

DBSpan: Density-Based Spanner for Clustering Complex Data, With an Application to Persistence Diagrams*

Brittany Terese Fasy[†], David L. Millman[†], Elliott Pryor[†], and Nathan Stouffer[‡]

Abstract. Since its introduction in the mid-1990s, DBSCAN has become one of the most widely used clustering algorithms. However, one of the steps in DBSCAN is to perform a range query, a task that is difficult in many spaces, including the space of persistence diagrams. In this paper, we introduce a spanner into the DBSCAN algorithm to facilitate range queries in such spaces. We provide a proof-of-concept implementation, and study time and clustering performance for two data sets of persistence diagrams.

Key words. clustering, distance approximation, DBSCAN, persistent homology, spanners

AMS subject classifications. 55N31 Persistent homology and applications, topological data analysis, 62H30 Classification and discrimination; cluster analysis (statistical aspects), 62R40 Topological data analysis, and 68T09 Computational aspects of data analysis and big data.

1. Introduction. A common task in data analysis is clustering. Since its introduction in the mid-1990s, DBSCAN [16] has become one of the most widely used clustering algorithms. Range queries—finding all data points that satisfy a certain (usually geometric) property, such as all points within distance ϵ of a point x —are a central step in the DBSCAN algorithm. However, computing range queries and pairwise distance matrices is a task that is often cumbersome in high-dimensions and in non-Euclidean spaces, such as spaces of: persistence diagrams, Reeb graphs, point clouds, and curves in space.

In this paper, we are interested in studying the clustering of persistence diagrams. The space of persistence diagrams under the bottleneck distance has infinite doubling dimension [17], making it not suitable for many clustering algorithms. To get around this issue, approximating distances (and defining new distances) for the task of clustering persistence diagrams has been the focus of several recent papers in topological data analysis (TDA) [8, 13, 14, 17, 32]. We present DBSpan, a modified version of the DBSCAN algorithm, that replaces the range query with an approximate range query by using a spanner technique introduced by Kerber and Nigmatov [25] for approximating expensive distances. When used to cluster persistence diagrams, we demonstrate performance improvements over the standard DBSCAN algorithm.

This paper begins, in Section 2, by introducing preliminaries needed to understand our algorithm, then presents the algorithm itself in Section 3. Furthermore, we provide a proof-of-concept implementation in Python¹ and demonstrate the utility of the algorithm with experimental results in the space of persistence diagrams under the bottleneck distance in Section 4.

2. Preliminaries. In this section, we introduce notation and definitions needed to understand the contributions of this short paper. Relevant references are provided where the interested reader can find more details on these preliminaries.

Notation. S is a finite metric space, with corresponding distance metric $d: S \times S \rightarrow \mathbb{R}$. For a point $x \in S$ and radius $\epsilon > 0$, we use $N_\epsilon(x)$ to denote the set of points in $S \setminus \{x\}$ within the closed ϵ -neighborhood or ϵ -neighborhood of x ; that is, $N_\epsilon(x) := \{y \neq x \in S \mid d(x, y) \leq \epsilon\}$.

Clustering (and Other Tasks) with Persistence Diagrams. While the algorithm presented in this paper is applicable to many types of data, we are particularly motivated by using persistence diagrams to compare data, and the difficulties that arise in clustering data using persistence diagrams.

*Submitted to the editors 8 March 2022.

Funding: BTF was partially supported by NSF grants DMS 1664858, DMS 1854336 and CCF 2046730.

[†]Montana State University, Bozeman, MT (brittany.fasy@montana.edu, david.millman@montana.edu, elliott.pryor@student.montana.edu).

[‡]onXmaps, Inc., Bozeman, MT

¹The implementation is publicly available via MIT license at: <https://github.com/compTAG/dbspan>

43 Geometrically, a persistence diagram is a multiset of points in the extended plane (the extended
 44 plan $\overline{\mathbb{R}^2}$ allows infinite coordinates). To compare two persistence diagrams D_1 and D_2 , we use the
 45 bottleneck distance between them, denoted $d_\infty(D_1, D_2)$. We defer formal definitions of persistence
 46 diagrams and distances between them to [15].

