro-clusters M_r and M_s are to be merged, the updates given in equation (1.4.5) are used to merge them into M_r , and M_s will be deleted. Again as all of these updates only involve addition steps, they are fast to implement.

$$\begin{aligned}
 \mathbf{CF1}_r^x &\leftarrow \mathbf{CF1}_r^x + \mathbf{CF1}_s^x, \\
 \mathbf{CF2}_r^x &\leftarrow \mathbf{CF2}_r^x + \mathbf{CF2}_s^x, \\
 CF1_r^t &\leftarrow CF1_r^t + CF1_s^t, \\
 CF2_r^t &\leftarrow CF2_r^t + CF2_s^t, \\
 n_r &\leftarrow n_r + n_s. \tag{1.4.5}
 \end{aligned}$$

² The micro-clustering update algorithm for Clustream is given in Algorithm 4.

³ Macro-clustering stage

⁴ The second stage of clustream is a macro-clustering stage, where we take the current micro-
⁵ cluster feature vectors, and use these as input into global clustering algorithm. The macro-
⁶ clustering step is where the general data summary is transformed into a snapshot of the true
⁷ underlying clusters at that point in the stream. The q micro-cluster centres \bar{M}_j , ($1 \leq j \leq q$)
⁸ are treated as representative points for the data stream S , and a standard clustering algorithm
⁹ can be used to determine clusters. The nature of this algorithm allows the user to get close
¹⁰ to online streaming and perform spectral clustering on a summary of the whole of the data
¹¹ set. The full algorithm is given in Algorithm 4.

¹² There are a number of possible ways to feed the micro-clusters into a spectral clustering
¹³ algorithm. Two of the options suggested in Zhang et al. (1996) are implemented here.

Algorithm 3 Clustream Microclustering

Input: Data Stream $S = \mathbf{x}_1, \dots, \mathbf{x}_n$, number of micro-clusters q

Output: Microclusters M_1, \dots, M_q

Initialise the micro-clusters kmeans(x_1, \dots, x_{init}, q) and equations (1.4.1)

for each new data point x_i **do**

 Find the closest micro-cluster to x_i , M_* using equation (1.4.4)

if x_i satisfies the MBF criterion for M_* **then**

 absorb x_i into micro-cluster M_* using equations (1.4.2)

else

 Use x_i to initialise it's own new micro-cluster using equations (1.4.1)

if any micro-cluster is suitably old **then**

 Remove it

else

 Merge the two closest micro-clusters using equation (1.4.5)

end if

end if

end for

1. Calculate the centre of each micro-cluster \bar{M}_j and use it as an object to be clustered by the macro-clustering algorithm.
2. Do the same as before, but weighting each micro-cluster centre \bar{M}_j proportionally to n_j , the number of points assigned to that micro-cluster, so that micro-clusters with more objects will have more influence on the final clustering.

Algorithm 4 Online Spectral Clustream

Input: Data Stream $S = \mathbf{x}_1, \dots, \mathbf{x}_n$, number of clusters k , number of micro-clusters q

Output: A k way clustering of the microclusters M_1, \dots, M_q .

- 1: Initialise the micro-clusters kmeans(x_1, \dots, x_{init}, q) and equations (1.4.1)
 - 2: **for** each new data point x_i **do**
 - 3: Apply Clustream update as in Algorithm 3
 - 4: **if** A Macro-clustering is required **then**
 - 5: Perform spectral clustering on M_1, \dots, M_q with k clusters.
 - 6: **end if**
 - 7: **end for**
-

- 6: In Section 1.5 we analyse the performance of both unweighted and weighted Online Spectral Clustering, but first we further investigate how to use weighted micro-clusters in Spectral Clustering.

- 9 **1.4.3 Weighting the Micro-Clusters**

- 10 In this section, we discuss how can create a weighted affinity matrix, look at the effect this has on the Laplacians, and note a spectral link between weighting in this manner, and a

¹ larger affinity of repeated points.
² Why would it be beneficial to weight micro-clusters? For example Figure 1.4.1 shows
³ how the distribution of points to microclusters may change as the stream progresses. Figure
⁴ 1.4.1a shows a histogram of the number of points assigned to micro-clusters at the start of
⁵ the stream, Figure 1.4.1b shows the middle of the stream, and Figure 1.4.1c shows the end
⁶ of the stream. We can see that the distribution is not uniform. Therefore some information
⁷ is contained in the number of points assigned to a micro-cluster.



Figure 1.4.1: CREATE PLOTS Histograms showing the number of points assigned to microclusters

⁸ In order to weight the the micro-clusters, we simply construct an affinity as described.
⁹ Let $A \in \mathbb{R}^{q \times q}$ be the affinity matrix of the micro-cluster centres with i, j -th element equal to
¹⁰ the similarity between micro-cluster M_i and M_j ,

$$A_{i,j} = \exp\left(-\frac{\|\bar{M}_i - \bar{M}_j\|^2}{2\sigma^2}\right), \quad i, j = 1, \dots, q.$$

¹¹ .
¹² Let $\tilde{A} \in \mathbb{R}^{q \times q}$ have i, j -th element equal to $n_i n_j A_{ij}$. \tilde{A} is a valid affinity matrix since it
¹³ is symmetric with non-negative entries. If we wish to have $\tilde{A}_{ij} \leq 1$ then simply divide \tilde{A} by
¹⁴ $\max_i n_i^2$, but this makes no difference to the spectral decomposition.

