

Clustering Optimisation Method for Highly Connected Biological Data

Richard Tjörnhammar^{a,b}

^a*KTH Royal Institute of Technology, SE-100 44 Stockholm, Sweden*

^b*SciLifeLab, Tomtebodavägen 23, SE-171 65 Solna, Sweden*

Abstract

Currently, data-driven discovery in biological sciences resides in finding segmentation strategies in multivariate data that produce sensible descriptions of the data. Clustering is but one of several approaches and sometimes falls short because of difficulties in assessing reasonable cutoffs, the number of clusters that need to be formed or that an approach fails to preserve topological properties of the original system in its clustered form. In this work, we show how a simple metric for connectivity clustering evaluation leads to an optimised segmentation of biological data.

The novelty of the work resides in the creation of a simple optimisation method for clustering crowded data. The resulting clustering approach only relies on metrics derived from the inherent properties of the clustering. The new method facilitates knowledge for optimised clustering, which is easy to implement.

We discuss how the clustering optimisation strategy corresponds to the viable information content yielded by the final segmentation. We further elaborate on how the clustering results, in the optimal solution, corresponds to prior knowledge of three different data sets.

Keywords: Clustering, Connectivity, Unimodal Optimisation, Dimensionality Reduction, Statistical Learning, Hierarchical Agglomerative Clustering

1. Introduction

One key feature of biological data is the excessive amount of viable information. Real-world systems describing a simple cellular or biochemical process are often large, containing many active reagents in a crowded environment [10, 2, 9, 26]. Depending on the nature of the interactions belonging to the constituents of

Email address: richardt@kth.se (Richard Tjörnhammar)

such a system various dimensionality reduction techniques can be employed to coarse grain the system and reduce the complexity of the studied problem [19, 17]. Determining an optimal number of clusters for more in-depth analysis or visualisation is an open problem with many different solutions [16, 12, 21, 20]. Here, we present a simplified optimisation approach and algorithms to achieve this task.

We employ clustering with a naming scheme where clustering segmentation produced via 'connection clustering', see methods, is referred to as Connectivity Clustering Algorithms (CCA) and 'linkage' based clustering, commonly employed in Agglomerative Hierarchical Clustering (AHC) [28], is referred to as Link Clustering Algorithms (LCA) [23]. CCA constitutes a point evaluation of the system distance matrix for a single distance cutoff ($D_{ij}, \epsilon \in \mathbb{R}^+$ and $D_{ij} = D_{ji}$) while LCA evaluates the linkage matrix describing the entire hierarchy. Here D_{ij} describes all the pairwise distances between the parts of the entire system. CCA will let you determine if distance matrix indices are connected at some distance or not. The connection-based methods establish the number of clusters in the binary Neighbour matrix without constructing an intermediate linkage matrix.

The Neighbour matrix is defined here as the pairwise distance between the parts i and j of the system (D_{ij}) with an applied cutoff ($N_{ij} = D_{ij} \leq \epsilon$) and is related to the adjacency matrix from graph theory by adding an identity matrix to the adjacency matrix ($A_{ij} = N_{ij} - I_{ij}$). The three boolean matrices that describe a system at some distance cutoff (ϵ) are: the Identity matrix ($I_{ij} = D_{ij} \equiv 0$), the Adjacency matrix ($A_{ij} = D_{ij} \leq \epsilon - I_{ij}$) and the Community matrix ($O_{ij} = D_{ij} > \epsilon$). We note that summing the three matrices will return 1 for any i, j pair. CCA determines the number of clusters by traversing N_{ij} and evaluates if there is any true overlap for a specific distance cutoff. Publically available CCA methods include the connectivity Algorithm 1 in the methods section as well as the Density-Based Algorithm for Discovering Clusters (DBSCAN) [5, 19] without point rejections.

Linkage algorithms determine the number of clusters for all unique distances by forming the linkage matrix reducing and ignoring some connections to already linked constituents of the system in accord with a chosen heuristic. A Linkage algorithm is not an unambiguous treatment of a system [6] where all the true connections in it are important, such as in a molecular water bulk system, when you want all your quantum-mechanical water molecules to be treated at the same level of theory based on their connectivity at a specific distance. If you are doing statistics on a complete hierarchy then this distinction is not important. You can construct hierarchies from both algorithm types but a connection algorithm, without point rejection criteria, will always produce a unique and well-determined structure while the link algorithms will be unique but structurally dependent on how ties are resolved and which heuristic is employed for construction. The connection

hierarchy is exact and unique, but slow to construct, while the link hierarchies are heuristic dependent, but fast to construct. The Linkage algorithms are more efficient at creating a hierarchy based on a distance matrix representation of the data but can be thought of as throwing away information at every linking step. The full link algorithm determines the new cluster distance to the rest of the unclustered points in a self-consistent fashion by employing two different heuristics. Minimal distances for assigning a cluster link are determined by finding the minimum non-diagonal element in the distance matrix and the new link group distance to all remaining points are determined using the second heuristic. Using simple linkage, or min value distance assignment, ensures that the same heuristic is employed both when finding links as well as when assigning the link group to system distances. We will see that it will also produce an equivalent clustering as compared to the one deduced by a connection algorithm for a specific distance ϵ . Except for some of the cases when there are distance ties in the link evaluation. This is a computational quirk that does not affect 'connection' based hierarchy construction.

