
Characterizing and Measuring the Similarity of Neural Networks with Persistent Homology

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Characterizing the structural properties of neural networks is crucial yet poorly
2 understood, and there are no well-established similarity measures between net-
3 works. In this work, we observe that neural networks can be represented as abstract
4 simplicial complex and analyzed using their topological 'fingerprints' via Persis-
5 tent Homology (PH). We then describe a PH-based representation proposed for
6 characterizing and measuring similarity of neural networks. We empirically show
7 the effectiveness of this representation as a descriptor of different architectures
8 in several datasets. This approach based on Topological Data Analysis is a step
9 towards better understanding neural networks and serves as a useful similarity
10 measure.

11 1 Introduction

12 Machine learning practitioners can train different neural networks for the same task. Even for the
13 same neural architecture, there are many hyperparameters, such as the number of neurons per layer
14 or the number of layers. Moreover, the final weights for the same architecture and hyperparameters
15 can vary depending on the initialization and the optimization process itself, which is stochastic. Thus,
16 there is no direct way of comparing neural networks accounting for the fact that neural networks
17 solving the same task should be measured as being similar, regardless of the specific weights. This
18 also prevents one from finding and comparing modules inside neural networks (e.g., determining if a
19 given sub-network does the same function as other sub-network in another model). Moreover, there
20 are no well-known methods for effectively characterizing neural networks.

21 This work aims to characterize neural networks such that they can be measured to be similar
22 once trained for the same task, with independence of the particular architecture, initialization, or
23 optimization process. We focus on Multi-Layer Perceptrons (MLPs) for the sake of simplicity. We
24 start by observing that we can represent a neural network as a directed weighted graph to which we
25 can associate certain topological concepts.¹ Considering it as a simplicial complex, we obtain its
26 associated Persistent Diagram. Then, we can compute distances between Persistent Diagrams of
27 different neural networks.

28 The proposed experiments aim to show that the selected structural feature, Persistent Homology,
29 serves to relate neural networks trained for similar problems and that such a comparison can be
30 performed by means of a predefined measure between the associated Persistent Homology diagrams.
31 To test the hypothesis, we study different classical problems (MNIST, Fashion MNIST, CIFAR-10,
32 and language identification and text classification datasets), different architectures (number and size
33 of layers) as well as a control experiment (input order).

34 In summary, the main contributions of this work are the following:

¹See Jonsson [21] for a complete reference on graph topology.

- 35 • We propose an effective graph characterization strategy of neural networks based on Persistent Homology.
- 36
- 37 • Based on this characterization, we suggest a similarity measure of neural networks.
- 38 • We provide empirical evidence that this Persistent Homology framework captures valuable
- 39 information from neural networks and that the proposed similarity measure is meaningful.

40 The remainder of this paper is organized as follows. In Section 2, we go through the related work.
 41 Then, in Section 3 we describe our proposal and the experimental framework to validate it. Finally, in
 42 sections 4 and 5 we report and discuss the results and arrive to conclusions, respectively.

43 2 Related Work

44 One of the fundamental papers of Topological Data Analysis (TDA) is presented in Carlsson [8]
 45 and suggests the use of Algebraic Topology to obtain qualitative information and deal with metrics
 46 for large amounts of data. For an extensive overview of simplicial topology on graphs, see Giblin
 47 [18], Jonsson [21]. Aktas et al. [2] provide a thorough analysis of PH methods.

48 More recently, a number of publications have dealt with the study of the capacity of neural networks
 49 using PH. Guss and Salakhutdinov [19] characterize learnability of different neural architectures by
 50 computable measures of data complexity. Rieck et al. [30] introduce the *neural persistence* metric, a
 51 complexity measure based on TDA on weighted stratified graphs. This work suggests a representation
 52 of the neural network as a multipartite graph and the filtering of the Persistent Homology diagrams
 53 are performed for each layer independently. As the filtration contains at most 1-simplices (edges),
 54 they only capture zero-dimensional topological information, i.e. connectivity information. Donier
 55 [14] propose the concept of *spatial capacity allocation analysis*. Konuk and Smith [22] propose an
 56 empirical study of how NNs handle changes in topological complexity of the input data.

57 In terms of pure neural network analysis, there are relevant works, like Hofer et al. [20], that study
 58 topological regularization. Clough et al. [11] introduce a method for training neural networks for
 59 image segmentation with prior topology knowledge, specifically via Betti numbers. Corneanu et al.
 60 [13] try to estimate (with limited success) the performance gap between training and testing via
 61 neuron activations and linear regression of the Betti numbers.

62 On the other hand, topological analysis of decision boundaries has been a very prolific area. Ra-
 63 mamurthy et al. [28] propose a labeled Vietoris-Rips complex to perform PH inference of decision
 64 boundaries for quantification of the complexity of neural networks.

65 Naitzat et al. [27] experiment on the PH of a wide range of point cloud input datasets for a binary
 66 classification problems to see that NNs transform a topologically rich dataset (in terms of Betti
 67 numbers) into a topologically simpler one as it passes through the layers. They also verify that
 68 the reduction in Betti numbers is significantly faster for ReLU activations than hyperbolic tangent
 69 activations.

70 Liu [25] obtain certain geometrical and topological properties of decision regions for neural models,
 71 and provide some principled guidance to designing and regularizing them. Additionally, they use
 72 curvatures of decision boundaries in terms of network weights, and the rotation index theorem
 73 together with the Gauss-Bonnet-Chern theorem.