1 There exists a link between the spectral decomposition of the Laplacian generated by \tilde{A}
 2 and the Laplacian arising from a data set of repeated points, which we define as follows. Let
 3 $A^* \in \mathbb{R}^{n \times n}$ be the repeated affinity matrix with the micro-cluster centres repeated based on
 4 the number of points assigned to them. Assume that the columns (and therefore rows) of A^*
 5 are ordered such that the first n_1 are associated with the data assigned to micro-cluster 1,
 6 which has size n_1 and the next n_2 with those assigned to micro-cluster 2, and so on.

7 Let D, \tilde{D}, D^* , be the corresponding degree matrices and L, \tilde{L}, L^* be the corresponding
 8 normalised symmetric Laplacians.

9 Lets look at the Affinity and Laplacian matrices more closely for a very simple case.

10 Assume we have two micro-clusters, M_1 and M_2 , which have n_1 and n_2 points assigned to
 11 them respectively. Let the similarity between the two micro-clusters centres be s , and assume
 12 that we are using the standard Gaussian kernel to generate affinity matrices, so therefore the
 13 diagonal elements will be equal to 1.

$$A = \begin{pmatrix} 1 & s \\ s & 1 \end{pmatrix}, \quad D = \begin{pmatrix} 1+s & 0 \\ 0 & 1+s \end{pmatrix}, \quad L_{\text{symm}} = \begin{pmatrix} \frac{1}{1+s} & \frac{s}{1+s} \\ \frac{s}{1+s} & \frac{1}{1+s} \end{pmatrix}$$

14 In order to create a weighted version of the affinity matrix, we simple multiply through
 15 by n_1 and n_2 . We can see how this is incorporated into the Laplacian. Note that we are
 16 working with the Symmetric Normalised Laplacian. In order to weight in this manner for
 17 the unnormalised Laplacian a scaling of the diagonal matrix is required.

$$\tilde{A} = \begin{pmatrix} n_1^2 & sn_1n_2 \\ sn_1n_2 & n_2^2 \end{pmatrix}, \quad \tilde{D} = \begin{pmatrix} n_1^2 + sn_1n_2 & 0 \\ 0 & n_2^2 + sn_1n_2 \end{pmatrix}$$

$$\tilde{L}_{\text{symm}} = \begin{pmatrix} \frac{n_1^2}{n_1^2 + sn_1 n_2} & \frac{sn_1 n_2}{\sqrt{(n_1^2 + sn_1 n_2)(n_2^2 + sn_1 n_2)}} \\ \frac{sn_1 n_2}{\sqrt{(n_1^2 + sn_1 n_2)(n_2^2 + sn_1 n_2)}} & \frac{n_2^2}{n_2^2 + sn_1 n_2} \end{pmatrix}$$

- ¹ Now lets look at the repeated affinity matrix. Here we can see the block nature in A^* ,
² D^* and L^* .

$$A^* = \begin{pmatrix} 1 & \dots & 1 & s & \dots & s \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 1 & \dots & 1 & s & \dots & s \\ s & \dots & s & 1 & \dots & 1 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ s & \dots & s & 1 & \dots & 1 \end{pmatrix}, \quad D^* = \begin{pmatrix} \star & 0 & 0 & 0 & \dots & 0 \\ 0 & \ddots & 0 & \vdots & \ddots & \vdots \\ 0 & 0 & \star & 0 & \dots & 0 \\ 0 & \dots & 0 & \Delta & 0 & 0 \\ \vdots & \ddots & \vdots & 0 & \ddots & 0 \\ 0 & \dots & 0 & 0 & 0 & \Delta \end{pmatrix}$$

- ³ where $\star = n_1 + n_2 s$ and $\Delta = n_1 s + n_2$.

$$L_{\text{symm}}^* = \begin{pmatrix} \frac{1}{\star} & \dots & \frac{1}{\star} & \frac{s}{\sqrt{\star\Delta}} & \dots & \frac{s}{\sqrt{\star\Delta}} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \frac{1}{\star} & \dots & \frac{1}{\star} & \frac{s}{\sqrt{\star\Delta}} & \dots & \frac{s}{\sqrt{\star\Delta}} \\ \frac{s}{\sqrt{\star\Delta}} & \dots & \frac{s}{\sqrt{\star\Delta}} & \frac{1}{\Delta} & \dots & \frac{1}{\Delta} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \frac{s}{\sqrt{\star\Delta}} & \dots & \frac{s}{\sqrt{\star\Delta}} & \frac{1}{\Delta} & \dots & \frac{1}{\Delta} \end{pmatrix}$$

- ⁴ If we evaluate these expressions for a particular case, we can see how the spectral decom-
⁵ position of the matrices are linked.

- ⁶ Let $s = 0.5$, $n_1 = 3$, $n_2 = 2$.

The 2nd smallest eigenvector of L_{symm}^* is

$$e_2^* = \begin{bmatrix} -0.350 & -0.350 & -0.350 & 0.562 & 0.562 \end{bmatrix}$$

The 2nd smallest eigenvector of \tilde{L}_{symm} is

$$\tilde{e}_2 = \begin{bmatrix} 0.607 & -0.795 \end{bmatrix}$$

- ¹ If we expand the eigenvector \tilde{e}_2 by repeating it's elements and block divide by $\sqrt{n_1}$ and
- ² $\sqrt{n_2}$ respectively, we get

$$\left(\overbrace{\frac{0.607}{\sqrt{3}} \frac{0.607}{\sqrt{3}} \frac{0.607}{\sqrt{3}}}^{n_1} \overbrace{\frac{-0.795}{\sqrt{2}} \frac{-0.795}{\sqrt{2}}}^{n_2} \right) = \begin{pmatrix} 0.350 & 0.350 & 0.350 & -0.562 & -0.562 \end{pmatrix}$$

- ³ which is the negative of the 2nd smallest eigenvector of L_{symm}^* , e_2^* . In fact, we will always
- ⁴ have that the expanded repeated eigenvector of the weighted Laplacian equal to either e_k^* or
- ⁵ $-e_k^*$ for all k .

