
Inter-layer Information Similarity Assessment of Deep Neural Networks Via Topological Similarity and Persistence Analysis of Data Neighbour Dynamics

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The quantitative analysis of information structure through a deep neural network
2 (DNN) can unveil new insights into the theoretical performance of DNN architec-
3 tures. Two very promising avenues of research towards quantitative information
4 structure analysis are: 1) layer similarity (LS) strategies focused on the inter-layer
5 feature similarity, and 2) intrinsic dimensionality (ID) strategies focused on layer-
6 wise data dimensionality using pairwise information. Inspired by both LS and ID
7 strategies for quantitative information structure analysis, we introduce two novel
8 complimentary methods for inter-layer information similarity assessment premised
9 on the interesting idea of studying a data sample’s neighbourhood dynamics as it
10 traverses through a DNN. More specifically, we introduce the concept of **Near-
11 est Neighbour Topological Similarity** (NNTS) for quantifying the information
12 topology similarity between layers of a DNN. Furthermore, we introduce the concept
13 of **Nearest Neighbour Topological Persistence** (NNTP) for quantifying the
14 inter-layer persistence of data neighbourhood relationships throughout a DNN.
15 The proposed strategies facilitates the efficient inter-layer information similarity
16 assessment by leveraging only local topological information, and we demonstrate
17 their efficacy in this study by performing analysis on a deep convolutional neural
18 network architecture on image data to study the insights that can be gained with
19 respect to the theoretical performance of a DNN.

20 1 Introduction

21 Deep neural networks (DNNs) are functions that map information from one domain to another [1].
22 These maps often consist of hundreds of sub-maps in the form of element-wise non-linear functions,
23 matrix multiplications, convolutions, etc. [1]. Each one of these sub-maps gradually warps the
24 underlying manifold of a dataset. Studying the properties of these sub-maps and the effects on a
25 dataset’s manifold across a DNN at a micro and macro level can lead to a better understanding of a
26 DNN’s internal workings and can potentially guide improvement to their design.

27 At the micro level, intrinsic dimensionality (ID) methods [2, 3] allow for approximations of a mani-
28 fold’s dimensionality. Lacking from ID analysis is a notation of distance between layers. Knowing
29 the number of dimensions required to represent a manifold does not illuminate the manifold’s internal
30 characteristics, and directly comparing the magnitude of the ID between layers provides limited
31 actionable information. On the macro level layer similarity (LS) measures [4, 5, 6, 7] are designed
32 to compare the similarity of information representations between layers. LS measures work by
33 comparing the features of one layer to all the other features of another layer across a set of input data.
34 As such, measuring how a local region of the dataset manifold changes is not possible.

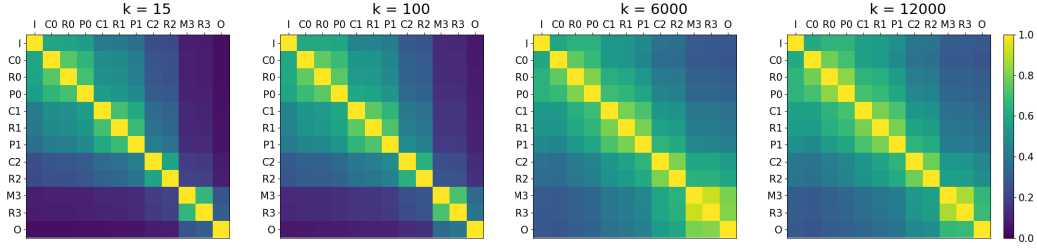


Figure 1: The Nearest Neighbour Topological Similarity between layers in a LeNet-5 model trained on the MNIST dataset for a different number of k nearest neighbours.

35 We propose a data centric approach to study the effects a DNN has on the local topological structure
 36 of a dataset’s manifold by taking inspiration from ID and LS methods. First we construct a nearest
 37 neighbour graph (NNG) to capture the topological structure of a dataset’s representation for each
 38 layer in a DNN. Then we compare each layers’ NNG using two novel forms of analysis, **Nearest**
 39 **Neighbour Topological Similarity** (NNTS) to measure the local topological similarity between
 40 layers, and **Nearest Neighbour Topological Persistence** (NNTP) to investigate inter-layer interacts
 41 on a pairwise data sample basis. These two proposed approaches open the door for fine-grained
 42 analysis of the complex dynamics present within a DNN. At a high level these methods compare the
 43 first degree relations between dataset samples within a layer to such relations in another layer.