74 Regarding neural network representations, one of the most related works to ours, Gebhart et al. [16],
 75 focuses on topological representations of neural networks. They introduce a method for computing PH
 76 over the graphical activation structure of neural networks, which provides access to the task-relevant
 77 substructures activated throughout the network for a given input.

78 Interestingly, in Watanabe and Yamana [35], authors work on neural network representations through
 79 simplicial complexes based on deep Taylor decomposition and they calculate the PH of neural
 80 networks in this representation. In Chowdhury et al. [10], they use directed homology to represent
 81 MLPs. They show that the path homology of these networks is non-trivial in higher dimensions and
 82 depends on the number and size of the network layers. They investigate homological differences
 83 between distinct neural network architectures.

84 As far as neural network similarity measures are concerned, the literature is not especially prolific. In
 85 Kornblith et al. [23], authors examine similarity measures for representations (meaning, outputs of

86 different layers) of neural networks based on canonical correlation analysis. However, note that this
87 method *compares neural network representations (intermediate outputs), not the neural networks*
88 *themselves*. Remarkably, in Ashmore and Gashler [3], authors *do* deal with the intrinsic similarity
89 of neural networks themselves based on Forward Bipartite Alignment. Specifically, they propose
90 an algorithm for aligning the topological structures of two neural networks. Their algorithm finds
91 optimal bipartite matches between the nodes of the two MLPs by solving the well-known graph
92 cutting problem. The alignment enables applications such as visualizations or improving ensembles.
93 However, the methods only works under very restrictive assumptions,² and this line of work does not
94 appear to have been followed up.

95 Finally, we note that there has been a considerable growth of interest in applied topology in the
96 recent years. This popularity increase and the development of new software libraries,³ along with the
97 growth of computational capabilities, have empowered new works. Some of the most remarkable
98 libraries are Ripser [32, 5], and Flagser [26]. They are focused on the efficient computation of PH.
99 For GPU-Accelerated computation of Vietoris-Rips PH, Ripser++ [37] offers an important speedup.
100 The Python library we are using, Giotto-TDA [31], makes use of both above libraries underneath.

101 We have seen that there is a trend towards the use of algebraic topology methods for having a better
102 understanding of phenomena of neural networks and having more principled deep learning algorithms.
103 Nevertheless, little to no works have proposed neural network characterizations or similarity measures
104 based on intrinsic properties of the networks, which is what we intend to do.

105 3 Methodology

106 In this section, we propose our method, which is heavily based on concepts from algebraic topology.
107 We refer the reader to the Supplementary Material for the mathematical definitions. In this section,
108 we also describe the conducted experiments.

109 Intrinsically characterizing and comparing neural networks is a difficult, unsolved problem. First, the
110 network should be represented in an object that captures as much information as possible and then it
111 should be compared with a measure depending on the latent structure. Due to the stochasticity of
112 both the initialization and training procedure, networks are parameterized differently. For the same
113 task, different functions that effectively solve it can be obtained. Being able to compare the trained
114 networks can be helpful to detect similar neural structures.

115 We want to obtain topological characterizations associated to neural networks trained on a given
116 task. For doing so, we use the Persistence Homology (from now on, PH) of the graph associated to a
117 neural network. We compute the PH for various neural networks learned on different tasks. We then
118 compare all the diagrams for each one of the task.

119 More specifically, for each of the studied tasks (image classification on MNIST, Fashion MNIST and
120 CIFAR-10; language identification, and text classification on the Reuters dataset),⁴ we proceed as
121 follows:

- 122 • We train several neural network models on the particular problem.
- 123 • We create a directed graph from the weights of the trained neural networks (after changing
124 the direction of the negative edges and normalising the weights of the edges).
- 125 • We consider the directed graph as a simplicial complex and calculate its PH, using the
126 weight of the edges as the filtering parameter, which range from 0 to 1. This way we obtain
127 the so-called Persistence Diagram.
- 128 • We compute the distances between the Persistence Diagrams (prior discretization of the
129 Persistence Diagram so that it can be computed) of the different networks.
- 130 • Finally, we analyze the similarity between different neural networks trained for the same
131 task, for a similar task, and for a completely different task, independently of the concrete
132 architecture, to see whether there is topological similarity.

²For example, the two neural networks "must have the same number of units in each of their corresponding layers", and the match is done layer by layer.

³<https://www.math.colostate.edu/~adams/advising/appliedTopologySoftware/>

⁴For more details, see Section 3.2.

133 As baselines, we set two standard matrix comparison methods that are the 1-Norm and the Frobenius
 134 norm. Having adjacency matrix A and B , we compute the difference as $norm(A - B)$. However, these
 135 methods only work for matrices of similar size and thus, they are not general enough. We could also
 136 have used the Fast Approximate Quadratic assignment algorithm suggested in Vogelstein et al. [34],
 137 but for large networks this method becomes unfeasible to compute.

138 3.1 Proposal

139 Our method is as follows. We start by associating to a neural network a weighted directed graph
 140 that is analyzed as an abstract simplicial complex consisting on the union of points, edges, triangles,
 141 tetrahedrons and larger dimension polytopes (those are the elements referred as simplices). Abstract
 142 simplicial complexes are used in opposition to geometric simplicial complexes, generated by a point
 143 cloud embedded in the Euclidean space \mathbb{R}^n .

144 Given a trained neural network, we take the collection of neural network parameters as directed and
 145 weighted edges that join neurons, represented by graph nodes. Biases are considered as new vertices
 146 that join target neurons with an edge having a given weight. Note that, in this representation, we lose
 147 the information about the activation functions, for simplicity and to avoid representing the network
 148 as a multiplex network. Bias information could also have been ignored because we want large PH
 149 groups that characterize the network, while these connections will not change the homology group
 150 dimension of any order.