- ⁶ We can see more generally that, [RD: What can we say generally?](#)

⁷ Let \tilde{u} be the second smallest eigenvector of \tilde{L} .

⁸ Define $u^* = \left(\overbrace{\frac{\tilde{u}_1}{\sqrt{n_1}}, \frac{\tilde{u}_1}{\sqrt{n_1}}, \frac{\tilde{u}_1}{\sqrt{n_1}}}^{n_1}, \dots, \overbrace{\frac{\tilde{u}_k}{\sqrt{n_k}}, \frac{\tilde{u}_k}{\sqrt{n_k}}}^{n_k} \right)$.

⁹ We can see that $\|u^*\|^2 = \sum_{i=1}^n (u_i^*)^2 = \sum_{i=1}^k n_i \cdot \frac{\tilde{u}_i^2}{n_i} = \sum_{i=1}^k \tilde{u}_i^2 = 1$.

¹⁰ 1.5 Experimentation

- ¹¹ In the experimentation section, we investigate the performance of Spectral Clustream, both
- ¹² unweighted and weighted variants, along with a simple windowed approach.

¹ **The Algorithms**

² Spectral Clustream is our implementation of Algorithm 3. As highlighted in Section 1.4,
³ there are two ways we can incorporate the micro-cluster centres into the macro-clustering.
⁴ Unweighted clustream takes the micro-cluster centres as direct input into the spectral cluster-
⁵ ing algorithm. Weighted clustream weights the micro-cluster centres by the number of data
⁶ points assigned to that micro-cluster. In both clustream algorithms, we use 150 micro-clusters
⁷ to monitor the streams.

⁸ Windowed Spectral Clustering is a simple algorithm stores a window of size w of the most
⁹ recently observed data points. When a new data point is observed, the oldest data point in
¹⁰ the window is discarded to make room for the new data point. Only the data points in the
¹¹ current window are available for input into a standard Spectral Clustering algorithm. We
¹² use a fixed window size of 150 data points for the duration of the stream.

¹³ **1.5.1 Methodology**

¹⁴ A simulated data stream $S = \mathbf{x}_1, \dots, \mathbf{x}_n$ is observed sequentially, with one data point \mathbf{x}_t ob-
¹⁵ served at each time point $t \in 1, \dots, n$. Clustream requires the micro-clusters to be initialised.
¹⁶ We do this by applying kmeans with 150 clusters to the first 500 points. The initial micro-
¹⁷ cluster components are then calculated using the equations described in equation (1.4.1).
¹⁸ Spectral Windowed uses 150 data points to initialise the data window.

¹⁹ At each time point t , a new data point \mathbf{x}_t is observed and is used to update the streaming
²⁰ algorithms. For clustream algorithms, this means updating your micro-clusters in the way
²¹ outlined in Algorithm 3. For windowed spectral clustering this is done by just shifting your

1 window along by one, forgetting the oldest data point $\mathbf{x}_{(t-150)}$ and including \mathbf{x}_t .
 2 In batch, every ten time steps, we use the current snapshot of the data to apply spectral
 3 clustering. In clustream this means feeding the centres of the micro-clusters (weight adjusted
 4 or not) into a standard spectral clustering algorithm. In windowed spectral, this means using
 5 all data points in your window as input into a standard spectral clustering algorithm.
 6 We then evaluate the spectral clustering performance. The following step is repeated 10
 7 times, and results are averaged out. Simulate a new test set of 200 data points using the
 8 underlying, unknown data stream parameters at time t . Assign each of the 200 points to
 9 their closest micro-cluster centre (or data point in the current window) using the nearest
 10 neighbour algorithm. These 200 points will then take on the macro-cluster assignment given
 11 to their nearest neighbour. Performance is measured in terms of purity and V-measure using
 12 the known true clusters of the simulated data.

13 Note that when comparing performance on the real data sets obviously we cannot simulate
 14 from the underlying distribution. Instead we look forward to the next 200 data points, assign
 15 them using nearest neighbour and continue in the manner above.

16 **Performance Measures**

17 The two measures we use to quantify cluster performance are purity and V-measure, both
 18 of which are well used in the clustering literature. Both measures require knowledge of the
 19 “true class” of the data points, which may not always be available for real data sets. Purity
 20 is the more intuitive of the two measures to understand as it quantifies how pure clusters are.
 21 Clusters which consist of just data points of one cluster will contribute towards a high purity
 22 score, whilst clusters consisting of data points from many different true classes will lower the

¹ purity. To calculate purity, for each cluster find the true class which is most prevalent in
² that cluster and count how many of that class are in that cluster. Do this for all clusters, sum
³ these counts and divide by the total number of data points. This is shown mathematically
⁴ in equation (1.5.1).

$$\text{Purity} = \frac{1}{N} \sum_k \max_j |\omega_k \in c_j|, \quad (1.5.1)$$

⁵ where we have clusters c_j and true classes ω_k .

⁶ V-measure (Rosenberg and Hirschberg (2007)) takes the harmonic mean of two other
⁷ performance measures, homogeneity and completeness. Homogeneity assesses if each cluster
⁸ contains members of only a single class, whilst completeness checks that all members of the
⁹ same class are assigned to the same cluster. This is shown in equation (1.5.2).

$$\text{V-measure} = 2 \frac{h \times c}{h + c}, \quad (1.5.2)$$

¹⁰ where h and c are homogeneity and completeness measures respectively defined in Rosen-
¹¹ berg and Hirschberg (2007). Both purity and V-measure take values between 0 and 1, where
¹² 1 indicates perfect performance.