For a specific ϵ we have C cluster segments ($\dim(C) = c$ and $c \in \mathbb{Z}^+$) with K_i parts in each for a system comprised of N parts. Furthermore, for a given minuscule ϵ , the clustering segmentation is exactly all the parts of the system ($\dim(K_i) = 1$, $c = N$). In the same fashion for a huge ϵ the entire system resides in a single cluster ($\dim(K_0) = N$, $c = 1$). It is clear that the number of clusters describing the system is a monotonically decreasing function of ϵ while the average size of the clusters is a monotonically increasing function of ϵ . We define the two functions as

$$M(\epsilon) = c(\epsilon) \quad (1)$$

for the decreasing function and

$$\overline{S(\epsilon)} = \langle \dim(K_i) \rangle_C = \frac{1}{c(\epsilon)} \sum_{i=0}^{c(\epsilon)} \dim(K_i). \quad (2)$$

for the increasing function.

To deduce an informative clustering cutoff ϵ we form the function G defined as

$$G(\epsilon) = \frac{\overline{S(\epsilon)}M(\epsilon)}{\overline{S(\epsilon)} + M(\epsilon)} - \frac{N}{N + 1} \quad (3)$$

We note that both functions ($\overline{S} \in [1.0 \cdots N]$, $M \in [N \cdots 1.0]$) contain complementary information and that G obtains the maximal value inside the interval. Since the G function is unimodal we can conduct a Golden Ratio Search (GRS) [11] to find the optimal ϵ . Changing the type of S function into any quantile value function, a min or max function should not change the unimodality of the

geometric G function. If we chose to employ a min (S^-) value function then the optimum should be pushed to a larger ϵ value while the max valued (S^+) should in general yield a smaller ϵ solution as compared to the mean. For highly structured data, with an extremely peaked distribution of pairwise distances, these assumptions are not true. One such example is an ideal $2D$ graphene mesh. The ideal graphene bond distances are the same for every connection in the hexagonal mesh. The S heuristic function will be a Heaviside function jumping from 1 to N at a single distance while the M will have the opposite behaviour. This is symmetric for the entire system and will thereby yield a G function response without structure and not exhibit any modality. For systems with a small persisting cluster of only a single constituent until the very last step will result in a unimodal G, min , but the mode will be a global minimum.

To calculate a compositional specificity metric of a cluster to non-binary multiclass labelled targets belonging to analytes in K_i with label counts x , we employ a metric defined by :

$$\gamma = \frac{x^+}{\sum_i x_i} \cdot (1 - e^{-(\dim(x)-1)}) \quad , \quad \dim(x) \geq 1 \quad (4)$$

Another such metric is the τ specificity [13], but we refrain from using it since it is not well defined for clustering solutions where there exist clusters with only one part. Here x^+ corresponds to the maximum value of the label counts in the vector x .

It is clear that a linkage method is more efficient for constructing complete agglomerative hierarchies while a single ‘connectivity’ the calculation might be more efficient if you only want the clusters at a predetermined distance. Searching for an optimal distance cutoff for the cluster representation will also be heuristic dependent. However using this approach, with S and M functions as CCA heuristics for the G metric, finding an optimal ϵ becomes a unimodal optimisation problem.

For completeness, we will include the description of an exhaustive connectivity method (Algorithm 1) in the methods section as well as a description of a GRS (Algorithm 2).

2. Method

For any data set we first construct the distance matrix using a distance measure and data axes, such as a euclidian distance or a spearman correlation distance [24] between all parts of interest in the data set. For molecular water, this is usually the euclidian distance between the atomic positions. For a microarray dataset, it can be a correlation (ρ) distance ($d = \sqrt{1 - \rho}$) between the transcripts sample positions.

In this work, we have chosen to employ euclidian pairwise distance metrics for all the studied data sets.

Given that there exists a segmentation of the system decomposed of c clusters so that the entire system is encoded into K cluster parts in accord with Algorithm 1: CCA¹ then Algorithm 2: GRS² ensures that we always find the optimum ϵ for a unimodal function. We use the heuristics of the CCA given by Equations 1 and 2 which are transformed into a unimodal function by Eq. 3.

We consistently employ an algorithm annotation scheme where subscripts denote index positions in the tensor of interest and where $\dim(A)$ or $\dim(A_i)$ denotes the number of elements in A along its first axis.