151 For negative edge weights, we reverse edge directions and maintain the absolute value of the weights.
 152 We discard the use of weight absolute value since neural networks are not invariant under weight sign
 153 transformations. This representation is consistent with the fact that every neuron can be replaced by a
 154 neuron from which two edges with opposite weights emerge and converge again on another neuron
 155 with opposite weights. From the point of view of homology, this would be represented as a closed
 156 cycle.

157 We then normalize the weights of all the edges as expressed in Equation 1 where w is the weight
 158 to normalize, W are all the weights and ζ is a smoothing parameter that we set to 0.000001. This
 159 smoothing parameter is necessary as we want to avoid normalized weights of edges to be 0. This is
 160 because 0 implies a lack of connection.

$$161 \quad \max\left(1 - \frac{|w|}{\max(|\max(W)|, |\min(W)|)}, \zeta\right) \quad (1)$$

162 Given this weighted directed graph, we then define a directed flag complex associated to it. Topology
 163 of this directed flag complex can be studied using homology groups H_n . In this work we calculate
 164 homology groups up to degree 3 (H_0 - H_3) due to computational complexity and our neural network
 165 representation method's layer connectivity limit.

166 The dimensions of these homology groups are known as Betti numbers. The i -th Betti number is
 167 the number of i -dimensional voids in the simplicial complex (β_0 gives the number of connected
 168 components of the simplicial complex, β_1 gives the number of non reducible loops and so on). For a
 169 deeper introduction to algebraic topology and computational topology, we refer to Edelsbrunner and
 170 Harer [15], Ghrist [17].

171 We work with a family of simplicial complexes, K_ε , for a range of values of $\varepsilon \in \mathbb{R}$ so that the complex
 172 at step ε_t is embedded in the complex at ε_{t+1} for $\varepsilon_t \leq \varepsilon_{t+1}$, i.e. $K_\varepsilon \subseteq K_{\varepsilon_{t+1}}$. In our case, ε is the
 173 minimum weight of included edges of our graph representation of neural networks.

174 The nested family of simplicial complexes is called a *filtration*. We calculate a sequence of homology
 175 groups by varying the ε parameter, obtaining a persistence homology diagram. PH calculations are
 176 performed on \mathbb{Z}_2 .

177 This filtration gives a collection of contained directed weighted graph or simplicial complex $K_{\varepsilon_{min}} \subseteq$
 178 $\dots \subseteq K_{\varepsilon_t} \subseteq K_{\varepsilon_{t+1}} \subseteq \dots \subseteq K_{\varepsilon_{max}}$, where $t \in [0, 1]$ and $\varepsilon_{min} = 0$, $\varepsilon_{max} = 1$ (recall that edge weights are
 179 normalized).

180 Given a filtration, one can look at the birth, when a homology class appears, and death, the time
 181 when the homology class disappears. The PH treats the birth and the death of these homological
 182 features in K_ε for different ε values. Lifespan of each homological feature can be represented as

182 an interval (*birth, death*), of the homological feature. Given a filtration, one can record all these
183 intervals by a Persistence Barcode (PB) [8], or in a Persistence Diagram (PD), as a collection of
184 multiset of intervals.

185 As mentioned previously, our interest in this work is to compare PDs from two different simplicial
186 complexes. There are two distances traditionally used to compare PDs, Wasserstein distance and
187 Bottleneck distance. Their stability with respect to perturbations on PDs has been object of different
188 studies [9, 12].

189 In order to make computations feasible and to obviate noisy intervals, we filter the PDs by limiting
190 the minimum PD interval size. We do so by setting a minimum threshold $\eta = 0.01$. Intervals with
191 a lifespan under this value are not considered. Additionally, for computing distances, we need to
192 remove infinity values. As we are only interested in the deaths until the maximum weight value, we
193 replace all the infinity values by 1.0.

194 Wasserstein distance calculations are computationally hard for large PDs (each PD of our NN models
195 has a million persistence intervals per diagram). Therefore we use a vectorized version of PDs instead,
196 also called PD discretization. This vectorized version summaries have been proposed and used on
197 recent literature [1, 6, 7, 24, 29].

198 For the persistence diagram distance calculation, we use the Giotto-TDA library [31] and compute
199 the following supported vectorized persistence summaries: 1. Persistence landscape. 2. Weighted
200 silhouette. 3. Heat vectorizations.

201 3.2 Experimental Framework

202 **Datasets** To determine the topological structural properties of trained NNs, we select different
203 kinds of datasets. We opt for four well-known benchmarks in the machine learning community
204 and one regarding language identification: (1) the MNIST⁵ dataset for classifying handwritten digit
205 images, (2) the Fashion MNIST [36] dataset for classifying clothing images into 10 categories, (3) the
206 CIFAR-10⁶ (CIFAR) dataset for classifying 10 different objects, (4) the Reuters dataset for classifying
207 news into 46 topics, and (5) the Language Identification Wikipedia dataset⁷ for identifying 7 different
208 languages.

209 We selected these datasets because, apart from being well-known benchmarks, the performances
210 without transfer learning are good enough and they have different data types and sizes. For CIFAR-
211 10 and Fashion MNIST datasets we train a Convolutional Neural Network (CNN) first, and the
212 convolutional layers are shared between all the models of the same dataset as a feature extractor.
213 Recall that in this work we are focusing on MLPs, so we do not consider that convolutional weights.
214 For the MNIST, Reuters and Language Identification datasets, we use an MLP. For Reuters and
215 Language identification datasets, we vectorize the sentences with character frequency.