You would need to define h and c.

¹³ 1.5.2 Simulated Results

¹⁴ The three simulated data sets we use are all normally distributed, one is static for the duration
¹⁵ of the stream, one exhibits an abrupt change, and another illustrates concept drift.

¹⁶ The static data set is a simulation from a multivariate normal distribution with parameters
¹⁷ $\pi = (0.5, 0, 5)$, $\mu_1 = (0, 0)$, $\mu_2 = (5, 5)$, and unit variance. There are 2 dimensions, 2 clusters,

1 and the data set remains static for the duration of the stream. This is shown in Figure
 2 1.5.1a. The abrupt change data set is initially distributed the same as the static data set, but
 3 at $t = 1000$, μ_1 changes from $(0, 0)$ to $(15, 15)$, and stays as such until the end of the stream.
 4 Figure 1.5.1b shows the abrupt data set before the change point and Figure 1.5.1c shows the
 5 data at a point in the stream after the change has occurred. The model used for simulating
 6 gradual concept drift is

$$(\mathbf{X}_t | Y_t = 0) \sim \mathcal{N}(\boldsymbol{\mu}_t^0, \boldsymbol{\Sigma}) \quad (1.5.3)$$

$$(\mathbf{X}_t | Y_t = 1) \sim \mathcal{N}(\boldsymbol{\mu}_t^1, \boldsymbol{\Sigma}), \quad (1.5.4)$$

7 where $Y_t \sim Bern(p)$ for some user defined p and $\boldsymbol{\Sigma}$ is a multiple of the identity matrix.
 8 For our data set $p = 0.5$ and $\boldsymbol{\Sigma} = 0.005 * \mathbf{I}$. In order to model gradual concept drift
 9 the mean vectors are chosen such that the paths $\boldsymbol{\mu}_1^{(0)}, \dots, \boldsymbol{\mu}_t^{(0)}$ and $\boldsymbol{\mu}_1^{(1)}, \dots, \boldsymbol{\mu}_t^{(1)}$ trace out
 10 complimentary orbits on a d dimensional hyper-sphere. We have $d = 3$, and set the means
 11 such that a quarter of a sphere is covered during the stream. The path traces can be viewed
 12 in Figures 1.5.1d, 1.5.1e and 1.5.1f. The simulation design needs to be revisited based on comments below.

13 Figure 1.5.2 shows performance of unweighted spectral clustream, weighted spectral clus-
 14 tream and windowed spectral clustering on a number of simulated data sets.

15 The first thing to note is that unweighted clustream and the windowed algorithm both
 16 generally perform excellently on all data sets for the duration of the stream. The gradual
 17 drift data set doesn't seem to be too challenging for these algorithms, and although there is
 18 a downward spike in the abrupt data set at the point of change, both algorithms adapt and
 19 recover quickly. The weighted version of clustream starts with good performance in all three

As we discussed in our meeting these are very easy and for this reason don't reveal anything really new/ interesting. The original clustering problem has to be harder, and/ or drift must be occur faster