47 The difficulty of dealing with persistence diagrams in tasks such as clustering or range queries
 48 largely boils down to two issues. First, the computation of the distance itself is expensive. To
 49 compute $d_\infty(D_1, D_2)$, one must minimize over the set of all matchings between point sets D_1
 50 and D_2 . While heuristic improvements on computation have been made to prune this search space
 51 in practice [23], it is still a time-consuming task. Thus, for large data sets, computing all pairwise
 52 distances is impractical. Second, the doubling dimension (that is, the log of the number of balls
 53 of radius $\frac{r}{2}$ needed to cover a ball of radius r ; see [11]) of the space of persistence diagrams is
 54 infinite [17, 24, 25]. Often, algorithms for computing nearest neighbors or range queries assume low
 55 doubling dimension. And, even in spaces of constant doubling dimension d , to use cover trees for
 56 computing nearest neighbors is $O(d^{12} \log n)$ for a data set with n points [2].

57 **DBSCAN.** The Density-Based Spatial Clustering of Applications with Noise (DBSCAN) Algo-
 58 rithm was first introduced in [16] and has seen many improvements and implementations [4, 6, 7,
 59 20, 27]. (Most of the implementations are optimized for datasets lying in \mathbb{R}^d). The intuition of
 60 this algorithm is that, given $\epsilon > 0$, we: (1) for each point x , compute $N_\epsilon(x)$, the set of all points
 61 whose distance is at most ϵ from x ; (2) construct a neighborhood graph on these neighborhoods,
 62 restricting to the “dense” points (i.e., the vertices of this graph are the points whose ϵ -neighborhoods
 63 meet a minimum threshold of neighbors, and edges exist between vertices if their distance is at
 64 most 2ϵ). The output of DBSCAN is a labeling of each vertex with its cluster, or as ‘noise’ (that is,
 65 non-dense points whose ϵ -neighborhoods do not contain any dense point). The vertices of the con-
 66 nected components of this graph (along with additional points “at the boundary” of these clusters)
 67 are the clusters that DBSCAN computes.

68 The central step of DBSCAN is the computation of the ϵ -neighborhoods. This step is an
 69 example of a *range query*: given a shape in the domain (such as a disc in \mathbb{R}^2), find all data points
 70 that are contained in that shape. In \mathbb{R}^2 , a data structure can be created to support range queries
 71 of ℓ_∞ balls with $O(n \log n)$ precomputation such that a range query whose range has k points
 72 will take $O(\log n + k)$ distance calculations [38]. Assuming the ϵ -neighborhoods in DBSCAN have
 73 constant size, we observe that the computation of all ϵ neighborhoods takes $O(n \log n)$. In spaces
 74 with low doubling dimension, rather than calculating the ϵ -neighborhoods exactly, approximating
 75 these using approximate nearest neighbor graphs can be used. For example, [39] replaces the range
 76 queries in DBSCAN with locality sensitive hashing. And, a paper posted to ArXiv in 2020 [9]
 77 uses approximate k -nn graphs using random projections to provide a speedy approximation of
 78 the ϵ -neighborhoods in DBSCAN. However, in spaces with large doubling dimension, a single range
 79 query (in practice) resorts to scanning through all data points; hence, the computation of the ϵ -
 80 neighborhood requires $O(n^2)$ distance calculations.

81 DBSCAN is familiar to the TDA community, as papers have been published using DBSCAN.
 82 Within this context, DBSCAN has been used to cluster points within a persistence diagram [26, 29],
 83 as a comparison against TDA-based techniques for clustering [12, 21, 28, 31], and as a clustering
 84 subroutine in a TDA pipeline [3]. In the latter, a pairwise distance matrix is required as input,
 85 which makes using DBSCAN prohibitive for large data sets of persistence diagrams.