216 **Experiments Pipeline** We study the following variables (hyperparameters): 1. Layer width,
217 2. Number of layers, 3. Input order⁸), 4. Number of labels (number of considered classes).

218 We define the *base* architecture as the one with a layer width of 512, 2 layers, the original features
219 order, and considering all the classes (10 in the case of MNIST, Fashion MNIST and CIFAR, 46 in
220 the case of Reuters and 7 in the case of the language identification task). Then, doing one change at a
221 time, keeping the rest of the base architecture hyperparameters, we experiment with architectures
222 with the following configurations:

- 223 • **Layer width:** 128, 256, 512 (*base*) and 1024.
- 224 • **Number of layers:** 2 (*base*), 4, 6, 8 and 10.
- 225 • **Input order:** 5 different randomizations (with *base* structure), the control experiment.
- 226 • **Number of labels** (MNIST, Fashion MNIST, CIFAR-10): 2, 4, 6, 8 and 10 (*base*).

⁵<http://yann.lecun.com/exdb/mnist/>

⁶<https://www.cs.toronto.edu/~kriz/cifar.html>

⁷<https://www.floydhub.com/floydhub/datasets/language-identification/1/data>

⁸Order of the input features, the control experiment. This one should definitely *not* affect the performance in the neural networks, so if our method is correct, it should be uniform as per the proposed topological distances.

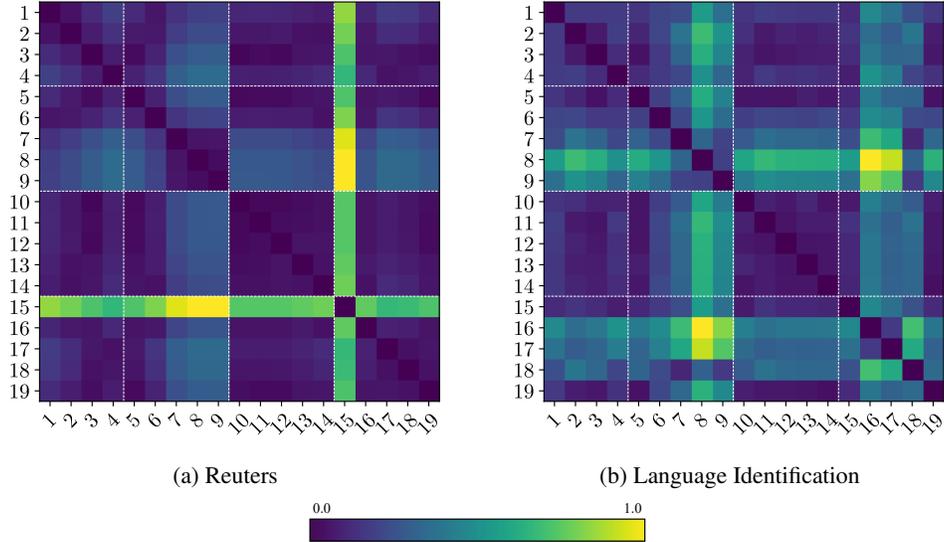


Figure 1: Distance matrices using Silhouette discretization.

- 227 • **Number of labels** (Reuters): 2, 6, 12, 23 and 46 (*base*).
- 228 • **Number of labels** (Language Identification): 2, 3, 4, 6 and 7 (*base*).

229 Note that this is *not* a grid search over all the combinations. We always modify one hyperparameter
 230 at a time, and keep the rest of them as in the base architecture. In other words, we experiment with all
 231 the combinations such that only one of the hyperparameters is set to a non-base value at a time.

232 For each dataset, we train 5 times (each with a different random weight initialization) each of these
 233 neural network configurations. Then, we compute the topological distances (persistence landscape,
 234 weighted silhouette, heat) among the different architectures. In total, we obtain $5 \times 5 \times 3$ distance
 235 matrices (5 datasets, 5 random initializations, 3 distance measures). Finally, we average the 5 random
 236 initializations, such that we get 5×3 matrices, one for each distance on each dataset. All the matrices
 237 have dimensions 19×19 , since 19 is the number of experiments for each dataset (corresponding to
 238 the total the number of architectural configurations mentioned above). Note that the base architecture
 239 appears 8 times (1, on the number of neurons per layer, 1 on the number of layers, 1 on the number of
 240 labels and the 5 randomizations of weight initializations).

241 All experiments were executed in a machine with 2 NVIDIA V100 of 32GB, 2 Intel(R) Xeon(R)
 242 Platinum 8176 CPU @ 2.10GHz, and of 1.5TB RAM, for a total of around 3 days.

243 The code and results are fully open source⁹ under MIT license.

244 4 Results & Discussion

245 Results from control experiments can be seen in the
 246 third group on Figures 1 and 4. In these figures, groups
 247 are separated visually using white dashed lines. Exper-
 248 iments groups are specified in Table 1. Control exper-
 249 iments in all the images appear very dimmed, which
 250 means that they are very similar, as expected. Recall
 251 that the control experiments consist of 5 (randomiza-
 252 tions) \times 5 (executions) and that 25 different neural
 253 networks have been trained; each one of the network
 254 has more than 690,000 parameters that have been ran-
 255 domly initialized. After the training, results show that
 256 these networks have very close topological distance, as expected.

