
Topological Deep Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Topological deep learning is a formalism that is aimed at introducing topological
2 language to deep learning for the purpose of utilizing the minimal mathematical
3 structures to formalize problems that arise in a generic deep learning problem. In
4 this article, we define and study the classification problem in machine learning in a
5 topological setting. Using this topological framework, we show when the classifica-
6 tion problem is possible or not possible in the context of neural networks. Finally,
7 we demonstrate how our topological setting immediately illuminates aspects of
8 this problem that are not as readily apparent using traditional tools.

9 1. Introduction

10 Recent years have witnessed increased interest in the role topology plays in machine learning and data
11 science [5]. Topology is a natural tool that allows the formulation of many longstanding problems in
12 these fields. For instance, *persistent homology* [10] has been overwhelmingly successful at finding
13 solutions to a vast array of complex data problems [1, 2, 3, 6, 7, 9, 12, 16, 17, 18, 19, 20, 23, 25, 26].

14 On the other hand, the role that topology plays in deep learning is still mostly restricted to techniques
15 that attempt to enhance machine learning models [14, 4, 28]. However, we believe that topology
16 can and will play a central role in deep learning and AI in general. Our purpose of this article is to
17 introduce *topological deep learning*, a formalism that is aimed at introducing topological language to
18 deep learning for the purpose of utilizing the minimal mathematical structures to formalize problems
19 that arise in a generic deep learning problem.

20 To this end we define and study the classification problem in a topological setting. Using this
21 topological machinery, we show when the classification problem is possible or not possible in the
22 context of neural networks. Finally, we show how the architecture of a neural network cannot be
23 chosen independently from the topology of the underlying data. To demonstrate these results, we
24 provide an example dataset and show how it is acted upon by a neural net from this topological
25 perspective.

26 2. Background

27 A *neural network*, or simply a *network*, is a function $Net : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R}^{d_{out}}$ defined by a composition
28 of the form:

$$Net := f_L \circ \dots \circ f_1 \tag{1}$$

29 where the functions f_i , $1 \leq i \leq L$ are called the *layer functions*. A layer function $f_i : \mathbb{R}^{n_i} \rightarrow \mathbb{R}^{m_i}$ is
30 typically a continuous, piece-wise smooth function of the following form: $f_i(x) = \sigma(W_i(x) + b_i)$
31 where W_i is an $m_i \times n_i$ matrix, b_i is a vector in \mathbb{R}^{m_i} , and $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is an appropriately chosen
32 nonlinear function that is applied coordinate-wise on an input vector (z_1, \dots, z_{m_i}) to get a vector
33 $(\sigma(z_1), \dots, \sigma(z_{m_i}))$.

34 3. Data In a Topological Setting

35 In the present article we clearly distinguished between data and the functions that operate on it. This
36 distinction is important because data as a separate mathematical object have complex properties that
37 intertwine non-trivially with the functions, that also have unique properties, that operate on the data.
38 The purpose of this section is define the notion of data using topological notions.

39 3.1. Topological Data

40 Denote by M^n to a manifold M of dimension n . Let $D = M_1^{i_1} \cup M_2^{i_2} \dots \cup M_k^{i_k}$ be a disjoint union of
41 k compact manifolds. Let $h : D \rightarrow E$ be a continuous function on D . We refer to the pair (D, h)
42 as *topological data* and refer to E as the *ambient space* of the topological data, or simply the
43 ambient space of the data.

44 A few remarks here must be made about the above definition. First note that the definition above is
45 consistent with the statistical version. The space E , usually some Euclidean space, represents the
46 ambient space of a probability distribution μ from which we sample the data. The support of μ is
47 $\mathcal{D} := h(D)$. The assumption that the data lives on a manifold-like structure is justified in the literature
48 [11, 21].¹

49 3.2. Topologically Labeled Data

50 Let (D, h) be topological data with $h : D \rightarrow \mathcal{D} \subset E$. Let $\mathcal{Y} = \{l_1, \dots, l_d\}$ be a finite set. A *topological*
51 *labeling* on \mathcal{D} is a closed subset $\mathcal{D}_L \subset \mathcal{D}$ along with a surjective continuous function $g : \mathcal{D}_L \rightarrow \mathcal{Y}$
52 where \mathcal{Y} is given the discrete topology. The triplet (D, h, g) will be called *topologically labeled data*.