86 **Spanners.** In DBSCAN, when spatial range queries or other ways to speed-up the calculation
 87 of ϵ -neighborhoods are not possible, we rely on the pairwise distance matrix; that is, we need to
 88 store (or at least compute) the distance between every pair of input data. Thinking of this matrix
 89 as a large graph (a vertex for each datum and an edge for each pairwise distance), spanners help
 90 to simplify the graph, while maintaining approximate distances. In particular, a $(1 + \delta)$ -spanner of
 91 a discrete metric space (S, d) is a weighted graph $G = (S, E, \omega)$ such that for all pairs $p, q \in S$,

$$92 \quad (2.1) \quad d(p, q) \leq d_G(p, q) \leq (1 + \delta)d(p, q),$$

93 where $d_G(p, q)$ is the length of the shortest path between p and q in G . Spanners were originally
 94 introduced to simplify the pairwise distance matrix for points in the plane [10], but have since been
 95 used in other metric spaces as well [1, 19, 25]. However, the underlying assumption is often that
 96 distances are fast to compute and that we can have access to the entire pairwise distance matrix
 97 to compute the spanner.

98 On the other hand, Kerber and Nigmatov [25] studied spanners in metric spaces where the
 99 distances are expensive to compute, such as in the space of persistence diagrams. Given a param-
 100 eter δ , they provide algorithms to compute a $(1 + \delta)$ -spanner for a pairwise distance graph. The
 101 spanners do not require computing all $|S|^2$ distances, and in practice, tend to be linear in size.
 102 While we defer the details of the algorithm to [25], we note that their algorithm is an iterative
 103 method that maintains, for each $p, q \in S$, a lower and an upper bound, $a(p, q)$ and $b(p, q)$ for the
 104 distance $d(p, q)$. When selecting which distance to compute next (that is, which edge to add to the
 105 spanner), they offer two heuristics, BLINDGREEDY and BLINDRANDOM. Both of these heuristics
 106 select which edge based on the ratio $\frac{b(p,q)}{a(p,q)}$, with BLINDGREEDY selecting the edge that maximizes
 107 that ratio and BLINDRANDOM selecting randomly among edges whose ratio is at least $1 + \delta$.

108 **3. Algorithm.** We present DBSpan, an approximation of DBSCAN for metric spaces where
 109 the distances are expensive to compute. The algorithm, given in Algorithm 3.1, mirrors DBSCAN,
 110 but differs in two key places:

- 111 1. We replace the ϵ -neighborhood $N_\epsilon(p)$ with an approximate nearest-neighbor ball $\tilde{N}_\epsilon^\delta(p)$
 112 computed using distances in the spanner graph from [25] using the BLINDRANDOM heuristic.
- 113 2. Rather than computing the ϵ -neighborhoods on the fly, the approximate ϵ -neighborhoods
 114 from the spanner are precomputed in Line 3.

115 Note that the ϵ -neighborhoods for DBSCAN could also be pre-computed; however, all $\Theta(|S|^2)$
 116 distances between input data will still need to be computed first. By using the spanner from [25],
 117 DBSpan reduces the number of pairwise distances that need to be computed.

118 **Implementation.** Along with this paper, we are releasing an open-source Python implementation
 119 (see link in Footnote 1). This code uses Scikit-TDA [34] for computing persistence diagrams on Rips
 120 filtrations and computing the bottleneck distance between persistence diagrams. For the most part,
 121 the Python code mirrors Algorithm 3.1. We note that, for clarity of exposition and for adhering
 122 to variable naming practices in Python, several variables have different names in the pseudocode
 123 in this paper and the Python code: S in the pseudocode is called `data` in the Python code, ϵ is
 124 called `eps`, m is called `min_samples`, and Q is called `queue`. In addition, the Python code has one
 125 performance improvement: rather than pushing p onto the queue in Line 13, p is directly processed
 126 and Q is initialized to $N_\epsilon(p)$. This saves one range query.