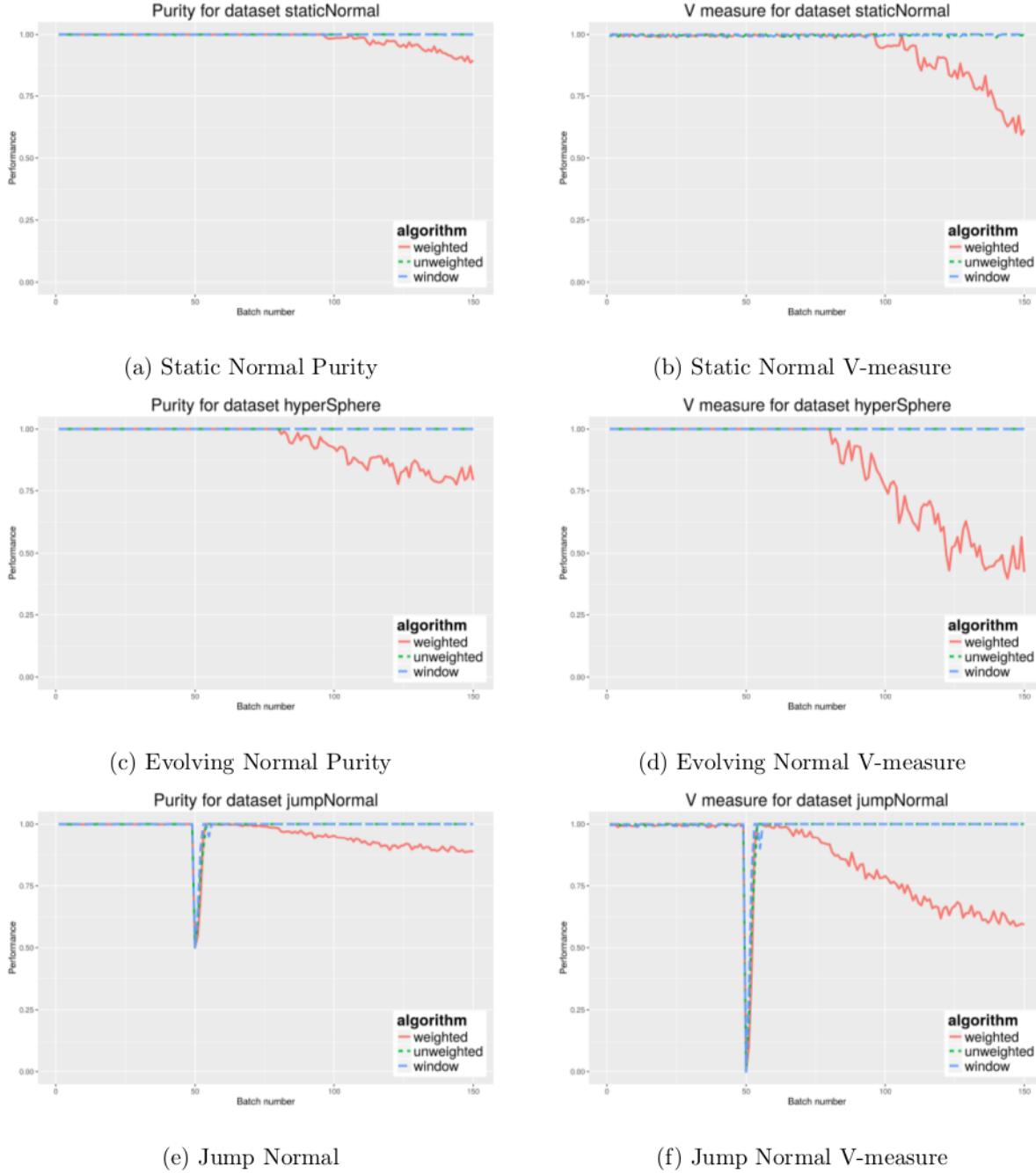


Figure 1.5.1: Descriptions of the simulated data sets. Figure 1.5.1a shows the Static Normal Data set. Figures (1.5.1b, 1.5.1c) shows the Jump data set before and after the changepoint. Figures (1.5.1d, 1.5.1e, 1.5.1f) show the evolving data set at the start, in the middle and at the end of the stream.

¹ sets, but rapidly deteriorate as the streams progress, at times dropping to V-measure values
² of below 0.5.

³ 1.5.3 Real Data Results

⁴ We investigate the performance of Clustream spectral clustering and Windowed Spectral
⁵ clustering on real data. The data set is taken from the UCI Pendigits data set and consists



These are two easy. From these experiments the only conclusion is that weighted is fundamentally wrong (for reasons we discussed). There is hardly any learning curve involved in the static dataset. Both methods do perfectly from the very step time-step. In the abrupt change case, the response seems to be almost instantaneous and there is not noticeable difference between the two methods. Is this real or is this an artifact of the simulation? We need to explore when the simple windowed approach fails. What is its sensitivity to the window size? There is also no indication of variability of

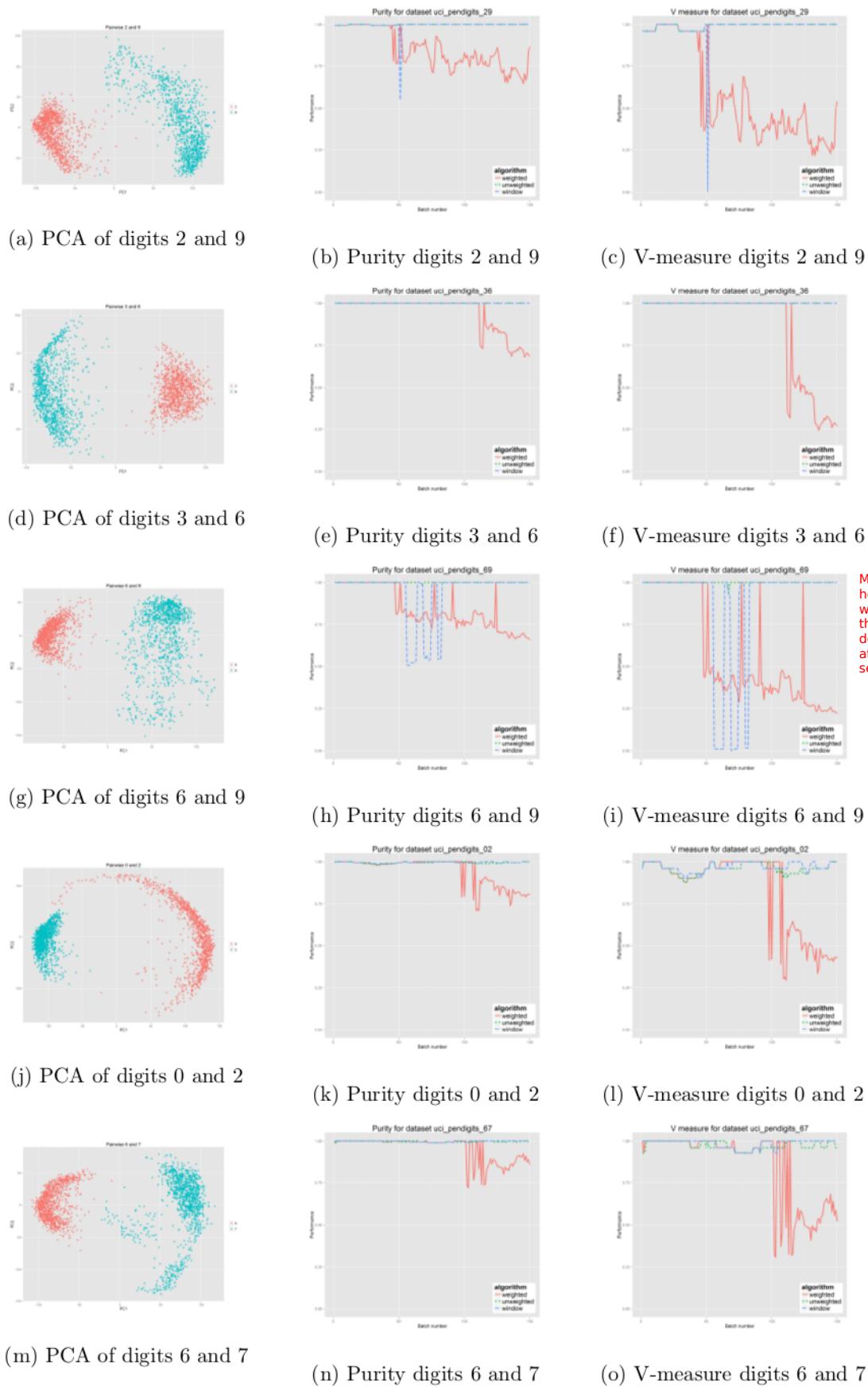
Figure 1.5.2: Performance of Spectral Clustream (weighted and unweighted) and windowed performance. This is very relevant for the entire chapter. If our results simply look like this then there is not much to learn and it will be hard to justify novelty