53 Topologically labeled data is a topological object that corresponds to labeled data in the typical
54 statistical setting for a supervised classification machine learning problem.

55 4. The Topological Classification Problem

56 With the above setting we now demonstrate how to realize the classification problem as a topological
57 problem. In what follows we set \mathcal{D}_k to denote $g^{-1}(l_k)$ for $l_k \in \mathcal{Y}$.

58 **Definition 1.** Let (D, h, g) be topologically labeled data with, $h : D \rightarrow \mathcal{D} \subset \mathbb{R}^n$ and $g : \mathcal{D}_L \rightarrow \mathcal{Y}$
59 where $|\mathcal{Y}| = d$. A *topological classifier* on (D, h, g) is a continuous function $f : \mathbb{R}^n \rightarrow \mathbb{R}^k$. We say that
60 f separates the topologically labeled data (D, h, g) if we can find k disjoint embedded k -dimensional
61 discs A_1, \dots, A_k in \mathbb{R}^k such that $f(\mathcal{D}_k) \subset A_k$.

62 In general, a topologically labeled data can be knotted, linked and entangled together in a non-trivial
63 manner by the embedding h , and the existence of a function f that separates this data is not immediate.
64 The preceding description is an topological rewording of the classification problem typically given in
65 a statistical setting. Indeed, a successful classifier tries to *separate* the labeled data by mapping the
66 raw input data into another space where this data can be separated easily according to the given class.

67 The function f is the learning function that we try to compute, in practice. The first question one
68 could ask in this context is one of existence: given topologically labeled data (D, h, g) when can we
69 find a function f that separates this data? We answer this question next.

70 4.1. Topological Classifiers and Separability of Topologically Labeled Data

71 We start with the binary classification problem, namely when $|\mathcal{Y}| = 2$. We have the following
72 proposition:

73 **Proposition 4.1.** Let (D, h, g) by a topologically labeled data with $h : D \rightarrow \mathcal{D} \subset \mathbb{R}^{d_{in}}$ and
74 $g : \mathcal{D}_L \rightarrow \{l_1, l_2\}$. Then there exists a topological classifier $f : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R}$ that separates (D, h, g) .

75 **Proof.** The label function $g : \mathcal{D}_L \rightarrow \{l_1, l_2\}$ induces a partition on \mathcal{D}_L into two disjoint closed sets
76 $\mathcal{D}_1 := g^{-1}(l_1)$ and $\mathcal{D}_2 := g^{-1}(l_2)$. By Urysohn's lemma there exists a function $f^* : \mathcal{D}_L \rightarrow [0, 1]$

¹While we make this assumption here, it not strictly necessary anywhere in our proofs.

77 such that $f^*(\mathcal{D}_1) = 0$ and $f^*(\mathcal{D}_2) = 1$. Since \mathcal{D}_L is closed in $\mathbb{R}^{d_{in}}$ then by Tietze extension theorem
 78 there exists an extension of f^* to a continuous function $\hat{f} : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R}$ such that $f^*(\mathcal{D}_L) = f(\mathcal{D}_L)$. In
 79 particular, $f(\mathcal{D}_1) = 0$ and $f(\mathcal{D}_2) = 1$. Hence the function f separates (D, h, g) .

80 Proposition 4.1 can be easily generalized to obtain functions that separate (D, h, g) in any Euclidean
 81 space \mathbb{R}^k . Namely, for any $k \geq 1$ there exists a continuous map $F : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R}^k$ that separates
 82 (D, h, g) . This can be done by defining $F = (f_1, f_2)$ where $f_1 : \mathbb{R}^{d_{in}} \rightarrow [0, 1]$ is the continuous
 83 function guaranteed by Urysohn's Lemma and $f_2 : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R}^{k-1}$ is an arbitrary continuous function.
 84 This function F clearly separates (X, h, g) . We record this fact in the following proposition.