44 2 Nearest Neighbour Topological Similarity

45 Below is a brief definition of the nearest neighbour graph (NNG) used within this work to capture
 46 properties of a dataset’s topological structure. See Appendix A for the motivation behind the graph’s
 47 design choices. Let $\mathbf{x}_i \in \mathbf{X}$ be a set of input samples, and let G represent a DNN. Let the output of
 48 some sub-function v_v for the \mathbf{x}_i sample be defined as $\mathbf{y}_{vi} = v_v(\mathbf{x}_i; G_v)$, where $G_v \subseteq G$ contains all
 49 required sub-functions, edges, and weights to calculate \mathbf{y}_{vi} . The main idea behind our approach is to
 50 use a graph of neighbours to capture the local structure between samples within a layer. For a given
 51 layer v_v with a set of outputs \mathbf{Y}_v , let $H_v = (\mathbf{Y}_v, D_v)$ be the graph of neighbours for layer v_v , where
 52 $\mathbf{Y}_v = v_v(\mathbf{X}; G_v)$ are the vertices of the graph, and D_v are the edges between two given samples
 53 $\mathbf{y}_{vi}, \mathbf{y}_{vj} \in \mathbf{Y}_v$. Let $K_{vi} \subseteq \mathbf{Y}_v$ be an ordered set of nearest neighbours of sample \mathbf{y}_{vi} . Directed edges
 54 are used for NNG construction.

55 Let $Q(H_a, H_b) = q_{ab}$ measure the Nearest Neighbour Topological Similarity (NNTS) between layers
 56 v_a and v_b where \mathbf{y}_{ai} and \mathbf{y}_{bi} are sample \mathbf{x}_i ’s representation in layers v_a and v_b , respectively. To
 57 compare a single sample across layers we propose a sample-wise similarity function $Q_s(\mathbf{y}_{ai}, \mathbf{y}_{bi})$
 58 where $Q(\cdot)$ is a function of all $Q_s(\cdot)$. Let $Q(H_a, H_b)$ be defined as

$$Q(H_a, H_b) = \frac{1}{n} \sum_i^n Q_s(\mathbf{y}_{ai}, \mathbf{y}_{bi}) \quad (1)$$

59 Then for a given sample \mathbf{x}_i for layers a and b we get neighbour sets K_{ai} and K_{bi} , respectively, for
 60 some given k . Let the per-sample inter-layer similarity function be defined as the IOU between layers.
 61 Note that this formulation also allows a sample to have different neighbour between layers.

$$Q_s(\mathbf{y}_{ai}, \mathbf{y}_{bi}) = \frac{|K_{ai} \cap K_{bi}|}{|K_{ai} \cup K_{bi}|} \quad (2)$$

62 $Q(\cdot)$ uses local information through first degree relations of a sample within a layer, and compares
 63 samples between layers by comparing the local characteristics of different representations of a sample.

64 We apply the notion of Nearest Neighbour Topological Similarity (NNTS) to a LeNet-5 [8] architec-
 65 ture trained on the MNIST [9] dataset to see how the local topological structure of a dataset changes
 66 across the model. Since LeNet-5 is a small architecture we break up what is normally considered
 67 a layer into their respective atomic operations before applying NNTS. We measure the NNNTS
 68 between all pairs of operations in the LeNet-5 model. The results for NNTS analysis are shown in
 69 Figure 1. We show the NNTS matrix for four different values of k , 15, 100, 6000, and 12000. The
 70 table headers along the top and left indicate the operation with the LeNet-5 model. The operations
 71 are causally aligned moving from left to right along the top, and top to bottom on the left. I stands
 72 for the input layer, C for convolutional operations, R for ReLU activation, P for max-pooling, M for
 73 matrix multiplication, and O for output (note that O is also a matrix multiplication operation). The
 74 number following the operation identifier indicates the layer number.

75 Each of the NNTS matrices are symmetric about the diagonal. The diagonal in each plot all have
76 values of 1 since layers are all self similar. Notice the block-like pattern staggered both vertically
77 and horizontally in all four plots every time the layer number changes (e.g., moving from P0 to
78 C1, or from R3 to O). This is a clear indication that sequences of operations which are normally
79 considered *layers* (within LeNet-5) are not arbitrary since internal representations within a layer are
80 more similar to one another than to operations outside the layer. The inter-layer similarity pattern
81 even persists when comparing the first operation in a layer to the first operation in the following layer.
82 This observation could be used to study other standard layer designs (e.g., a layer designed with
83 batchnorm) to determine if such designs follow the same inter-layer similarity pattern. Notice that
84 the similarity between P1 and M3 is marginally smaller (about 0.05) when compared the similarity
85 between R2 and M3. The marginal difference is a good indication that removing layer 3 will have
86 little effect on network performance. The drastic change in plots $k = 6000$ and $k = 12000$ is less
87 obvious due to each neighbour being connected to 10% and 20% of the whole dataset, respectively.