We have chosen three data sets. The first is a small molecular water system in the liquid state [1] containing 32 H₂O and a single Hydronium ion. The system is a single time frame of a CPMD [3] water simulation³. The second set contain microarray transcript readings of adipocytes with accession information: *Expression profiling in adipocytes of obese humans* (GSE2508), that employed the GPL8300 platform, describing 20 obese and lean men and women [15]. The third data is the first 35000 Modified National Institute of Standards and Technology (MNIST) digit images [14]⁴.

For all three data sets the full AHC solutions, employing single linkages, were also computed for comparison. The water coordinates were used as is while the two larger data sets were processed by transforming to standardised values, by removing the mean and dividing with the standard deviation across samples or pixels. The microarray and digits were further transformed using Uniform Manifold Approximation and Projection (UMAP) [17] prior to clustering.

The microarray transcripts were also analysed using a two-way ANalysis Of VAriance (ANOVA) [22, 8] modulation for the body type class and biological sex. The significance for body type p-values was employed while sex was considered as a blocking variable. The generated p-values were adjusted using q-value rank correction [25] and used as input to calculate cluster significances. Cluster significances were determined by using a Fisher exact test [7] where all analytes with q-values < 0.05 were considered significant. The data consists of significant analytes A , analytes in a cluster B , insignificant analytes $\neg A$ and analytes not in the cluster $\neg B$. The contingency table (T_{ij}) was populated by the amount of signif-

¹The CCA, Connectivity algorithm have been implemented by the author in the publically available Python package 'impetuous-gfa'

²One such search function has been implemented by the author in the publically available Python 'impetuous-gfa' package in the 'optimisation' module.

³<http://www.theochem.ruhr-uni-bochum.de/legacy.akohlmey/files/32spce-h3op-1ns.xyz>

⁴<http://yann.lecun.com/exdb/mnist/>

Algorithm 1: CCA, Connectivity

input : A Distance matrix D_{ij} , a float ϵ

output: list of cluster sizes Q , cluster towards part index list R

```
1 begin
2    $L \leftarrow \dim(D_{0j})$ 
3    $R, w, P, Q, I$  are empty lists
4    $C_0 \leftarrow 0$ 
5   for  $i \in [0, L)$  do
6      $w_i \leftarrow i + 1$ 
7      $R_{2 \cdot i} \leftarrow 0$ 
8      $R_{2 \cdot i + 1} \leftarrow 0$ 
9      $I_i \leftarrow i$ 
10  end
11  while  $\dim(I) > 0$  do
12     $i \leftarrow I_{\dim(I)-1}$ 
13     $I \leftarrow I_{j \in [0, \dim(I)-1]}$ 
14     $P \leftarrow$  empty list
15    if  $w_i > 0$  then
16       $C_0 \leftarrow C_0 - 1$ 
17      for  $j \in [0, L)$  do
18        if  $D_{ij} \leq \epsilon$  then
19           $P_{\dim(P)} \leftarrow j$ 
20        end
21      end
22      while  $\dim(P) > 0$  do
23         $k \leftarrow P_{\dim(P)-1}$ 
24         $P \leftarrow P_{j \in [0, \dim(P)-1]}$ 
25         $w_k \leftarrow C_0$ 
26        for  $j \in [0, L)$  do
27          if  $D_{ij} \leq \epsilon$  then
28            for  $q \in [0, L)$  do
29              if  $w_q = j + 1$  then
30                 $P_{\dim(P)} \leftarrow q$ 
31              end
32            end
33          end
34        end
35      end
36    end
37  end
```

```

38
39   for  $i \in [0, -1 \cdot C_0)$  do
40     |  $Q_i \leftarrow 0$ 
41   end
42   for  $q \in [0, L)$  do
43     |  $R_{2 \cdot q+1} \leftarrow q$ 
44     |  $R_{2 \cdot q} \leftarrow w_q - C_0$ 
45     |  $Q_{R_{2 \cdot q}} \leftarrow Q_{R_{2 \cdot q}} + 1$ 
46   end
47 end

```

Algorithm 2: GRS

```

input : A Distance matrix  $D_{ij}$ , a float  $\epsilon$  , a float tolerance
output: float  $\epsilon_{optimal}$ 
1 begin
2    $a \leftarrow \min(D_{ij})$ 
3    $b \leftarrow \max(D_{ij})$ 
4    $\psi \leftarrow \frac{\sqrt{5}-1}{2}$ 
5    $c \leftarrow b - \psi \cdot (b - a)$ 
6    $d \leftarrow a + \psi \cdot (b - a)$ 
7   while  $d - c > tolerance$  do
8     |  $fc \leftarrow ( G( heuristics(CCA(D_{ij}, c)) ) )^2$ 
9     |  $fd \leftarrow ( G( heuristics(CCA(D_{ij}, d)) ) )^2$ 
10    | if  $fc \geq fd$  then
11      |  $b \leftarrow d$ 
12      |  $d \leftarrow c$ 
13      |  $c \leftarrow b - \psi \cdot (b - a)$ 
14    | else
15      |  $a \leftarrow c$ 
16      |  $c \leftarrow d$ 
17      |  $d \leftarrow a + \psi \cdot (b - a)$ 
18    | end
19    |  $\epsilon_{optimal} \leftarrow \frac{c+d}{2}$ 
20  end
21 end