85 **Proposition 4.2.** *Let (D, h, g) be a topologically labeled data with $h : D \rightarrow \mathcal{D} \subset \mathbb{R}^{d_{in}}$ and $g : \mathcal{D}_L \rightarrow$
 86 $\{l_1, l_2\}$. Then for any $k \geq 1$ there exists a continuous map $f : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R}^k$ that separates (D, h, g) .*

87 Proposition 4.2 can be generalized to the case when the set \mathcal{Y} has an arbitrary finite size. This can
 88 be done by because Urysohn's Lemma remains valid when we start with n disjoint sets instead of 2.
 89 The following theorem, which generalizes 4.2, asserts the existence of a topological classifier f that
 90 separates any given topologically labeled data.

91 **Theorem 4.3.** *Let (D, h, g) be topologically labeled data with $h : D \rightarrow \mathcal{D} \subset \mathbb{R}^{d_{in}}$ and $g : \mathcal{D}_L \rightarrow \mathcal{Y}$.
 92 Then there exists a continuous map $f : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R}^k$ that separates (D, h, g) for any integer $k \geq 1$.*

93 5. Neural Networks as Topological Classifiers

94 Let (D, h, g) be a topologically labeled data with, $h : D \rightarrow \mathcal{D} \subset \mathbb{R}^{d_{in}}$ and $g : \mathcal{D}_L \rightarrow \mathcal{Y} = \{l_1, \dots, l_n\}$.
 95 Can we find a neural network defined on $\mathbb{R}^{d_{in}}$ that separates the data (D, h, g) ? We start by framing
 96 the softmax classification networks using topological terminologies.

97 Typical, classification neural networks have a special layer function at the end where one uses the
 98 *softmax activation function*². Denote by Δ_n the n^{th} simplex as the convex hull of the vertices
 99 $\{v_0, \dots, v_n\}$ where $v_i = (0, \dots, 1, \dots, 0) \in \mathbb{R}^{n+1}$ with the lone 1 in the $(i+1)^{th}$ coordinate.

100 The *softmax function* on n vertices $softmax : \mathbb{R}^n \rightarrow \text{Int}(\Delta_{n-1}) \subset \mathbb{R}^n$, is defined by the compo-
 101 sition $S \circ Exp$ where $Exp : \mathbb{R}^n \rightarrow (\mathbb{R}^+)^n$ is defined by : $Exp(x_1, \dots, x_n) = (\exp(x_1), \dots, \exp(x_n))$,
 102 and $S : \mathbb{R}^n \rightarrow \Delta_{n-1}$ is defined by : $S(x_1, \dots, x_n) = (x_1 / \sum_{i=1}^n x_i, \dots, x_n / \sum_{i=1}^n x_i)$.

103 A network Net is said to be a *softmax classification neural network* with n labels if the final layer
 104 of Net is softmax function with n vertices. Usually n is the number of labels in the classification
 105 problem. Each vertex v_i in Δ_{n-1} corresponds to precisely one label $l_{i+1} \in \mathcal{Y}$ for $0 \leq i \leq n-1$.

106 For an input $x \in \mathcal{D}$ the point $Net(x)$ is an element of Δ_{n-1} . By definition, the point x is assigned to
 107 the label l_{i+1} by the neural network if and only if $Net(x) \in \text{Int}(VC(v_i))$ where $VC(C)$ denotes
 108 the Voronoi cell of the set C and $\text{Int}(A)$ denotes the interior of a set A . This immediately yields the
 109 following theorem.

110 **Theorem 5.1.** *Let (D, h, g) be a topologically labeled data with, $h : D \rightarrow \mathcal{D} \subset \mathbb{R}^{d_{in}}$ and $g : \mathcal{D}_L \subset$
 111 $\mathbb{R}^{d_{in}} \rightarrow \{l_1, \dots, l_n\}$. A softmax classification neural network $Net : \mathbb{R}^{d_{in}} \rightarrow \text{Int}(\Delta_{n-1})$ separates
 112 (D, h, g) if and only if $Net(\mathcal{D}_{i+1}) \subset \text{Int}(VC(v_i))$ for $0 \leq i \leq n-1$.*

113 Finally, to answer the question about the ability of a neural network to separate a topologically labeled
 114 data, we combine the result we obtained from Theorem 4.3 with the universality of neural networks
 115 [8, 13, 22]³. The universality of neural networks essentially states that for any continuous function f
 116 we can find a network that approximates it to an arbitrary precision⁴. Hence we conclude that any
 117 topologically labeled data can effectively be separated by a neural network.