Number	Experiment	Index
1	Layer size	1-4
2	Number of layers	5-9
3	Input order	10-14
4	Number of labels	15-19

Table 1: Indices of the experiments of the distance matrices.

⁹See Supplementary Material.

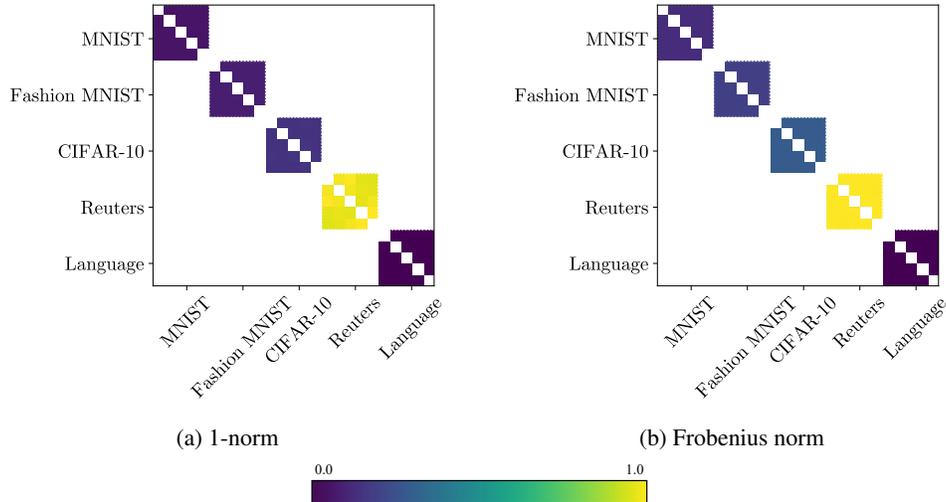


Figure 2: Control experiments using norms.

Norm	Minimum	Maximum	Mean	Standard deviation
1-Norm	0.6683	4.9159	1.9733	1.5693
Frobenius	0.0670	0.9886	0.4514	0.3074

Table 2: Normalized difference comparison of self-norm against the maximum mean distance of the experiment.

257 For Figure 2 we computed both 1-norm and Frobenius norm (the baselines) for graphs' adjacency
 258 matrices of control experiments. Note that as we ran the experiment five times, we make the mean
 259 for each value of the matrix. In order to show whether the resulting values are positive or negative,
 260 we subtract to the maximum difference of each dataset the norm of each cell separately, we take the
 261 absolute value and we divide by the maximum difference of each dataset. Therefore, we obtain five
 262 values per dataset. Table 2 shows the statistics reflecting that the distance among the experiments are
 263 large and, thus, they are not characterizing any similarity but rather an important dissimilarity.

264 In contrast, Figure 3, with our method (Silhouette), shows perfect diagonal of similarity blocks. In
 265 the corresponding numeric results, we obtained show small distances, as shown in Table 3. We can
 266 appreciate that each dataset has its own hub. This confirms the validity of our proposed similarity
 267 measure.

268 The method we present also seems to capture
 269 some parts of hyperparameter setup. For instance,
 270 in Figure 4 we can observe gradual increase of distances
 271 in the first group regarding layer size meaning that, as layer size increases,
 272 the topological distance increases too. Similarly, for the number of layers (second group) and
 273 number of labels (fourth group) the same situation holds. Note that in Fashion MNIST and
 274 CIFAR-10, the distances are dimmer because we are not dealing with the weights of the CNNs.
 275 Recall that the CNN acts as a frozen extractor and are pretrained for all runs (with the same
 276 weights), such that the MLP layers themselves are the only potential source of dissimilarity
 277 between runs.
 278
 279
 280
 281
 282
 283

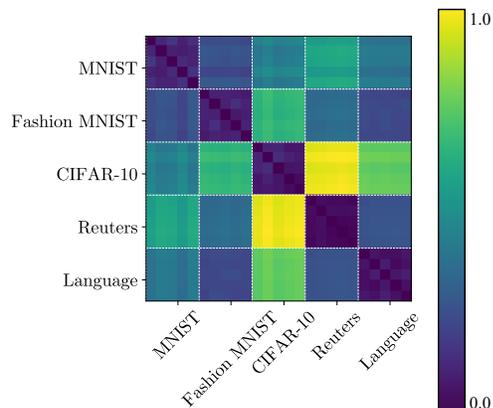


Figure 3: Control experiment comparison matrix using Silhouette discretization.