```

	G_{max}	c_{max}	G_{mean}	c_{mean}	G_{min}	c_{min}
Water	1.782	17	1.836	11	1.132	33
Pima	0.0630	1313	0.0961	110	0.1688	3
MNIST	0.04713	3388	0.09573	188	0.45622	6

Table 1: Optimal ϵ values as determined via CCA GRS Algorithm 2

icant analytes in the cluster ($T_{00} = \dim(A \cap B)$), significant not in the cluster ($T_{01} = \dim(A \cap \neg B)$), insignificant in the cluster ($T_{10} = \dim(\neg A \cap B)$) and all non significant analytes not in the cluster ($T_{11} = \dim(\neg A \cap \neg B)$). The top cluster transcripts were further analysed using the STRING Database⁵ (string-db) [26] for context.

We did not train a model to infer digits depending on the standardised image data but only describe the cluster content of the optimal solution. The MNIST data clusters were subjected to compositional analysis and benchmarked with the γ metric (Equation 4).

3. Results

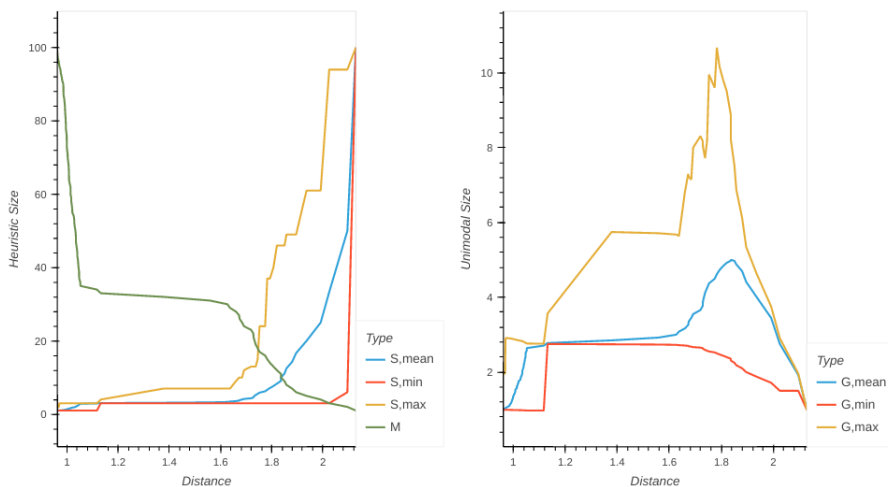
We confirmed that CCA were numerically equivalent to LCA employing single linkage construction. The DBSCAN method without point rejections produced identical clustering solutions with equivalent clustering labels for all distances checked as compared to suggested with (Algorithm 1) and LCA (with single linkage) for the water system.

We summarize the optimisation results from the GRS (Algorithm 2) employing CCA in Table 1.

3.1. The Water coordinates

In Figure 1a we note that the S_{min} heuristic increases early and persists through much of the interval. This causes the G_{min} function in Figure 1b to obtain an early max shifted closer to the smallest distances than the mean. Liquid molecular water is structured and forms hydrogen bonds with its four tetrahedral neighbouring water molecules at distances smaller than 2 [Å]. The atomic hydrogen is bound to the oxygen at distances smaller than 1.1 [Å] and we expect hydrogen reactions between different water molecules to occur in the range $\epsilon \in [1, 2]$ [Å]. We observe that our results, in Table 1, for the optimal clustering solutions are all in this range.

⁵<https://string-db.org/>



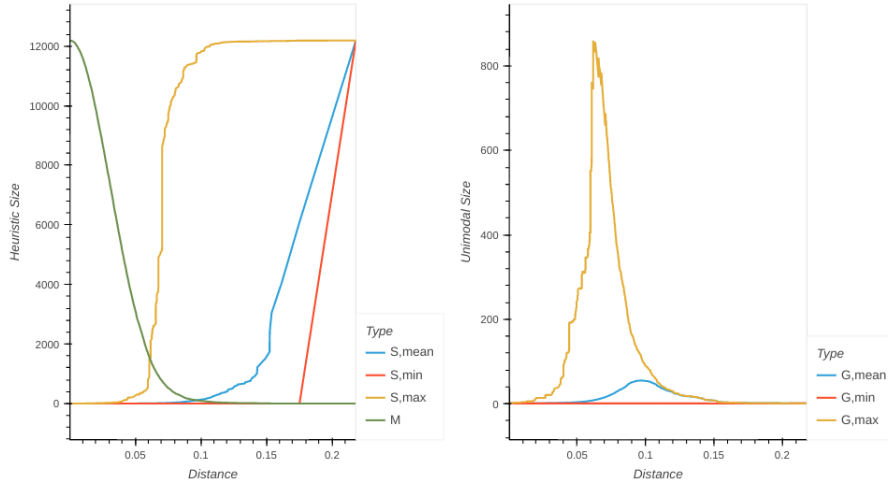
(a) The heuristics described by Eq. 1 and Eq. 2 where the $S, mean$ corresponds to \bar{S} while S, min and S, max corresponds to S^- and S^+ respectively for different ϵ distance values. (b) The G function values for different ϵ distances. The geometric $G, label$ values have been calculated using their corresponding $S, label$ value heuristic.