spectral

- ¹ of hand drawn digits of the numbers 0-9. There are 250 samples taken from 44 writers. The
- ² data was collected using a pressure sensitive tablet. There are 16 features, each relating to

1 the co-ordinate information taken from the input tablet. We restrict our analysis to pairwise
2 comparison of digits. For example we attempt to cluster the digits 0 and 1, and treat the
3 data as if it is arriving in a constant data stream.

4 The results for a selection of the pairwise digits are shown in Figure 1.5.3. The first
5 column displays the digit data in PCA space, the second column shows the purity, and the
6 third column shows the V-measure. The pendigit data is more challenging than the simulated
7 data, however both windowed spectral and the unweighted spectral clustream generally do
8 well. Again we observe that weighted Clustream has poor performance, and it seems that once
9 the performance drops, it struggles to improve again. Interestingly, occasionally windowed
1 spectral clustering seems to heavily drop in performance, making a bad clustering, but it
2 recovers quickly due to the nature of the relatively small window size.



More or less the same comments apply here. This is the only case where we have something interesting, the window method completely breaks down. Why is this happening? Can we attempt to replicate this in different settings?

these problems appear to be too easy

You can use the pendigits data to simulate both gradual and abrupt drift and this can be very helpful

Figure 1.5.3: Pendigits Pairwise - Clustream and windowed spectral

3 Texture data

4 We also tested there methods on another set of real data based on textured images. The
 5 Kylberg (Kylberg, 2011) texture data set consists of 28 texture classes with 160 unique
 6 texture patches per class. The patches consist of 576×576 pixel images.

7 In order to extract features to cluster, I used the LS2W (Eckley and Nason, 2011) method
 8 to extract wavelet features of from the textured images.

9 I selected a subset of 6 classes to run the algorithms on, examples of all 6 textures are
 10 shown in Figure 1.5.4 . The classes selected represent, two types of blanket, canvas, a ceiling,
 1 some lentils and a screen.

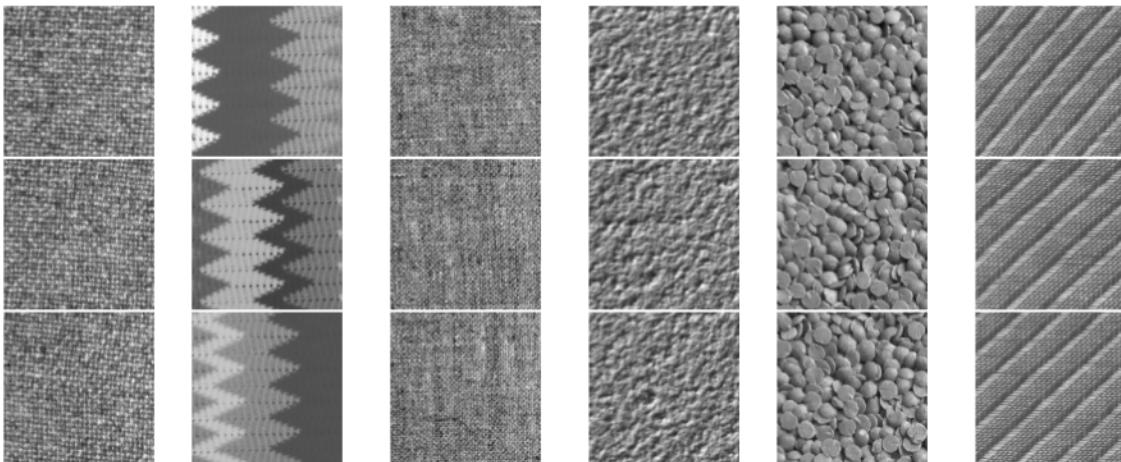


Figure 1.5.4: Three examples from each of the 6 different texture tiles used in the streaming data set. The texture classes are (L-R) Blanket 1, Blanket 2, Canvas, Ceiling, Lentils and Screen

2 The performance plots for the texture data are shown in Figure 1.5.5. Here we can see
 3 a difference between windowed spectral clustering and unweighted clustream, the windowed
 4 approach is generally performing better. Once again weighted clustream does not perform

5 well, and performance seems to decline as the stream progresses.

This makes us look very bad. Have you tried experimenting with window size, number of micro-clusters?