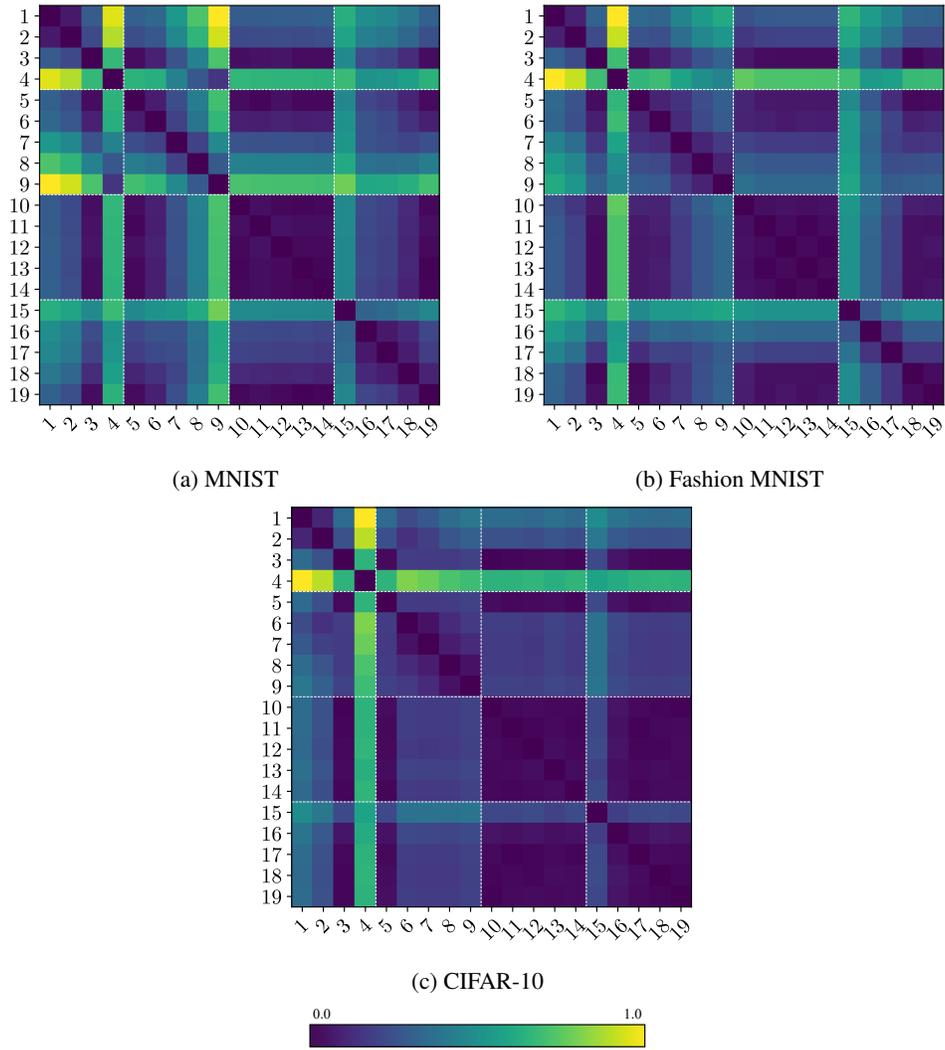


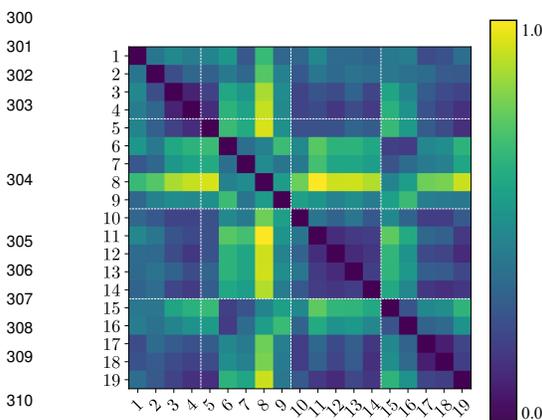
Figure 4: Distance matrices using Heat discretization.

Dataset	Heat distance		Silhouette distance	
	Mean	Deviation	Mean	Deviation
MNIST	0.0291	0.0100	0.1115	0.0364
F. MNIST	0.0308	0.0132	0.0824	0.0353
CIFAR-10	0.0243	0.0068	0.0769	0.0204
Language I.	0.0159	0.0040	0.0699	0.0159
Reuters	0.0166	0.0051	0.0387	0.0112

Table 3: PH distances across input order (control) experiments, normalized by dataset.

284 Thus, our characterization is sensitive to the architecture (e.g., if we increase the capacity, distances
 285 vary), but at the same time, as we saw before, it is not dataset-agnostic, meaning that it also captures
 286 whether two neural networks are learning the same problem or not.

287 In Figure 4, Fashion MNIST (Figure 4b) and CIFAR (Figure 4c) dataset results are interestingly
 288 different from those of MNIST (Figure 4a) dataset. This is, presumably, because both Fashion
 289 MNIST and CIFAR use a pretrained CNN for the problem. Thus, we must analyze the results taking
 290 into account this perspective. The first fully connected layer size is important as it can avoid a
 291 bottleneck from the previous CNN output. Some works in the literature show that adding multiple
 292 fully connected layers does not necessarily enhance the prediction capability of CNNs [4], which
 293 is congruent with our results when adding fully connected layers (experiments 5 to 9) that result in
 294 dimmer matrices than the one from. Concerning the experiments on input order, there is slightly
 295 more homogeneity than in MNIST, again showing that the order of sample has negligible influence.
 296 Moreover, there could have been even more homogeneity taking into account that the fully connected
 297 network reduced its variance thanks to the frozen weights of the CNN. This also supports the fact
 298 that the CNN is the main feature extractor of the network. As in MNIST results, CIFAR results show
 299 that the topological properties are, indeed, a mapping of the practical properties of neural networks.



312 Figure 5: Language Identification dataset PH Land-
 313 scape distance matrix.

314
 315
 316 In other words, it is not helpful for comparing PH diagrams associated to neural networks.

317 The most remarkable conclusion comes from the control experiments. The corresponding neural
 318 networks, with different input order but the same architecture, are very close to each other. The PH
 319 framework does, indeed, abstract away the specific weight values, and captures latent information
 320 from the networks, allowing comparisons to be based on the function they approximate. The selected
 321 neural network representation is reliable and complete, and yields coherent and meaningful results.
 322 Instead, the baseline measures, the 1-Norm and the Frobenius norm, implied an important dissimilarity
 323 between the experiments in the control experiments, meaning that they did not capture the fact that
 324 these neural networks were very similar in terms of the solved problem.