127 **Relationship between Clusters Found by DBSCAN and DBSpan.** By Equation (2.1), several rela-
 128 tionships between the clusters found by DBSCAN and DBSpan arise that allow us to ensure that
 129 ‘nearby’ points are clustered together and ‘far away’ points are not. In particular:

- 130 1. If $p, q \in S$ such that $\text{DBSpan}(S, \epsilon, m, \delta)$ clusters them into the same cluster and their ϵ -
 131 neighborhoods have at least m points, then p and q will be in the same cluster using
 132 $\text{DBSCAN}(S, \epsilon, M)$ for all $M \geq m$.
- 133 2. If $p, q \in S$ such that $\text{DBSpan}(S, \epsilon, m, \delta)$ clusters them into different clusters, C_p and C_q , if
 134 there are no outliers, and if the minimum inter-cluster distance between all pairs of clusters
 135 is greater than $(1 + \delta)\epsilon$, then p and q will not be in the same cluster using $\text{DBSCAN}(S, \epsilon, m)$.

136 These statements follow from the fact that the approximate ϵ -neighborhoods are computed from
 137 an $(1 + \delta)$ -spanner over S .

138 4. Experimental Results.

139 **Experimental Methods.** In the experiments, we compare clustering using the scikit-learn [30]
 140 implementation of DBSCAN to the clustering obtained using DBSpan (Algorithm 3.1). To validate
 141 the clustering, we treat the clusters from DBSCAN as the *ground truth* and compute the adjusted
 142 Rand index (ARI) [37] for DBSpan. Similar to the standard Rand index [33], if the clusters match

Algorithm 3.1 DBSpan(S, ϵ, m, δ)**Input:** Dataset S ; neighborhood radius $\epsilon \geq 0$; core size $m \geq 0$; and spanner parameter δ **Output:** A cluster label is attributed to each point in S

```

1: For all  $p \in S$ , initialize  $p.cluster = -1$  {All points are considered unlabeled initially}
2: Compute spanner of  $S$ 
3: For all  $p \in S$ , compute  $\tilde{N}_\epsilon^\delta(p)$  using the spanner
4:  $c = 0$  {number of clusters; also current cluster number}
5: for  $p \in S$  do
6:   if  $p$  is labeled then
7:     continue
8:   else if  $|\tilde{N}_\epsilon^\delta(p)| < m$  then
9:      $p.cluster \leftarrow 0$  { $p$  is labeled as noise}
10:    continue
11:  end if
12:   $c \leftarrow c + 1$ 
13:   $Q.push(p)$  {initialize queue of points to be added to current cluster}
14:  while  $Q$  is not empty do
15:     $q \leftarrow Q.pop()$ 
16:    if  $q.cluster \geq 1$  then
17:      continue
18:    end if
19:     $q.cluster \leftarrow c$  { $q$  is now considered labeled}
20:    if  $|\tilde{N}_\epsilon^\delta(q)| \geq m$  then
21:      Add  $N_\epsilon(q)$  to  $Q$ 
22:    end if
23:  end while
24: end for
25: return

```

143 exactly between DBSCAN and DBSpan (up to relabeling), then the adjusted index is one. The
 144 more disagreement between the two clusterings, the closer the index is to zero.

145 *Experiment 1: Shape Dataset.* In this experiment, we study how in the shape dataset, varying
 146 the approximation of the true distance, δ , changes the quality of the clustering (as measured by ARI).

147 We consider a simple data set of shapes in \mathbb{R}^4 . Using TaDAsets from Scikit-TDA [34], we
 148 create 74 shapes from three different classes: 30 torii, 30 spheres, and 14 Swiss rolls. For each
 149 shape, we sample 100 points from the shape, apply Gaussian noise with $\sigma = .1$, and compute the
 150 persistence diagram corresponding to the Vietoris-Rips filtration. Then, the 1d diagram is used as
 151 input into DBSCAN and DBSpan with common parameters $\epsilon = 0.3$ and $m = 15$.

152 In Table 1, we see the result of varying the spanner parameter δ and comparing the speed and
 153 accuracy of DBSpan. The first column, δ , is the approximation for the spanner. The second column
 154 is the ARI of the DBSpan output as compared to the DBSCAN ground truth. The third column is
 155 the number of edges in the resulting spanner. The fourth column is the proportion of edges in the
 156 spanner to the complete graph. Note that the complete graph has 2701 edges. The fifth column
 157 is the time (in seconds) for running DBSpan. The sixth column is the speedup of DBSpan. Note
 158 that the DBSCAN took 74.99s to complete.