Figure 1: Overview of all distances in the AHC single linkage solution and their corresponding value functions employed in the clustering optimisation of 100 atoms belonging to water and hydronium molecules.

It is also clear from the M heuristic in, Figure 1a, that going to larger distances would cause the entire system to become connected in a single cluster.

Using the S_{min} metric our unimodal function retains a single water molecule cluster as the heuristic reference through most of the ϵ search range and is the reason for the early jump in the S_{min} . The M heuristic is also flat in a large range starting from the mean O-H position until the first atom species position of the first coordinating water molecule. This causes the G_{min} to obtain its extremum early and clustering at this distance $\epsilon \approx 1.13$ causes the system to decompose into 32 complete water molecules and a single Hydronium molecule.

Both the S_{max} and the S_{mean} solutions obtain similar $\epsilon \approx 1.8$ and form cluster sizes of 17 and 11 segments respectively. One corresponds to the Hydronium centred cluster with 3 and 6 coordinating water molecules respectively as well as several smaller clusters with fewer water molecules. This is the expected outcome since the polarising Hydronium ion will cause a locally denser liquid medium in its first coordination shell [27].



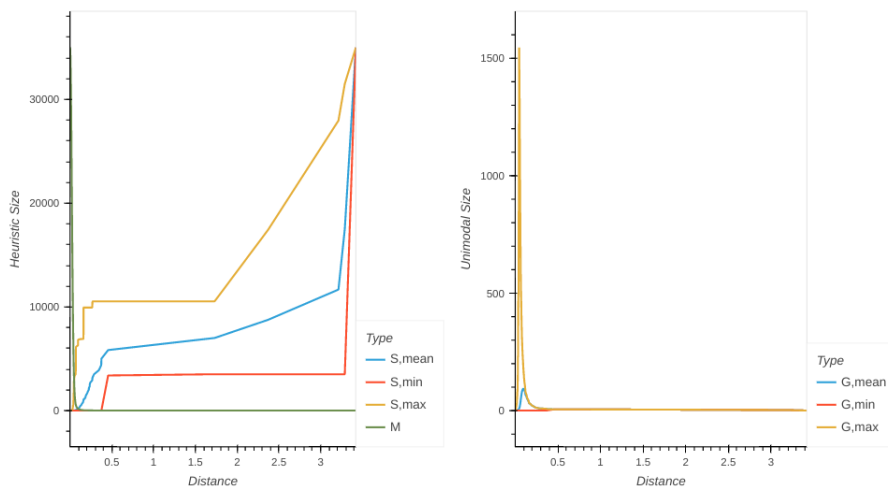
(a) The heuristics described by Eq. 1 and Eq. 2 where the $S, mean$ corresponds to \bar{S} while S, min and S, max corresponds to S^- and S^+ respectively for different ϵ distance values. (b) The G function values for different ϵ distances. The geometric $G, label$ values have been calculated using their corresponding $S, label$ value heuristic.

Figure 2: Overview of all distances in the AHC single linkage solution and their corresponding value functions employed in the clustering optimisation of 12185 gene transcripts belonging to lean and obese men and women.

3.2. The Pima data

The optimisation protocol was evaluated for S^+ , \bar{S} and S^- generating successively fewer c , see Figure 2b. The min optimisation generated a single cluster with 12183 transcripts and two smaller clusters with a single transcript in each. This rendered the decomposition uninformative. The optimisation result was caused by the outlier transcripts causing the persistent single component clusters to form. The fast decay of the total number of clusters causes the G_{min} to obtain a late minimum extremum. The other two solutions, S^+ and \bar{S} both obtain early optima with a larger number of clusters. The \bar{S} solution contains a single huge cluster and several smaller single or few transcript clusters. The S^+ solution contains several larger clusters and many trailing few component clusters. The optimal clustering results are visible in Figure 4.

Since the microarray data has annotated groups, for all the samples, we assessed how well the cluster formation corresponded to traditional ANOVA results. While the clustering metrics by themselves clearly show that the useful solutions correspond to the S^+ and \bar{S} solution. The evaluation of enrichment for significant transcripts explaining the lean-obesity variation showed that the S^+ had the largest amount of significant clusters. The \bar{S} dependent solution results in 1 significant



(a) The heuristics described by Eq. 1 and Eq. 2 where the $S, mean$ corresponds to \bar{S} while S, min and S, max corresponds to S^- and S^+ respectively for different ϵ distance values. (b) The G function values for different ϵ distances. The geometric $G, label$ values have been calculated using their corresponding $S, label$ value function.