325 We conclude that our proposed characterization, does, indeed, capture meaningful information from
 326 neural network, and the computed distances can serve as an effective similarity measure between
 327 networks. To the best of our knowledge, this similarity measure between neural networks is the first
 328 of its kind.

329 As future work, we suggest adapting the method to different deep learning libraries and make it
 330 support popular neural architectures such as CNNs, Recurrent Neural Networks, and Transformers
 331 [33]. Finally, we suggest performing more analysis regarding the learning of a neural network, and
 332 trying to topologically answer the question of how a neural network learns.

We refer to the Supplementary Material for all distance matrices for all datasets and all distances, as well as for the standard deviations matrices and experiment group statistics.

5 Conclusions & Future Work

Results from different experiments, in five different datasets from computer vision and natural language, lead to similar topological properties and are trivially interpretable, which yields to general applicability.

The best discretizations chosen for this work are the Heat and Silhouette. They show better separation of experiment groups, and are effectively reflecting changes in a sensitive way. We also explored the Landscape discretization but it offers a very low interpretability and clearance.

333 **Checklist**

- 334 1. For all authors...
- 335 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
336 contributions and scope? [Yes]
- 337 (b) Did you describe the limitations of your work? [Yes] See second paragraph of the
338 introduction and last paragraph of the conclusions.
- 339 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 340 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
341 them? [Yes]
- 342 2. If you are including theoretical results...
- 343 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 344 (b) Did you include complete proofs of all theoretical results? [N/A]
- 345 3. If you ran experiments...
- 346 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
347 mental results (either in the supplemental material or as a URL)? [Yes] Both code and
348 outputs.
- 349 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
350 were chosen)? [Yes] Check the Experimental Framework Section and the code.
- 351 (c) Did you report error bars (e.g., with respect to the random seed after running ex-
352 periments multiple times)? [Yes] We include means, standard deviations and raw
353 outputs.
- 354 (d) Did you include the total amount of compute and the type of resources used (e.g., type
355 of GPUs, internal cluster, or cloud provider)? [Yes] Check Experimental Framework
356 Section.
- 357 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 358 (a) If your work uses existing assets, did you cite the creators? [Yes] In the case of the
359 datasets. We do not use any other additional asset.
- 360 (b) Did you mention the license of the assets? [No]
- 361 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
362 Code, results and pictures we have made for explanations.
- 363 (d) Did you discuss whether and how consent was obtained from people whose data you're
364 using/curating? [N/A]
- 365 (e) Did you discuss whether the data you are using/curating contains personally identifiable
366 information or offensive content? [N/A]
- 367 5. If you used crowdsourcing or conducted research with human subjects...
- 368 (a) Did you include the full text of instructions given to participants and screenshots, if
369 applicable? [N/A]
- 370 (b) Did you describe any potential participant risks, with links to Institutional Review
371 Board (IRB) approvals, if applicable? [N/A]
- 372 (c) Did you include the estimated hourly wage paid to participants and the total amount
373 spent on participant compensation? [N/A]

374 **References**

- 375 [1] H. Adams, T. Emerson, M. Kirby, R. Neville, C. Peterson, P. Shipman, S. Chepushtanova,
 376 E. Hanson, F. Motta, and L. Ziegelmeier. Persistence images: A stable vector representation of
 377 persistent homology. *J. Mach. Learn. Res.*, 18:8:1–8:35, 2017.
- 378 [2] M. Aktas, E. Akbaş, and A. E. Fatmaoui. Persistence homology of networks: methods and
 379 applications. *Applied Network Science*, 4:1–28, 2019.
- 380 [3] S. Ashmore and M. Gashler. A method for finding similarity between multi-layer perceptrons
 381 by forward bipartite alignment. In *2015 International Joint Conference on Neural Networks*
 382 (*IJCNN*), pages 1–7, 2015. doi: 10.1109/IJCNN.2015.7280769.
- 383 [4] S. H. S. Basha, S. R. Dubey, V. Pulabaigari, and S. Mukherjee. Impact of fully connected
 384 layers on performance of convolutional neural networks for image classification. *CoRR*,
 385 abs/1902.02771, 2019. URL <http://arxiv.org/abs/1902.02771>.
- 386 [5] U. Bauer. Ripser: efficient computation of victoris-rips persistence barcodes, 2021.
- 387 [6] E. Berry, Y.-C. Chen, J. Cisewski-Kehe, and B. T. Fasy. Functional summaries of persistence
 388 diagrams. *Journal of Applied and Computational Topology*, 4:211–262, 2020.
- 389 [7] P. Bubenik. Statistical topological data analysis using persistence landscapes. *J. Mach. Learn.*
 390 *Res.*, 16:77–102, 2015.
- 391 [8] G. Carlsson. Topology and data. *Bulletin of the American Mathematical Society*, 46:255–308,
 392 2009.
- 393 [9] F. Chazal, V. D. Silva, and S. Oudot. Persistence stability for geometric complexes. *Geometriae*
 394 *Dedicata*, 173:193–214, 2012.
- 395 [10] S. Chowdhury, T. Gebhart, S. Huntsman, and M. Yutin. Path homologies of deep feedforward
 396 networks. *2019 18th IEEE International Conference On Machine Learning And Applications*
 397 (*ICMLA*), pages 1077–1082, 2019.
- 398 [11] J. Clough, I. Öksüz, N. Byrne, V. Zimmer, J. A. Schnabel, and A. P. King. A topological
 399 loss function for deep-learning based image segmentation using persistent homology. *IEEE*
 400 *transactions on pattern analysis and machine intelligence*, PP, 2020.
- 401 [12] D. Cohen-Steiner, H. Edelsbrunner, and J. Harer. Stability of persistence diagrams. *Proceedings*
 402 *of the twenty-first annual symposium on Computational geometry*, 2005.
- 403 [13] C. Corneanu, M. Madadi, S. Escalera, and A. Martínez. Computing the testing error without a
 404 testing set. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
 405 pages 2674–2682, 2020.
- 406 [14] J. Donier. Capacity allocation analysis of neural networks: A tool for principled architecture
 407 design. *ArXiv*, abs/1902.04485, 2019.
- 408 [15] H. Edelsbrunner and J. Harer. *Computational Topology - an Introduction*. American Mathemat-
 409 ical Society, 2009.
- 410 [16] T. Gebhart, P. Schrater, and A. Hylton. Characterizing the shape of activation space in deep neu-
 411 ral networks. *2019 18th IEEE International Conference On Machine Learning And Applications*
 412 (*ICMLA*), pages 1537–1542, 2019.
- 413 [17] R. Ghrist. *Elementary Applied Topology*. Self-published, 2014.
- 414 [18] P. Giblin. *Graphs, surfaces, and homology : an introduction to algebraic topology*. Chapman
 415 and Hall, 1977.
- 416 [19] W. H. Guss and R. Salakhutdinov. On characterizing the capacity of neural networks using
 417 algebraic topology. *ArXiv*, abs/1802.04443, 2018.
- 418 [20] C. Hofer, F. Graf, M. Niethammer, and R. Kwitt. Topologically densified distributions. *ArXiv*,
 419 abs/2002.04805, 2020.
- 420 [21] J. Jonsson. *Simplicial complexes of graphs*. PhD thesis, KTH Royal Institute of Technology,
 421 2007.
- 422 [22] E. Konuk and K. Smith. An empirical study of the relation between network architecture and
 423 complexity. *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*,
 424 pages 4597–4599, 2019.