159 Observe that for $\delta = \{.1, 1\}$, DBSpan and DBSCAN have a ARI index of 1, which means
 160 that they produced the same output. Unfortunately, when $\delta = .1$, DBSpan is slightly slower than
 161 DBSCAN. In this case, even though we only compute 83% of all distances, the cost of determining
 162 which distances can be omitted overwhelms the gains (as could be expected). For this problem, $\delta =$
 163 10 strikes a nice balance between speed and accuracy as Rand index is .98 with a 4x speedup.

Table 1: Performance metrics comparing DBSCAN and varying δ values for DBSpan. For this experiment, DBSCAN ran in 74.99 seconds and the complete graph contains 2701 edges. We compare accuracy between the clusterings with the adjusted Rand index (ARI) in column 2 and measure DBSpan time in column 5 in seconds.

δ	ARI	Num Edges	% Possible Edges	DBSpan time	Speedup
0.1	1.00	2243	83.0	80.91	0.9
1.0	1.00	1114	41.2	48.00	1.6
10.0	0.98	442	16.4	18.74	4.0
50.0	0.84	398	14.7	16.58	4.5
100.0	0.86	375	13.9	15.66	4.8
500.0	0.95	421	15.6	18.62	4.0
1000.0	0.84	388	14.4	17.44	4.3

Table 2: Performance metrics comparing DBSCAN and DBSpan on varying size inputs. We compare accuracy between the clusterings with the adjusted Rand index (ARI) in column 2. The time to run DBSCAN and DBSpan are measured in seconds and reported in columns 4 and 5, respectively.

Num Dgms	ARI	% Possible Edges	DBSCAN time	DBSpan time	Speedup
30	1.00	66.5	283.25	185.14	1.5
45	1.00	66.3	654.45	429.65	1.5
60	0.79	64.5	1007.53	656.37	1.5
75	0.94	55.9	1762.89	998.35	1.8
90	0.90	56.2	2435.76	1369.58	1.8
105	0.96	54.3	3476.13	1942.88	1.8
120	0.86	49.3	4457.90	2269.63	2.0

164 For $\delta > 10$, however, the sparsification (and along with it the loss of accuracy and speed up) seems
 165 to plateau at computing $15\% \pm 1.5\%$ of the possible edges. We have seen this plateau in other tests
 166 cases as well and would like to investigate further.

167 *Experiment 2: Prostate Cancer Dataset.* In this experiment, we study the scaling of DBSpan. For
 168 a fixed distance approximation, we compare various size datasets. As in the previous experiment,
 169 we treat the clusters from DBSCAN as ground truth and assess accuracy with ARI.

170 We consider a subset of images from the Kaggle Prostate cANcer graDe Assessment (PANDA)
 171 Challenge [5], partitioned into small 512×512 pixel regions of interest (ROIs). For each ROI, we
 172 use the Histocartography Python library [22] to extract nuclei centroids, and we use GUDHI [36]
 173 to compute a persistence diagram using the Vietoris-Rips filtration on these centroids. So that
 174 we consider interesting diagrams, we remove any diagram with less than 100 1d-persistence points.
 175 The resulting set of 8595 1d-diagrams are available on the Open Science Framework [18] in the
 176 DBSpan Test Dataset project [35]. The distribution of the number of points in these diagrams are
 177 imbalanced, so when we sample from this set, we create three levels: diagrams with less than 150
 178 points, diagrams with between 150 and 200 points, and diagrams with more than 200 points; and
 179 perform disproportionate stratified random sampling. Then, the sampled diagrams are used as
 180 input into DBSCAN and DBSpan with common parameters $\epsilon = 10$ and $m = 5$.

181 Then, we pick a reasonable delta value for this experiment by sampling 30 diagrams from our
 182 dataset and identifying a δ in which the ARI is still close to 1. For this dataset, at $\delta = 1$, we have
 183 a ARI of .867 while only computing about 64% of the pairwise distance matrix.