88 As the number of k nearest neighbours increases from $k = 5$ to $k = 6000$ there is a gradual increase
89 in the similarity between all layer pairs. The transition between from $k = 6000$ to $k = 12000$ sees
90 a decrease in similarity, and most noticeably between the last couple of operations (bottom right
91 hand corner), in the LeNet-5 model. This decrease is to be expected considering that MNIST has
92 ten classes with approximately 6000 samples per class. Near the end of the network samples from
93 the same class should be clustered near one another. At a $k = 6000$ a sample’s connections will
94 mostly consist of all samples from within the class. Any operation performed on samples from the
95 same class would likely have the same effect and thus not effect the inter-layer neighbour relations.
96 But when $k = 12000$ half a sample’s neighbours will be from other classes. While operations are
97 unlikely to effect intra-class neighbours, they can still effect inter-class neighbours, and thus resulting
98 in the decrease in similarity from $k = 6000$ to $k = 12000$. It is expected that the inter-layer similarity
99 converges to 1 as the number of connections approaches the number of samples in the dataset.

100 3 Nearest Neighbour Topological Persistence

101 Reducing the similarity between two layers to a single value provides a useful measure for high level
102 measure for topological similarity. On the other hand, such reduction also removes most of the local
103 inter-sample relationship information, thereby reducing one’s ability to study the complex interactions
104 between layers throughout a network. In this section we introduce an approach from which higher
105 order analysis can be performed. Specially, we investigate when pairs of samples become neighbours
106 in a DNN, properties of the pairs while they are neighbours, and when pairs of samples are no longer
107 neighbours. We call such analysis Nearest Neighbour Topological Persistence (NNTP).

108 Consider a network where layers follow a sequential design $v_{in} > \dots > v_a > \dots > v_b > \dots >$
109 $v_{out} \in V$, where layer v_{in} is the input layer and v_{out} is the output layer of a DNN. Let e_{ij} be an
110 abstract un-directed connection between samples $(\mathbf{x}_i, \mathbf{x}_j)$, and let e_{ij}^v be the un-directed connection
111 between samples $(\mathbf{x}_i, \mathbf{x}_j)$ in layer v . Let $e_{ij}^v \in H_v$ iff either of the corresponding directional
112 connections are in H_v , where H_v is the NNG of layer v .

113 Let e_{ij} be α -persistent between two layers v_a and v_b if there exists no more than α contiguous
114 layers in the chain of layers $v_a > \dots > v_b$ where $e_{ij}^v \notin H_v$ for all v ’s within the chain of layers.
115 α -persistent is a whole family of measures. In this work we only investigate 0-persistent sample
116 pairs. 0-persistent can be interpreted as a measure for local network stability. If a connection persists
117 across a series of layers then it is reasonable to assume that the pair is located in a local region of the
118 dataset’s manifold that share specific features. Analysing when samples are no longer neighbours
119 may help illuminate what specific features a given network layer is detecting. When considering the
120 entire dataset using this approach one can see the interactions between layers. For example, aspects
121 of a network like connection-cancellation would become evident (i.e., if one layer moves a lot of
122 samples near each other and a down stream layer moves those samples apart). By studying how
123 layers interact with each other on a more granular level (when compared to scalar LS measures) one
124 can tailor a DNN’s design at both a macro-architecture resolution and a micro-architecture resolution.

125 We apply the notion of persistence to a LeNet-5 architecture trained on the MNIST dataset. We break
126 the LeNet-5 model into the same atomic operations as done in the previous section. For 0-persistent
127 analysis we count how many pair-wise nearest neighbour connections are 0-persistent between all
128 pairs of layers in the LeNet-5 model. The results for 0-persistent analysis are shown in Table 1.
129 The headers along the top and left indicate the operation with the LeNet-5 model. The operations are

Table 1: LeNet-5 0-persistent matrix. Each operation pair (v_{first}, v_{last}) counts the number of 0-persistent pairwise samples (in the thousands) that start at operation v_{first} and last appear at operation v_{last} .