Figure 3: Overview of all distances in the AHC single linkage solution and their corresponding value functions employed in the clustering optimisation of 35000 grayscale 28×28 pixel images from the MNIST digits data.

(q -value < 0.05) cluster with a huge amount of transcripts while the S^+ solution correspond to 15 significant clusters with sizes in the range $[10, 400]$. The cluster with the highest significance contains 296 transcripts. The corresponding protein coding activity can be assessed via string-db and relate to: extracellular response to stimulus and organic substances, rheumatoid arthritis [2] as well as abnormal MAPK/ERK pathway signalling [10], which can lead to uncontrolled growth [4]. Changes to cytokine signalling, inflammation and immune system response [9] were also enriched for the cluster. These results are in line with common knowledge of the expressions of obesity [18] and prior knowledge for this data set [15].

3.3. The MNIST digits

The UMAP transformation of the standardised MNIST images obtains a clear structure, see Figure 4c, belonging to the $G, mean$ solution in Figure 3b.

The S^+ dependent solution forms a large number of clusters with high specificity to specific targets, as calculated with Equation 4. The top 20 most specific clusters, exhibiting specificity above 92%, all contain around 100 to 1000 images each but are also redundant in that several digits reappear in several clusters and that the digit 4 is missing. It appears at $\gamma = 0.86$ in a cluster with 237 images.

The downside of this solution is a large number of clusters with $\gamma < 0.5$ which corresponds to 30% of the entire data set.

Both of the solutions corresponding to the choices S^- and \bar{S} are more similar as compared to the S^+ dependent solution. The S^- dependent solution forms 6 clusters, see Table 1. Of which 4 has a specificity of over 96% corresponding to the digits 0, 1, 2, 6. The remaining 2 clusters have specificities of 34% and 35% and are comprised of the digits 4, 7, 9 and 3, 5, 8 respectively. The low specificity of the two larger mixed clusters means that almost 60% of the data segmentation can be assumed to be unreliable.

For the \bar{S} solution the number of clusters is large but there are only 13 clusters with specificity above 50% of which 9 have sizes above 3000 images each. The largest cluster has obtained the lowest specificity of 51% and is dominated by the digits 5, 8. The remaining clusters contain the digits 0, 1, 2, 3, 4, 6, 7, 9 with specificities above 93%. The data set clustering with \bar{S} is more exhaustive in that less than 1% of the data has specificity lower than 50%.

4. Conclusion and Result summary

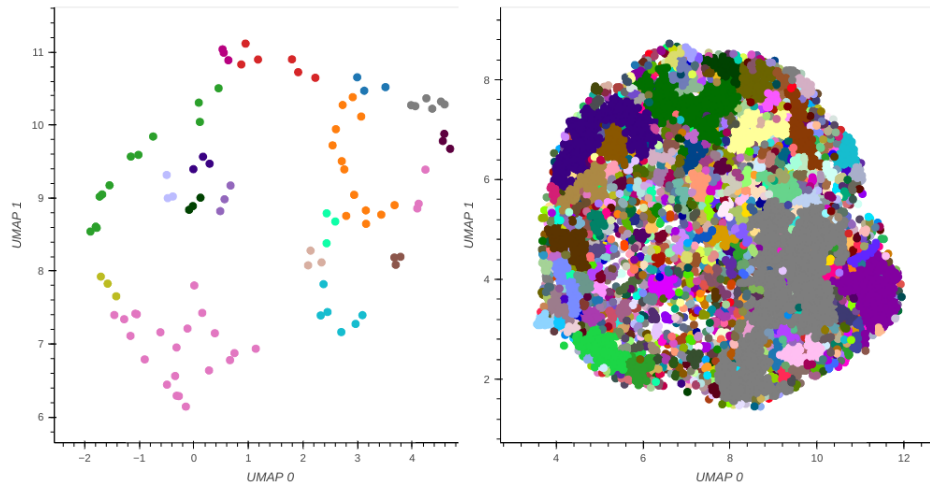
In this work, we evaluated different metrics for unimodal optimisation for determining cluster segments. All the studied metrics resulted in connectivity-based clustering results that corresponded to different distance cuts through a AHC using single linkage. The optimisation strategy determined the global extremum in all the tested cases. The optimisation method suggests a distance cutoff using only the metrics from the distance matrix and the connectivity clustering results.

The choice for the S function value has got a large influence on the optimal solution. For the water system using S^- resulted in segments corresponding to the molecules in the system while the larger \bar{S} and S^+ resulted in chemical environments consisting of several complete molecules.