184 In Table 2, we see the result of varying the number of diagrams and compare the speed and
 185 accuracy of DBSpan to DBSCAN. The first column is the number of diagrams to cluster. The
 186 second column is the ARI of the DBSpan output as compared to the DBSCAN ground truth. The
 187 third column is the proportion of edges in the spanner to the complete graph. The fourth column is
 188 the time (in seconds) for running DBSCAN. The fifth column is the time (in seconds) for running
 189 DBSpan. The sixth column is the speedup of DBSpan.

190 Observe that as the input gets larger we increase the speedup from 1.5 to 2. As one would
 191 expect, speed seems to increase with the reduction of the number of edges. Perhaps, unexpectedly,
 192 the proportion of edges computed decreases as the as we increase the number of nodes. For all but
 193 two trials, the ARI is above .9. Note that from the previous experiment, one would not expect the
 194 ARI to change substantially for a fixed δ . It seems that further study is required to understand
 195 how to pick a more predictable δ for clustering.

196 **5. Conclusions.** In this paper, we introduced a spanner into the DBSCAN algorithm to facil-
 197 itate range queries in spaces with expensive distance computations. In the practical side, we saw
 198 that even when we remove the computation of a small fraction of distances, we see improvements
 199 in performance with little deviation from the output of DBSCAN. We believe that we could see
 200 an even larger speedups by optimizing the implementation and by parallelizing the construction of
 201 the spanner. Kerber and Nigmatov [25] proposed multiple heuristics for building up the spanner
 202 without computing all distances. In this paper, we only considered one method. Are there better
 203 ways to add edges to the spanner when the task is approximating neighborhoods? Finally, our
 204 process for picking δ is done in a somewhat ad-hock fashion. It would be useful to develop a more
 205 theoretically guided method.

206 **Acknowledgments.** We thank users of the R package TDA² for questions on how to improve
 207 pairwise distance computations that inspired us to work in this direction. We thank Jordan Schup-
 208 bach for pre-computing persistence diagrams, and Ben Holmgren and Braeden Sopp for their
 209 thoughtful feedback. We also thank Michael Kerber for conversations after we read [25], a key
 210 result needed for DBSpan.

211 REFERENCES

- 212 [1] I. ALTHÖFER, G. DAS, D. DOBKIN, D. JOSEPH, AND J. SOARES, *On sparse spanners of weighted graphs*,
 213 *Discrete & Computational Geometry*, 9 (1993), pp. 81–100.
- 214 [2] A. BEYGEZIMMER, S. KAKADE, AND J. LANGFORD, *Cover trees for nearest neighbor*, in Proceedings of the 23rd
 215 international conference on Machine learning, 2006, pp. 97–104.
- 216 [3] S. BIASOTTI, A. CERRI, D. GIORGI, AND M. SPAGNUOLO, *PHOG: Photometric and geometric functions for*
 217 *textured shape tetrieval*, in Computer Graphics Forum, vol. 32, Wiley Online Library, 2013, pp. 13–22.
- 218 [4] D. BIRANT AND A. KUT, *ST-DBSCAN: An algorithm for clustering spatial-temporal data*, *Data & knowledge*
 219 *engineering*, 60 (2007), pp. 208–221.
- 220 [5] W. BULTEN, K. KARTASALO, P.-H. C. CHEN, P. STRÖM, H. PINCKAERS, K. NAGPAL, Y. CAI, D. F. STEINER,
 221 H. VAN BOVEN, R. VINK, ET AL., *Artificial intelligence for diagnosis and gleason grading of prostate cancer:*
 222 *the panda challenge*, *Nature medicine*, (2022), pp. 1–10.
- 223 [6] R. J. CAMPELLO, D. MOULAVI, AND J. SANDER, *Density-based clustering based on hierarchical density esti-*
 224 *mates*, in Pacific-Asia conference on knowledge discovery and data mining, Springer, 2013, pp. 160–172.
- 225 [7] R. J. CAMPELLO, D. MOULAVI, A. ZIMEK, AND J. SANDER, *Hierarchical density estimates for data clustering,*
 226 *visualization, and outlier detection*, *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 10
 227 (2015), pp. 1–51.
- 228 [8] S. CHEN AND Y. WANG, *Approximation algorithms for 1-Wasserstein distance between persistence diagrams*,
 229 in 19th International Symposium on Experimental Algorithms (SEA), 2021.
- 230 [9] Y. CHEN, W. RUYS, AND G. BIROS, *KNN-DBSCAN: A DBSCAN in high dimensions*, arXiv preprint
 231 arXiv:2009.04552, (2020).
- 232 [10] L. P. CHEW, *There is a planar graph almost as 600d as the complete graph*, 2 (1986), pp. 169–177.
- 233 [11] K. L. CLARKSON, *Nearest-neighbor searching and metric space dimensions*, in Nearest-Neighbor Methods for
 234 Learning and Vision: Theory and Practice, G. Shakhnarovich, T. Darrell, and P. Indyk, eds., MIT Press,
 235 2006, p. 1559.
- 236 [12] R. COTSAKIS, J. SHAW, J. TIERNY, AND J. A. LEVINE, *Implementing persistence-based clustering of point*
 237 *clouds in the Topology ToolKit*, in Topological Methods in Data Analysis and Visualization VI, Springer,
 238 2021, pp. 343–357.
- 239 [13] T. K. DEY AND S. ZHANG, *Approximating 1-Wasserstein distance between persistence diagrams by graph spar-*
 240 *sification*, in 2022 Proceedings of the Symposium on Algorithm Engineering and Experiments (ALENEX),
 241 SIAM, 2022, pp. 169–183.