Layer of 0-persistent beginning	Layer of 0-persistent end											
	I	C0	R0	P0	C1	R1	P1	C2	R2	M3	R3	O
I	455	122	6.99	5.95	2.51	3.02	36.0	0.99	0.59	3.94	0.32	28.7
C0		427	2.94	5.30	0.65	0.52	2.27	1.23	0.06	0.18	0.01	0.62
R0			550	2.72	0.57	0.61	10.2	0.13	0.02	0.09	0.01	0.46
P0				494	19.0	22.1	6.08	0.54	0.35	2.50	0.42	0.53
C1					110	114	286	1.71	0.17	0.82	0.10	2.83
R1						67.2	57.7	1.33	0.23	0.29	0.04	0.39
P1							205	4.53	2.54	0.57	0.12	0.55
C2								528	36.3	8.79	0.57	0.46
R2									291	104	10.4	168
M3										170	19.1	141
R3											149	128
O												180

130 causally aligned when moving from left to right a long the top, and top to bottom on the left. We use
 131 $k = 15$ for the number of neighbours each sample has.

132 Notice the large number of connections present along the diagonal. These connections only sequen-
 133 tially exist for one layer (note that they may reappear in other layers). Let connections along the
 134 diagonal be called transient connections. Many layers in the LeNet-5 model have a plurality of their
 135 pairwise sample connections existing as transient connections, with the first layer (i.e. layer 0) being
 136 especially transient heavy. This may indicate that the first couple of operations are mainly responsible
 137 for placing the samples in approximately their final location in the data’s manifold for classification,
 138 with the rest of the layers being responsible for fine tuning.

139 Another interesting observation is the number of connections present in the top right layer pair (I,O).
 140 These connections persist throughout all operations in the LeNet-5 model, indicating that they are
 141 likely to be true neighbours on the data’s intrinsic manifold. Studying the relationship between such
 142 neighbours would be useful in a number of areas including building more robust datasets, tracking
 143 clusters of strongly persistent neighbours (i.e., connections that are persistent across many layers),
 144 and training a model on a reduced number of samples.

145 From this matrix one can see that C2 and R2 seem to have little effect on the data manifold as
 146 they largely add persistent connections while allowing most other connections to pass though. For
 147 applications like layer reduction, C2 and R2 are potentially strong candidates for layer removal,
 148 and even more so considering that C2 has the largest number of parameters when compared to the
 149 other convolutional operations. One anomaly with C2 is that it largely kills connections created by
 150 C1 as indicated by the operation pair (C1, P1) of 286000. Notice that 286000 is by far the largest
 151 non-transient group of connections in Table 1. In a sense, C2 is undoing the alterations to the data
 152 manifold made by C1. In addition, such a relationship does not exist between C0 and C1, or C0 and
 153 C2. Further research is required to understand such behavior.

154 4 Conclusion and Future Work

155 We propose two complementary data centric analytic methods for studying the complex dynamics
 156 of a dataset’s manifold as it moves through a DNN using a set nearest neighbour graphs. The first
 157 proposed approach, Nearest Neighbour Topological Similarity, measures the local similarity between
 158 two NNGs, and second proposed approach Nearest Neighbour Topological Persistence captures the
 159 complex local interactions between layers. We demonstrate that both these approaches have the
 160 potential for providing a better understanding interactions between layers on a local topological level,
 161 and how such insights can be used to built better DNNs. Future directions of research include, but not
 162 limited to, using the proposed approach to study local clusters of data throughout a DNN, studying
 163 how a family of operations (e.g., activation functions) effects local characteristics of a dataset’s
 164 manifold, and measuring how a manifold changes throughout training a DNN.

165 **References**

- 166 [1] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- 167 [2] E. Facco, M. d’Errico, A. Rodriguez, and A. Laio, “Estimating the intrinsic dimension of datasets by a
168 minimal neighborhood information,” *Scientific Reports*, 2017.
- 169 [3] A. Ansuini, A. Laio, J. H. Macke, and D. Zoccolan, “Intrinsic dimension of data representations in deep
170 neural networks,” in *Advances in Neural Information Processing Systems (NIPS)*, 2019.
- 171 [4] A. Gretton, O. Bousquet, A. Smola, and B. Schölkopf, “Measuring statistical dependence with hilbert-
172 schmidt norms,” in *International Conference on Algorithmic Learning Theory*, Springer, 2005.
- 173 [5] D. R. Hardoon, S. Szedmak, and J. Shawe-Taylor, “Canonical correlation analysis: An overview with
174 application to learning methods,” *Neural Computation*, 2004.
- 175 [6] M. Raghu, J. Gilmer, J. Yosinski, and J. Sohl-Dickstein, “Svcca: Singular vector canonical correlation
176 analysis for deep learning dynamics and interpretability,” in *Advances in Neural Information Processing
177 Systems (NIPS)*, 2017.
- 178 [7] S. Kornblith, M. Norouzi, H. Lee, and G. Hinton, “Similarity of neural network representations revisited,”
179 *arXiv preprint arXiv:1905.00414*, 2019.
- 180 [8] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,”
181 *Proceedings of the IEEE*, 1998.
- 182 [9] Y. LeCun, “The mnist database of handwritten digits,” <http://yann.lecun.com/exdb/mnist/>, 1998.