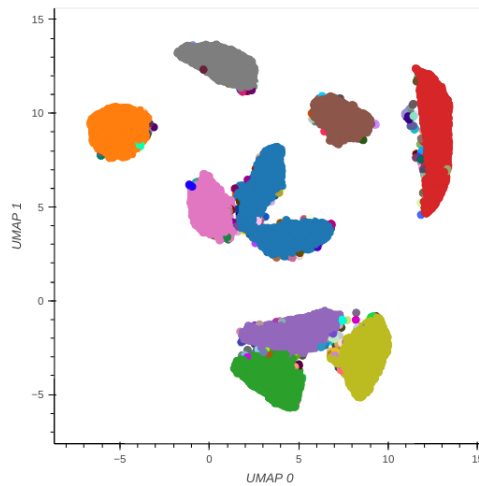
The remaining data sets both had labels corresponding to either feature targets or sample descriptor labels for the MNIST and Pima data sets respectively. This facilitated further evaluation of the usefulness of the clustering solutions. The MNIST data obtained the most informative clustering solution using the \bar{S} heuristic. The Pima data set is the densest and obtained significant and informative clusters by employing the S^+ heuristic.

5. Discussion

For denser data, the S heuristic function probably needs to correspond to larger quantile value functions to yield a useful decomposition. For microarray and



(a) Water optimum solution employing G_{max} calculated with S^+ (b) Pima optimum solution employing G_{max} calculated with S^+



(c) MNIST optimum solution employing G_{mean} calculated with \bar{S}

Figure 4: Overview of the clustering solutions for the G optimum for the three data sets.

Ribonucleic Acid (RNA) sequencing data the S^+ or a high quantile value function is the suggested heuristic for finding a useful optimal ϵ cut. This method will not work well with highly structured data, where the system is entirely connected at one distance, but functions as expected when the data is connected via an ensemble of distances. Then finding a useful cluster representation for visualizing the data becomes a simple optimisation task.

References

- [1] M. P. Allen and D. J. Tildesley. *Computer simulation of liquids*. Clarendon press, Oxford, 1994. ISBN 978-0-19-855645-9.
- [2] M. Ashburner, C. A. Ball, J. A. Blake, D. Botstein, H. Butler, J. M. Cherry, A. P. Davis, K. Dolinski, S. S. Dwight, J. T. Eppig, M. A. Harris, D. P. Hill, L. Issel-Tarver, A. Kasarskis, S. Lewis, J. C. Matese, J. E. Richardson, M. Ringwald, G. M. Rubin, and G. Sherlock. Gene Ontology: tool for the unification of biology. *Nature Genetics*, 25(1):25–29, May 2000. ISSN 1061-4036, 1546-1718. doi: 10.1038/75556.
- [3] R. Car and M. Parrinello. Unified Approach for Molecular Dynamics and Density-Functional Theory. *Physical Review Letters*, 55(22):2471–2474, Nov. 1985. ISSN 0031-9007. doi: 10.1103/PhysRevLett.55.2471.
- [4] J. Downward. Targeting RAS signalling pathways in cancer therapy. *Nature Reviews Cancer*, 3(1):11–22, Jan. 2003. ISSN 1474-175X, 1474-1768. doi: 10.1038/nrc969.
- [5] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, KDD’96, page 226–231. AAAI Press, 1996.
- [6] A. Fernández and S. Gómez. Solving Non-Uniqueness in Agglomerative Hierarchical Clustering Using Multidendrograms. *Journal of Classification*, 25(1):43–65, June 2008. ISSN 0176-4268, 1432-1343. doi: 10.1007/s00357-008-9004-x.
- [7] R. A. Fisher. On the Interpretation of χ^2 from Contingency Tables, and the Calculation of P. *Journal of the Royal Statistical Society*, 85(1):87, Jan. 1922. ISSN 09528385. doi: 10.2307/2340521.
- [8] A. Gelman. Analysis of variance—why it is more important than ever. *The Annals of Statistics*, 33(1), Feb. 2005. ISSN 0090-5364. doi: 10.1214/009053604000001048.
- [9] M. Gillespie, B. Jassal, R. Stephan, M. Milacic, K. Rothfels, A. Senff-Ribeiro, J. Griss, C. Sevilla, L. Matthews, C. Gong, C. Deng, T. Varusai, E. Ragueneau, Y. Haider, B. May, V. Shamovsky, J. Weiser, T. Brunson, N. Sanati, L. Beckman, X. Shao, A. Fabregat, K. Sidiropoulos, J. Murillo, G. Viteri, J. Cook, S. Shorser, G. Bader, E. Demir, C. Sander, R. Haw, G. Wu,