- 425 [23] S. Kornblith, M. Norouzi, H. Lee, and G. E. Hinton. Similarity of neural network representations
426 revisited. *CoRR*, abs/1905.00414, 2019. URL <http://arxiv.org/abs/1905.00414>.
- 427 [24] P. Lawson, A. Sholl, J. Brown, B. T. Fasy, and C. Wenk. Persistent homology for the quantitative
428 evaluation of architectural features in prostate cancer histology. *Scientific Reports*, 9, 2019.
- 429 [25] B. Liu. Geometry and topology of deep neural networks' decision boundaries. *ArXiv*,
430 abs/2003.03687, 2020.
- 431 [26] D. Lütgehetmann, D. Govc, J. Smith, and R. Levi. Computing persistent homology of directed
432 flag complexes. *arXiv: Algebraic Topology*, 2019.
- 433 [27] G. Naitzat, A. Zhitnikov, and L. Lim. Topology of deep neural networks. *J. Mach. Learn. Res.*,
434 21:184:1–184:40, 2020.
- 435 [28] K. Ramamurthy, K. R. Varshney, and K. Mody. Topological data analysis of decision boundaries
436 with application to model selection. *ArXiv*, abs/1805.09949, 2019.
- 437 [29] B. A. Rieck, F. Sadlo, and H. Leitte. Topological machine learning with persistence indicator
438 functions. *ArXiv*, abs/1907.13496, 2019.
- 439 [30] B. A. Rieck, M. Togninalli, C. Bock, M. Moor, M. Horn, T. Gumbsch, and K. Borgwardt.
440 Neural persistence: A complexity measure for deep neural networks using algebraic topology.
441 *ArXiv*, abs/1812.09764, 2019.
- 442 [31] G. Tauzin, U. Lupo, L. Tunstall, J. B. Pérez, M. Caorsi, A. Medina-Mardones, A. Dassatti,
443 and K. Hess. giotto-tda: A topological data analysis toolkit for machine learning and data
444 exploration, 2020.
- 445 [32] C. Tralie, N. Saul, and R. Bar-On. Ripser.py: A lean persistent homology library for python.
446 *The Journal of Open Source Software*, 3(29):925, Sep 2018. doi: 10.21105/joss.00925. URL
447 <https://doi.org/10.21105/joss.00925>.
- 448 [33] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and
449 I. Polosukhin. Attention is all you need. *CoRR*, abs/1706.03762, 2017. URL <http://arxiv.org/abs/1706.03762>.
- 451 [34] J. T. Vogelstein, J. M. Conroy, V. Lyzinski, L. J. Podrazik, S. G. Kratzer, E. T. Harley, D. E.
452 Fishkind, R. J. Vogelstein, and C. E. Priebe. Fast approximate quadratic programming for
453 graph matching. *PLOS ONE*, 10(4):1–17, 04 2015. doi: 10.1371/journal.pone.0121002. URL
454 <https://doi.org/10.1371/journal.pone.0121002>.
- 455 [35] S. Watanabe and H. Yamana. Topological measurement of deep neural networks using persistent
456 homology. In *ISAIM*, 2020.
- 457 [36] H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking
458 machine learning algorithms, 2017.
- 459 [37] S. Zhang, M. Xiao, and H. Wang. Gpu-accelerated computation of vietoris-rips persistence
460 barcodes. In *Symposium on Computational Geometry*, 2020.