²<https://cran.r-project.org/web/packages/TDA/vignettes/article.pdf>

- 242 [14] P. DŁOTKO AND N. HELLMER, *Bottleneck profiles and discrete prokhorov metrics for persistence diagrams*,
243 (2021). arXiv preprint arXiv:2106.02538.
- 244 [15] H. EDELSBRUNNER AND J. HARER, *Computational Topology: An Introduction*, American Mathematical Society,
245 2010.
- 246 [16] M. ESTER, H.-P. KRIEGEL, J. SANDER, X. XU, ET AL., *A density-based algorithm for discovering clusters in*
247 *large spatial databases with noise*, in kdd, vol. 96, 1996, pp. 226–231.
- 248 [17] B. T. FASY, X. HE, Z. LIU, S. MICKA, D. L. MILLMAN, AND B. ZHU, *Approximate nearest neighbors in the*
249 *space of persistence diagrams*, 2018. arXiv preprint arXiv:1812.11257.
- 250 [18] E. D. FOSTER AND A. DEARDORFF, *Open science framework (OSF)*, Journal of the Medical Library Association,
251 105 (2017), pp. 203–206.
- 252 [19] L.-A. GOTTLIEB AND L. RODITTY, *An optimal dynamic spanner for doubling metric spaces*, in European
253 Symposium on Algorithms, Springer, 2008, pp. 478–489.
- 254 [20] M. HAHLER, M. PIEKENBROCK, AND D. DORAN, *dbscan: Fast density-based clustering with R*, Journal of
255 Statistical Software, 91 (2019), pp. 1–30.
- 256 [21] U. ISLAMBEKOV AND Y. R. GEL, *Unsupervised space–time clustering using persistent homology*, Environmetrics,
257 30 (2019), p. e2539.
- 258 [22] G. JAUME, P. PATI, V. ANKLIN, A. FONCUBIERTA, AND M. GABRANI, *Histocartography: A toolkit for graph*
259 *analytics in digital pathology*, in MICCAI Workshop on Computational Pathology, 2021, pp. 117–128.
- 260 [23] M. KERBER, D. MOROZOV, AND A. NIGMETOV, *Geometry helps to compare persistence diagrams*, Journal of
261 Experimental Algorithmics (JEA), 22 (2017), pp. 1–4.
- 262 [24] M. KERBER AND A. NIGMETOV, *Spanners for topological summaries*, June 2018. CG Week Young Researcher’s
263 Forum.
- 264 [25] M. KERBER AND A. NIGMETOV, *Metric spaces with expensive distances*, International Journal of Computational
265 Geometry & Applications, 30 (2020), pp. 141–165.
- 266 [26] A. LAWSON, T. HOFFMAN, Y.-M. CHUNG, K. KEEGAN, AND S. DAY, *A density-based approach to feature*
267 *detection in persistence diagrams for firn data*, Foundations of Data Science, (2021).
- 268 [27] L. MCINNES, J. HEALY, AND S. ASTELS, *HDBSCAN: Hierarchical density based clustering*, Journal of Open
269 Source Software, 2 (2017), p. 205.
- 270 [28] M. MICHEL, J. LE DEUNF, N. DEBESE, L. BAZINET, AND L. DEJOIE, *Multibeam outlier detection by clustering*