- L. Stein, H. Hermjakob, and P. D'Eustachio. The reactome pathway knowledgebase 2022. *Nucleic Acids Research*, 50(D1):D687–D692, Jan. 2022. ISSN 0305-1048, 1362-4962. doi: 10.1093/nar/gkab1028.
- [10] M. Kanehisa, M. Furumichi, Y. Sato, M. Ishiguro-Watanabe, and M. Tanabe. KEGG: integrating viruses and cellular organisms. *Nucleic Acids Research*, 49(D1):D545–D551, Jan. 2021. ISSN 0305-1048, 1362-4962. doi: 10.1093/nar/gkaa970.
- [11] J. Kiefer. Sequential minimax search for a maximum. *Proceedings of the American Mathematical Society*, 4(3):502–506, 1953. ISSN 0002-9939, 1088-6826. doi: 10.1090/S0002-9939-1953-0055639-3.
- [12] V. Y. Kiselev, K. Kirschner, M. T. Schaub, T. Andrews, A. Yiu, T. Chandra, K. N. Natarajan, W. Reik, M. Barahona, A. R. Green, and M. Hemberg. SC3: consensus clustering of single-cell RNA-seq data. *Nature Methods*, 14(5):483–486, May 2017. ISSN 1548-7091, 1548-7105. doi: 10.1038/nmeth.4236.
- [13] N. Kryuchkova-Mostacci and M. Robinson-Rechavi. A benchmark of gene expression tissue-specificity metrics. *Briefings in Bioinformatics*, page bbw008, Feb. 2016. ISSN 1467-5463, 1477-4054. doi: 10.1093/bib/bbw008.
- [14] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Backpropagation Applied to Handwritten Zip Code Recognition. *Neural Computation*, 1(4):541–551, Dec. 1989. ISSN 0899-7667, 1530-888X. doi: 10.1162/neco.1989.1.4.541.
- [15] Y. H. Lee, S. Nair, E. Rousseau, D. B. Allison, G. P. Page, P. A. Tataranni, C. Bogardus, and P. A. Permana. Microarray profiling of isolated abdominal subcutaneous adipocytes from obese vs non-obese Pima Indians: increased expression of inflammation-related genes. *Diabetologia*, 48(9):1776–1783, Sept. 2005. ISSN 0012-186X. doi: 10.1007/s00125-005-1867-3.
- [16] S. Liu, A. Thennavan, J. P. Garay, J. S. Marron, and C. M. Perou. MultiK: an automated tool to determine optimal cluster numbers in single-cell RNA sequencing data. *Genome Biology*, 22(1):232, Dec. 2021. ISSN 1474-760X. doi: 10.1186/s13059-021-02445-5.
- [17] L. McInnes, J. Healy, and J. Melville. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. 2018. doi: 10.48550/ARXIV.1802.03426. Publisher: arXiv Version Number: 3.

- [18] R. Monteiro and I. Azevedo. Chronic Inflammation in Obesity and the Metabolic Syndrome. *Mediators of Inflammation*, 2010:1–10, 2010. ISSN 0962-9351, 1466-1861. doi: 10.1155/2010/289645.
- [19] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, pages 2825–2830, 2011.
- [20] M. Ronen, S. E. Finder, and O. Freifeld. DeepDPM: Deep Clustering With an Unknown Number of Clusters. Technical Report arXiv:2203.14309, arXiv, Mar. 2022. arXiv:2203.14309 [cs, stat] type: article.
- [21] R. Satija, J. A. Farrell, D. Gennert, A. F. Schier, and A. Regev. Spatial reconstruction of single-cell gene expression data. *Nature Biotechnology*, 33(5): 495–502, May 2015. ISSN 1087-0156, 1546-1696. doi: 10.1038/nbt.3192.
- [22] S. Seabold and J. Perktold. statsmodels: Econometric and statistical modeling with python. In *9th Python in Science Conference*, 2010.
- [23] R. Sibson. SLINK: An optimally efficient algorithm for the single-link cluster method. *The Computer Journal*, 16(1):30–34, Jan. 1973. ISSN 0010-4620, 1460-2067. doi: 10.1093/comjnl/16.1.30.
- [24] V. Solo. Pearson Distance is not a Distance. Technical Report arXiv:1908.06029, arXiv, Aug. 2019. arXiv:1908.06029 [stat] type: article.
- [25] J. D. Storey and R. Tibshirani. Statistical significance for genomewide studies. *Proceedings of the National Academy of Sciences*, 100(16):9440–9445, Aug. 2003. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.1530509100.
- [26] D. Szklarczyk, A. L. Gable, K. C. Nastou, D. Lyon, R. Kirsch, S. Pyysalo, N. T. Doncheva, M. Legeay, T. Fang, P. Bork, L. J. Jensen, and C. von Mering. The STRING database in 2021: customizable protein–protein networks, and functional characterization of user-uploaded gene/measurement sets. *Nucleic Acids Research*, 49(D1):D605–D612, Jan. 2021. ISSN 0305-1048, 1362-4962. doi: 10.1093/nar/gkaa1074.
- [27] R. Tjörnhammar and O. Edholm. Molecular dynamics simulations of Zn²⁺ coordination in protein binding sites. *The Journal of Chemical Physics*, 132(20):205101, May 2010. ISSN 0021-9606, 1089-7690. doi: 10.1063/1.3428381.

- [28] J. H. Ward. Hierarchical Grouping to Optimize an Objective Function. *Journal of the American Statistical Association*, 58(301):236–244, Mar. 1963. ISSN 0162-1459, 1537-274X. doi: 10.1080/01621459.1963.10500845.