²There are other types of classification neural networks but this is beyond the scope of our discussion here

³The universal approximation theorem is available in many flavors : one may fix the depth of the network and vary the width or the other way around.

⁴The closeness between functions is with respect to an appropriate functional norm. See [8, 22] for more details.

118 6. Shape of Data and Neural Networks

119 We end our discussion by briefly showing how the shape of input data is essential when deciding on
 120 the architecture of the neural network. Theorem 6.1 that if we are not careful about the choice of the
 121 first layer function of a network then we can always find a topologically labeled data that cannot be
 122 separated by this network.

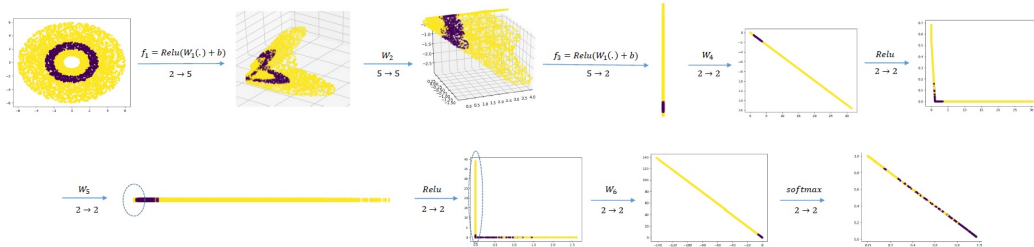
123 **Theorem 6.1.** *Let Net be neural network of the form $Net = Net_1 \circ f_1$ with $f_1 : \mathbb{R}^n \rightarrow \mathbb{R}^k$ such
 124 that $f_1(x) = \sigma(W(x) + b)$ and $k < n$ and $Net_1 : \mathbb{R}^k \rightarrow \mathbb{R}^d$ is an arbitrary net. Then there exists a
 125 topologically labeled data (D, h, g) with $h : D \rightarrow \mathcal{D} \subset \mathbb{R}^n$ and $g : \mathcal{D}_L \subset \mathcal{D} \rightarrow \mathbb{R}^d$ that is not separable
 126 by Net .*

127 **Proof.** Let $D = \mathcal{D} = \{x \in \mathbb{R}^n, \|x\| \leq 2\}$. Let $\mathcal{D}_L = \mathcal{D}_1 \cup \mathcal{D}_2$ where $\mathcal{D}_1 = \{x \in \mathbb{R}^n, \|x\| \leq 0.9\}$ and
 128 $\mathcal{D}_2 = \{x \in \mathbb{R}^n, 1 \leq \|x\| \leq 2\}$. Choose $g : \mathcal{D}_L \rightarrow \{l_1, l_2\}$ such that $g(\mathcal{D}_1) = l_1$ and $g(\mathcal{D}_2) = l_2$.
 129 Let f_1 be a function as defined in the Theorem. The matrix $W : \mathbb{R}^n \rightarrow \mathbb{R}^k$ where $k < n$ has a
 130 nontrivial kernel. Hence, there is a non-trivial vector $v \in \mathbb{R}^n$ such that $W(v) = 0$. Choose a point
 131 $p_1 \in \mathcal{D}_1$ and a point $p_2 \in \mathcal{D}_2$ on the line that passes through the origin and has the direction of v .
 132 We obtain $W(p_1) = W(p_2) = 0$. In other words, $f_1(p_1) = f_1(p_2)$. Hence $Net(p_1) = Net(p_2)$ and
 133 hence $Net(\mathcal{D}_1) \cap Net(\mathcal{D}_2) \neq \emptyset$ and so we cannot find two embedded disks that separate the sets
 134 $Net(\mathcal{D}_1), Net(\mathcal{D}_2)$.

135 Note that in Theorem 6.1 the statement is independent of the depth of the neural network. This is also
 136 related to the work [15] which shows that skinny neural networks are not universal approximators.
 137 This is also related to the work in [24] where it was shown that a network has to be wide enough in
 138 order to successfully classify the input data.