271 *and topological persistence approach, ToMATo algorithm*, in OCEANS 2021: San Diego–Porto, IEEE, 2021,
272 pp. 1–8.
- 273 [29] R. R. MORALEDA, W. XIONG, N. A. VALOUS, AND N. HALAMA, *Segmentation of biomedical images based on*
274 *a computational topology framework*, in Seminars in Immunology, vol. 48, Elsevier, 2020, p. 101432.
- 275 [30] F. PEDREGOSA, G. VAROQUAUX, A. GRAMFORT, V. MICHEL, B. THIRION, O. GRISEL, M. BLONDEL, P. PRET-
276 TENHOFER, R. WEISS, V. DUBOURG, J. VANDERPLAS, A. PASSOS, D. COURNAPEAU, M. BRUCHER,
277 M. PERROT, AND E. DUCHESNAY, *Scikit-learn: Machine learning in Python*, Journal of Machine Learning
278 Research, 12 (2011), pp. 2825–2830.
- 279 [31] J. A. PIKE, A. O. KHAN, C. PALLINI, S. G. THOMAS, M. MUND, J. RIES, N. S. POULTER, AND I. B. STYLES,
280 *Topological data analysis quantifies biological nano-structure from single molecule localization microscopy*,
281 Bioinformatics, 36 (2020), pp. 1614–1621.
- 282 [32] Y. QIN, B. T. FASY, C. WENK, AND B. SUMMA, *A domain-oblivious approach for learning concise repre-*
283 *sentations of filtered topological spaces for clustering*, IEEE Transactions on Visualization and Computer
284 Graphics, 28 (2021), pp. 302–312.
- 285 [33] W. M. RAND, *Objective criteria for the evaluation of clustering methods*, Journal of the American Statistical
286 association, 66 (1971), pp. 846–850.
- 287 [34] N. SAUL AND C. TRALIE, *Scikit-TDA: Topological data analysis for python*, 2019, [https://doi.org/10.5281/](https://doi.org/10.5281/zenodo.2533369)
288 [zenodo.2533369](https://doi.org/10.5281/zenodo.2533369), <https://doi.org/10.5281/zenodo.2533369>.
- 289 [35] J. SCHUPBACH, *Dbspan test dataset*, Apr 2022, osf.io/uj32d.
- 290 [36] THE GUDHI PROJECT, *GUDHI User and Reference Manual*, GUDHI Editorial Board, 3.5.0 ed., 2022, <https://gudhi.inria.fr/doc/3.5.0/>.
- 291
292 [37] N. X. VINH, J. EPPS, AND J. BAILEY, *Information theoretic measures for clusterings comparison: Variants,*
293 *properties, normalization and correction for chance*, Journal of Machine Learning Research, 11 (2010),
294 p. 28372854.
- 295 [38] D. S. WILLARD, *New data structures for orthogonal range queries*, SIAM Journal on Computing, 14 (1985),
296 pp. 232–253.
- 297 [39] Y.-P. WU, J.-J. GUO, AND X.-J. ZHANG, *A linear DBSCAN algorithm based on LSH*, in 2007 International
298 Conference on Machine Learning and Cybernetics, vol. 5, IEEE, 2007, pp. 2608–2614.