183 A Nearest Neighbour Graph

184 Let $\mathbf{x}_i \in \mathbf{X}$ be a set of input samples of shape $n \times d$, and let $G = (V, E, W)$, represent a DNN, where
185 $V = \{v\}$ is a set of sub-functions, $E = \{e\}$ is a set of edges that represent the sub-function's i/o
186 relationships, and $W = \{w\}$ is a set of weights that parameterize the sub-functions. Let the output of
187 some sub-function v_v for the \mathbf{x}_i sample be defined as $\mathbf{y}_{vi} = v_v(\mathbf{x}_i; G_v)$, where $G_v \subseteq G$ contains all
188 required sub-functions, edges, and weights to calculate \mathbf{y}_{vi} .

189 The main idea behind our approach is to use a graph of neighbours to capture the local structure
190 between samples within a layer. More formally, let $\mathbf{Y}_v = v_v(\mathbf{X}; G_v)$ be a one-to-one mapping of
191 samples from the input space to the space of layer v_v of a DNN. For a given layer v_v with a set of
192 outputs \mathbf{Y}_v , let $H_v = (\mathbf{Y}_v, D_v)$ be the graph of neighbours for layer v_v , where \mathbf{Y}_v are the vertices
193 of the graph, and D_v are the edges between two given samples $\mathbf{y}_{vi}, \mathbf{y}_{vj} \in \mathbf{Y}_v$. Let $K_{vi} \subseteq \mathbf{Y}_v$ be an
194 ordered set of nearest neighbours of sample \mathbf{y}_{vi} .

195 The goal of the graph is to represent localized information from the samples. As such, a metric
196 for measuring distance between two samples in a given layer is required. In general there are two
197 common methods used. The first approach uses a distance threshold to find all samples $\mathbf{y}_{vj} \in K_{vi}$
198 that are within some fixed radius r_v of sample \mathbf{y}_{vi} , where r_v is constant for the entire graph. Note that
199 each set K_{vi} for a single layer v_v can contain a variable number of neighbours. The second approach
200 uses a variable radius but a fixed number of samples k in K_{vi} for each sample \mathbf{y}_{vi} . Such an approach
201 is called a k nearest neighbour (k -nn) graph. For this work a k -nn based approach is used to ensure
202 that each sample $\mathbf{y}_{vj} \in \mathbf{Y}_v$ has a neighbour (i.e., $|K_{vi}| > 0$). Note that it would be possible to find
203 the smallest radius such that every sample has at least one neighbour, but this would also allow for
204 samples to be connected to the entire graph (e.g., when there is one extreme outlier).

205 To build a k -nn graph one must choose if connections are directed or un-directed, what distance metric
206 to use, and the number of neighbours. One of the features a distance metric requires is that the metric
207 be commutative (i.e., $\langle x, y \rangle = \langle y, x \rangle$). From a k -nn graph's perspective this requires that connections
208 between samples be un-directed. That is, if sample \mathbf{y}_{vi} is a neighbour of \mathbf{y}_{vj} , then \mathbf{y}_{vj} must also be a
209 neighbour of \mathbf{y}_{vi} . However, the un-directed nature of connections would require a loosening of the
210 fixed number of neighbours inherent to k -nn graphs as a k -nn graph with k un-directed edges per
211 sample may not exist.

212 One way to loosen the neighbourhood criteria is to perform an intersection where by two samples
213 are un-directed neighbours iff both samples are directed neighbours of each other; this effectively
214 sets an upper bound to the number of neighbours to k . Such an approach undermines the choice of
215 a k -nn graph in that some samples might not have neighbours. Another way to solve the issue is
216 to perform a union where by two samples are un-directed neighbours iff either sample is a directed
217 neighbour of one another; effectively setting k as the lower bound to the number of neighbours. This
218 approach can result in some samples having orders of magnitude more neighbours than other samples.
219 A third option to loosen the neighbourhood criteria is to just use directed edges, thereby ensuring
220 every sample has the same number of neighbours. In this proposal directed edges are used for nearest
221 neighbour graph (NNG) construction, other graph representations will be studied in the future work.
222 For this work a euclidean based distance metric is used.