139 To demonstrate the role that the topology of data may play in regard to the architecture of a neural
 140 network we end our discussion by considering the following example. Let Net be a neural network
 141 given by the composition $Net = f_6 \circ f_5 \circ f_4 \circ f_3 \circ f_2 \circ f_1$. For $1 \leq i \leq 5$ maps are given by
 142 $f_i := Relu(W_i(x) + b_i)$ such that $W_1 : \mathbb{R}^2 \rightarrow \mathbb{R}^5, W_2 : \mathbb{R}^5 \rightarrow \mathbb{R}^5, W_3 : \mathbb{R}^5 \rightarrow \mathbb{R}^2$ and $W_j : \mathbb{R}^2 \rightarrow \mathbb{R}^2$
 143 for $4 \leq j \leq 5$. Finally, the function, $f_5 = softmax(W_6(x) + b_6)$ where $W_6 : \mathbb{R}^5 \rightarrow \mathbb{R}^2$.

144 We train this network on the annulus dataset given in the top left Figure in 1. In Figure 1 we
 145 also trace the activations as demonstrated in Figure 1. In the Figure we visualize the activations
 146 in higher dimension by projecting them using Isomap [27] to \mathbb{R}^3 . Our choice of this algorithm as
 147 a dimensionality reduction algorithm is driven by the fact that the dataset we work with here is
 148 essentially a manifold; as such, projecting the space to a lower dimension with the Isomap algorithm
 149 should preserve most of the topological and geometric structure of the this space.



150 **Figure 1:** The topological operations performed by a network on data sampled from the annulus and colored by two labels.

150 Inspecting the activations in Figure 1 we make the following observation:

- 151 1. A neural network can collapse the topological space either using the nonlinear $Relu$ or by
 152 utilizing the linear part of a given layer function. This is the case with the map $f_3 : \mathbb{R}^5 \rightarrow \mathbb{R}^2$.
 153 While the linear component is a projection onto \mathbb{R}^2 , the network chose "to project the space
 154 into 1-manifold since the second dimension is not needed for the final classification.
- 155 2. Note that the yellow components are separated by the purple one, and in order to map both
 156 of these parts to the same part of the space, the net has to glue these two parts together.
 157 Indeed, the neural network quotients parts of the space as it sees it necessary. This is visible
 158 in W_5 , which acts as a projection, and again W_6 .