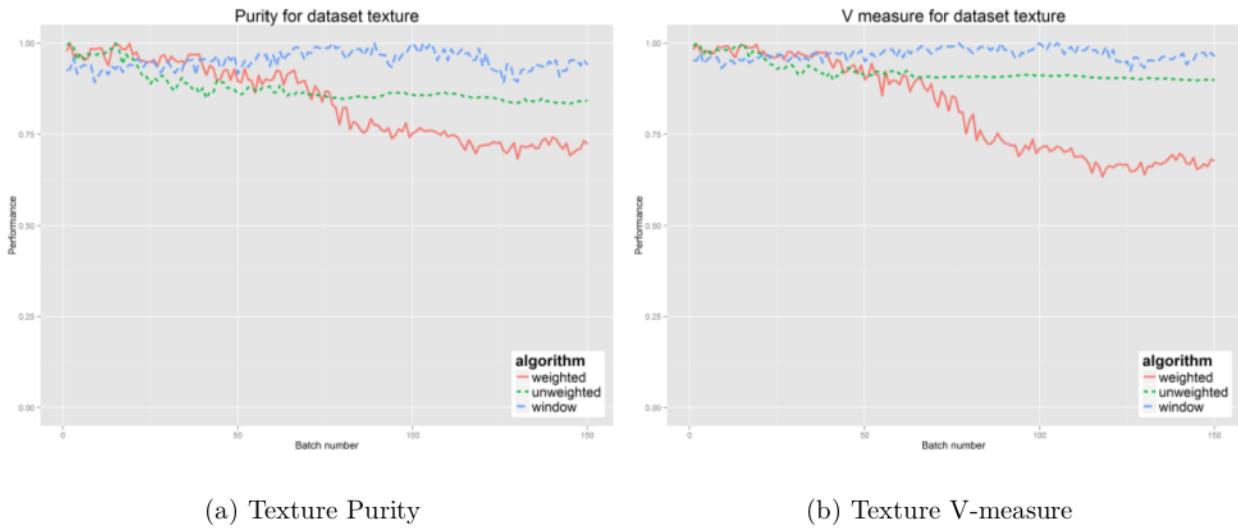


Figure 1.5.5: Texture Results

6 1.5.4 Optional Extras

- 7 • Evaluate the number of micro-clusters needed for good performance. Give some insight
into choosing the number of micro-clusters q . How does performance change based on
the value of q ? Aggarwal et al. (2003) suggests using the biggest value of q computa-
tionally possible. I suspect diminishing returns when increasing q . Can I say something
about this?
- 12 • Show the distribution of micro-cluster sizes. Often it is very non-uniform, is it this
which is effecting the performance of weighted spectral clustering?
- 1 • Explain why Ning was not used in experimental results - too slow!

³ **Spectral Clustream vs Ning (without re-clustering)**

⁴ In order to implement the Ning algorithm, the authors recommend a regular full re-clustering
⁵ of the data sets to stop errors in the eigenvector estimates growing too large.

⁶ In Figure 1.5.6 I run both Ning and Clustream on the “Norm overlap” data set. We do
⁷ not perform re-clustering for the Ning algorithm at all, so we can view how bad performance
⁸ may get over time. This was run over a fairly short stream of just 200 data points. Cluster
⁹ performance was requested every 10 time steps, resulting in 20 performance points.

¹⁰ Although performance does seem to deteriorate over time, it isn’t as dramatic as I
¹¹ expected. It would be nice to see this for a much longer stream (5000 data points?)

¹² We can also see that Spectral Clustream is generally performing much better in terms of
¹³ both purity and V-measure. It would be good to see something similar on real data.

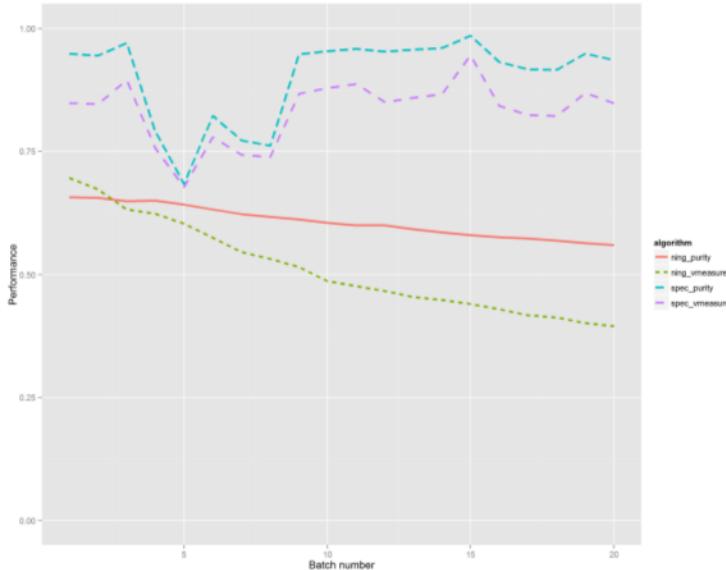


Figure 1.5.6: Ning (no re clustering) vs Spectral Clustream

¹ However, I have found Ning very slow to run. The update method can be very expensive
² for fully connected networks, which makes it not feasible to use as it is in this study. However

3 by using an ϵ neighbourhood to force the updates to be minimal I should be able to achieve a
4 vast speed up. This will be tried on my return. I will also implement the regular re-clustering
5 as recommended by the authors but this re-clustering will have to be over a window as it is
1 not possible to performance a truly full re-clustering in a streaming setting.

2 1.6 Conclusion

- 3 • The way to compare clustering algorithms in a streaming setting is not obvious
4 • We may want a clustering to accurately represent the data stream in it's current state,
5 however we may also wish to incorporate older stream information into the clustering.

The main question that needs to be considered in this chapter is what is the novelty. Given the state of the experimental results and the fact that we don't really prove or propose something very novel this is a severe limitation. In my view we need to identify situations where the obvious method (window) fails, and then to explore how we can improve over that. There seems to be large room to do this. So far not even different window sizes, # of micro-clusters etc are explored.

The write-up in the introduction will need to reflect the main novelty, so this can be postponed.

For the literature review and the exposition of Spectral Clustering: You need to be more precise and to use more mathematical notation wherever possible. See the comments there and consider how they apply to the rest of the text.

⁶ Bibliography

- ⁷ Aggarwal, C. C., Watson, T. J., Ctr, R., Han, J., Wang, J., and Yu, P. S. (2003). A
⁸ Framework for Clustering Evolving Data Streams. *Proceedings of the 29th international*
⁹ *conference on Very large data bases*, pages 81–92.
- ¹⁰ Alpert, S.-Z. Y. C. J. (1995). Spectral Partitioning: The More Eigenvectors, The Better. In
¹¹ *32nd Design Automation Conference*, pages 195–200. ACM.
- ¹² Cao, J., Chen, P., Dai, Q., and Ling, W. K. (2014). Local information-based fast approximate
¹³ spectral clustering. *Pattern Recognition Letters*, 38(1):63–69.
- ¹⁴ Chung, F. R. (1997). *Spectral graph theory*, volume 92. American Mathematical Soc.
- ¹⁵ Dhanjal, C., Gaudel, R., and Clémenton, S. (2011). Incremental spectral clustering with the
¹ normalised laplacian.
- ² Donath, W. E. and Hoffman, A. J. (1973). Lower Bounds for the Partitioning of Graphs.
³ *IBM Journal of Research and Development*, 17(5):420–425.
- ⁴ Eckley, I. and Nason, G. (2011). LS2W: Implementing the locally stationary 2d wavelet
⁵ process approach in R. *Journal of Statistical Software*, 43(3):1–23.

- 6 Ester, M., Kriegel, H.-P., Sander, J., and Xu, X. (1996). A density-based algorithm for
7 discovering clusters in large spatial databases with noise. *Kdd*.
- 8 Fieldler, M. (1973). Algebraic Connectivity of Graphs. *Czechoslovak Mathematical Jurnal*,
9 23(2):298–305.
- 10 Fowlkes, C., Belongie, S., Chung, F., and Malik, J. (2004). Spectral grouping using the
11 Nyström method. *IEEE transactions on pattern analysis and machine intelligence*,
12 26(2):214–225.
- 13 Gama, J. (2010). *Knowledge discovery from data streams*, volume 12. CRC Press.
- 14 Kong, T., Tian, Y., and Shen, H. (2011). A Fast Incremental Spectral Clustering for Large
15 Data Sets. In *2011 12th International Conference on Parallel and Distributed Computing,*
16 *Applications and Technologies*, pages 1–5. IEEE.
- 17 Kylberg, G. (2011). The kylberg texture dataset v. 1.0. External report (Blue series) 35, Cen-
18 tre for Image Analysis, Swedish University of Agricultural Sciences and Uppsala University,
19 Uppsala, Sweden.
- 1 Langone, R., Mauricio Agudelo, O., De Moor, B., and Suykens, J. A. (2014). Incremental
2 kernel spectral clustering for online learning of non-stationary data. *Neurocomputing*,
3 139:246–260.
- 4 Ng, A. Y., Jordan, M. I., and Weiss, Y. (2001). On Spectral Clustering: Analysis and an
5 algorithm. *Advances in Neural Information Processing Systems*, pages 849–856.

- 6 Ning, H., Xu, W., Chi, Y., Gong, Y., and Huang, T. (2007). Incremental Spectral Cluster-
7 ing With Application to Monitoring of Evolving Blog Communities. *SIAM International
8 Conference on Data Mining*, pages 261–272.
- 9 Ning, H., Xu, W., Chi, Y., Gong, Y., and Huang, T. S. (2010). Incremental spectral clustering
10 by efficiently updating the eigen-system. *Pattern Recognition*, 43(1):113–127.
- 11 Rosenberg, A. and Hirschberg, J. (2007). V-measure: A conditional entropy-based exter-
12 nal cluster evaluation measure. *Proceedings of the 2007 Joint Conference on Empirical
13 Methods in Natural Language Processing and Computational Natural Language Learning*,
14 1(June):410–420.
- 15 Shi, J. and Malik, J. (2000). Normalized cuts and image segmentation. *IEEE Transactions
16 on Pattern Analysis and Machine Intelligence*, 22(8):888–905.
- 17 Silva, J. a., Faria, E. R., Barros, R. C., Hruschka, E. R., Carvalho, A. C. P. L. F. D., and
18 Gama, J. (2013). Data stream clustering. *ACM Computing Surveys*, 46(1):1–31.
- 19 Valgren, C., Duckett, T., and Lilienthal, A. (2007). Incremental Spectral Clustering and Its
20 Application To Topological Mapping. In *Proceedings 2007 IEEE International Conference
1 on Robotics and Automation*, pages 4283–4288. IEEE.
- 2 von Luxburg, U., Belkin, M., and Bousquet, O. (2008). Consistency of spectral clustering.
3 *The Annals of Statistics*, 36(2):555–586.
- 4 Williams, C. K. and Seeger, M. (2001). Using the Nystrom Method to Speed up Kernel
5 Machines. In *Advances in Neural Information Processing Systems 13*, pages 682–688.

- 6 Wu, Z. and Leahy, R. (1993). An optimal graph theoretic approach to data clustering: theory
7 and its application to image segmentation. *IEEE Transactions on Pattern Analysis and*
8 *Machine Intelligence*, 15(11):1101–1113.
- 9 Yan, D., Huang, L., and Jordan, M. I. (2009). Fast approximate spectral clustering. In *Work*,
578 page 907, New York, New York, USA. ACM Press.
- 579 Zelnik-manor, L., Zelnik-manor, L., Perona, P., and Perona, P. (2004). Self-tuning spectral
580 clustering. *Advances in Neural Information Processing Systems 17*, pages 1601–1608.
- 581 Zhang, T., Ramakrishnan, R., and Livny, M. (1996). BIRCH: An Efficient Data Clustering
582 Method for Very Large Databases. *ACM SIGMOD Record*, 25(2):103–114.