159 Referencias

- 160 [1] M. Attene, S. Biasotti, and M. Spagnuolo. Shape understanding by contour-driven retiling. *The*
161 *Visual Computer*, 19(2):127–138, 2003.
- 162 [2] C. L. Bajaj, V. Pascucci, and D. R. Schikore. The contour spectrum. In *Proceedings of the 8th*
163 *Conference on Visualization '97*, pages 167–ff. IEEE Computer Society Press, 1997.
- 164 [3] R. L. Boyell and H. Ruston. Hybrid techniques for real-time radar simulation. In *Proceedings*
165 *of the November 12-14, 1963, fall joint computer conference*, pages 445–458. ACM, 1963.
- 166 [4] R. Brüel-Gabrielsson, B. J. Nelson, A. Dwaraknath, P. Skraba, L. J. Guibas, and G. Carlsson. A
167 topology layer for machine learning. *arXiv preprint arXiv:1905.12200*, 2019.
- 168 [5] G. Carlsson. Topology and data. *Bulletin of the American Mathematical Society*, 46(2):255–308,
169 2009.
- 170 [6] H. Carr, J. Snoeyink, and M. van de Panne. Simplifying flexible isosurfaces using local
171 geometric measures. In *IEEE Visualization*, pages 497–504, 2004.
- 172 [7] C. Curto. What can topology tell us about the neural code? *Bulletin of the American Mathema-*
173 *tical Society*, 54(1):63–78, 2017.
- 174 [8] G. Cybenko. Approximations by superpositions of a sigmoidal function. *Mathematics of*
175 *Control, Signals and Systems*, 2:183–192, 1989.
- 176 [9] Y. Dabaghian, F. Mémoli, L. Frank, and G. Carlsson. A topological paradigm for hippocampal
177 spatial map formation using persistent homology. *PLoS Computational Biology*, 8(8):e1002581,
178 2012.
- 179 [10] H. Edelsbrunner and J. Harer. *Computational topology: an introduction*. American Mathematical
180 Soc., 2010.
- 181 [11] C. Fefferman, S. Mitter, and H. Narayanan. Testing the manifold hypothesis. *Journal of the*
182 *American Mathematical Society*, 29(4):983–1049, 2016.
- 183 [12] C. Giusti, R. Ghrist, and D. S. Bassett. Two’s company, three (or more) is a simplex: Algebraic-
184 topological tools for understanding higher-order structure in neural data. *Journal of computati-*
185 *onal neuroscience*, 41:1, 2016.
- 186 [13] B. Hanin and M. Sellke. Approximating continuous functions by relu nets of minimal width.
187 *arXiv preprint arXiv:1710.11278*, 2017.
- 188 [14] C. Hofer, R. Kwitt, M. Niethammer, and A. Uhl. Deep learning with topological signatures. In
189 *Advances in Neural Information Processing Systems*, pages 1634–1644, 2017.
- 190 [15] J. Johnson. Deep, skinny neural networks are not universal approximators. *arXiv preprint*
191 *arXiv:1810.00393*, 2018.
- 192 [16] I. S. Kweon and T. Kanade. Extracting topographic terrain features from elevation maps.
193 *CVGIP: image understanding*, 59(2):171–182, 1994.
- 194 [17] H. Lee, M. K. Chung, H. Kang, B.-N. Kim, and D. S. Lee. Computing the shape of brain
195 networks using graph filtration and gromov-hausdorff metric. *International Conference on*
196 *Medical Image Computing and Computer Assisted Intervention*, pages 302–309, 2011.
- 197 [18] H. Lee, M. K. Chung, H. Kang, B.-N. Kim, and D. S. Lee. Discriminative persistent homology
198 of brain networks. *IEEE International Symposium on Biomedical Imaging: From Nano to*
199 *Macro*, pages 841–844, 2011.
- 200 [19] H. Lee, H. Kang, M. K. Chung, B.-N. Kim, and D. S. Lee. Persistent brain network homology
201 from the perspective of dendrogram. *IEEE Transactions on Medical Imaging*, 31(12):2267–
202 2277, 2012.
- 203 [20] H. Lee, H. Kang, M. K. Chung, B.-N. Kim, and D. S. Lee. Weighted functional brain network
204 modeling via network filtration. *NIPS Workshop on Algebraic Topology and Machine Learning*,
205 2012.
- 206 [21] N. Lei, D. An, Y. Guo, K. Su, S. Liu, Z. Luo, S.-T. Yau, and X. Gu. A geometric understanding
207 of deep learning. *Engineering*, 2020.
- 208 [22] Z. Lu, H. Pu, F. Wang, Z. Hu, and L. Wang. The expressive power of neural networks: A view
209 from the width. In *Advances in neural information processing systems*, pages 6231–6239, 2017.

- 210 [23] P. Lum, G. Singh, A. Lehman, T. Ishkanov, M. Vejdemo-Johansson, M. Alagappan, J. Carlsson,
211 and G. Carlsson. Extracting insights from the shape of complex data using topology. *Scientific*
212 *reports*, 3:1236, 2013.
- 213 [24] Q. Nguyen, M. C. Mukkamala, and M. Hein. Neural networks should be wide enough to learn
214 disconnected decision regions. *arXiv preprint arXiv:1803.00094*, 2018.
- 215 [25] M. Nicolau, A. J. Levine, and G. Carlsson. Topology based data analysis identifies a subgroup
216 of breast cancers with a unique mutational profile and excellent survival. *Proceedings of the*
217 *National Academy of Sciences*, 108(17):7265–7270, 2011.
- 218 [26] P. Rosen, B. Wang, A. Seth, B. Mills, A. Ginsburg, J. Kamenetzky, J. Kern, and C. R. Johnson.
219 Using contour trees in the analysis and visualization of radio astronomy data cubes. *arXiv*
220 *preprint arXiv:1704.04561*, 2017.
- 221 [27] J. B. Tenenbaum, V. De Silva, and J. C. Langford. A global geometric framework for nonlinear
222 dimensionality reduction. *science*, 290(5500):2319–2323, 2000.
- 223 [28] F. Wang, H. Liu, D. Samaras, and C. Chen. Topogan: A topology-aware generative adversarial
224 network.