Contents lists available at ScienceDirect

Computers & Graphics

journal homepage: www.elsevier.com/locate/cag

Vectorization of persistence barcode with applications in pattern classification of porous structures

Zhetong Dong^a, Chuanfeng Hu^a, Chi Zhou^a, Hongwei Lin^{a,b,*}

^aSchool of Mathematical Sciences, Zhejiang University, Hangzhou, Zhejiang Province, China ^bState Key Lab. of CAD&CG, Zhejiang University, Hangzhou, Zhejiang Province, China

ARTICLE INFO

Article history: Received May 3, 2020

Keywords: Computational Topology, Machine Learning, Persistent Homology, Porous Classification

ABSTRACT

Persistence barcode is a topological summary for persistent homology to exhibit topological features with different persistence. Persistence rank function (PRF), derived from persistence barcode, organizes persistence Betti numbers in the form of an integer-valued function. To obtain topological patterns of objects such as point clouds represented by finite-dimensional vectors for machine learning classification tasks, the vectorizing representations of barcodes is generated via decomposing PRF on a system of Haar basis. Theoretically, the generated vectorizing representation is proved to have 1-Wasserstein stability. In practice, to reduce training time and achieve better results, a technique of dimensionality reduction through out-of-sample mapping in supervised manifold learning is used to generate a low-dimensional vector. Experiments demonstrate that the representation of porous structures has become an essential problem in the fields such as material science in recent decades. The proposed method is successfully applied to distinguish porous structures on a novel data set of porous models.

1. Introduction

Persistent homology (PH) [1], an effective tool for obtaining topological features of spatial objects, has been used to discover topological patterns of practical data in recent decades. Topological patterns can be understood as representations re-5 lated to topological invariants, such as connected components in zero-dimension, loops in one-dimension, and spatial voids in two-dimension. Patterns represented by topology and geometry have been widely adopted in numerous studies, for instance, image processing [2], and network analysis [3]. Especially, in 10 computer graphics, the topological features of PH are used to 11 solve problems such as point cloud recognition [4], and surface 12 reconstruction [5]. 13

> *Corresponding author: *e-mail:* hwlin@zju.edu.cn (Hongwei Lin)

PH can be used to infer possible 'shape' of a point cloud in 14 Euclidean space. Assume that on each point, a ball with its cen-15 ter to be the point and the initial radius to be zero is assigned. 16 As the radius of each ball increases, simplexes are determined 17 by the intersections of these balls, and a simplicial complex is 18 generated corresponding to the value of the radius. In this dy-19 namical procedure, the homological invariants in different di-20 mension appear (or are born) at a moment when the value of 21 the radius is b, and it might disappear (or die) when some high-22 dimensional simplex appears at the value of d. Therefore, the 23 pair (b, d) is the life span of the topological invariants, and the 24 persistence is its lifetime d - b. One can obtain a set of these 25 pairs in each dimension. A persistence barcode [6, 7], shown in 26 Figure 1 as an instance, is ordered linear segments obtained by 27 embedding each pair into 2D Euclidean space. The endpoints 28 of each 'bar' represent the birth time b_i and the death time d_i 29 of the topological invariant, respectively. And a persistence di-30





2

agram is another topological summary embedding the set into 1 2D Euclidean space as points. Because the death time is always 2 greater than the birth time, the points are above the diagonal. To 3 measure the similarity of two barcodes or persistence diagrams, p-Wasserstein distance and the bottleneck distance are defined 5 in [8]. Readers can refer to [9] for more details. 6

In point cloud classification using topological features with 7 machine learning (ML) tools, the problem is how to transform 8 these persistence barcodes into representations compatible with 9 ML tools. Persistence barcode does not satisfy the multifold 10 requirements of practical applications in ML tasks such as clas-11 sification. It is because (1) 'bars' in a barcode are unordered, 12 and the number of 'bars' is not fixed; (2) to compute the dis-13 tance between two barcodes is not straightforward. Therefore, 14 to use topological features with ML tools, we need to identify a 15 vectorizing representation of barcodes. 16

The new representation of barcode should follow three prin-17 18 ciples: (1) the representation has an explainable meaning; (2) transformation should preserve the information which a bar-19 code contains, and (3) representation should be stable with re-20 spect to the metrics of barcode. The persistence rank func-21 tion (PRF), also known as persistence Betti number function 22 [10][11], is an ideal alternative to match these principles. First, 23 PRF is a bivariate and non-negative integer-valued function that 24 summarizes persistence Betti numbers (ranks) [12]. And the 25 PRF can be induced from barcode [13]. Second, no external 26 information is introduced to generate PRF. Third, it is shown in 27 Section 4 that PRF in L^2 in the sense of a particular measure has 28 1-Wasserstein stability in the condition of a small perturbation. 29 Consequently, we propose a framework to represent a barcode 30 as a stable finitely dimensional vector by vectorizing PRF. 31

As an application of the framework for point cloud classifi-32 cation, the vectorizing representation is used to classify porous 33 structures. Porous structure plays an important role in analyzing 34 the function of nanomaterials in material science [14]. In recen-35 t decades, as the databases of porous materials create, such as 36 the database in the material genome initiative [15], the classi-37 fication of porous materials has become a novel and important 38 problem. However, it is shown in [16] that traditional geomet-39 ric descriptors do not encode enough topological information 40 to detect materials that have similar global porous structures. 41 Therefore, persistent homology is adopted to capture the overall 42 porous features for quantifying similarity of nanoporous mate-43 rials. Inspired by this, we generate a data set of porous models 44 with category labels designed by triply periodic minimal sur-45 faces (TPMSs). On this data set, we extend the framework of 46 vectorizing barcode with the technique of dimensionality reduc-47 tion via out-of-sample mapping in supervised manifold learning 48 to speed up the training process of classifiers while maintaining 49 or even improving the classification accuracy. 50

Our contribution: In this paper, to vectorize persistence bar-51 codes for ML classification tasks of spatial point clouds using 52 topological features, we propose a finitely dimensional vector-53 izing representation of barcodes based on Haar basis decompo-54 sition of PRF in L^2 space with a limited domain, shown in Fig-55 ure 1. Theoretically, the generated vectorizing representation is 56 proved to be stable with respect to the 1-Wasserstein distance 57

in the condition of a small perturbation. For practical classi-58 fication tasks, because of the relatively high dimension of the 59 generated vector, a technique of dimensionality reduction via 60 out-of-sample mapping in supervised manifold learning is em-61 ployed to extend the vectorization framework. On a novel data 62 set of various porous materials, the proposed vectorizing repre-63 sentation is applied to classify the models, and it has the best 64 performance compared with other vectorizing methods. 65

2. Related Work

In this section, the vectorizing representations of barcode and persistence diagram for ML classification tasks are first introduced. Then, because the proposed framework is applied to classify the generated data set of porous structures, the generation approaches of 3D porous models are presented.

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

Vectorizing representations: From an historical perspective, size functions [17] were first used for the vectorization of 0-th homology, and the very first work about algebraic representations of 0-th persistence diagrams is in [18]. The vectorizing representations of a barcode or persistence diagram can be divided into non-learnable methods and kernel methods. There are numerous non-learnable representations for practical problems [19, 20, 21, 22, 23, 24]. For example, persistence images were proposed in [25] based on the integral on each mesh patch of a persistence surface, which is produced in the form of summation of weighted Gaussian functions related to the persistence of each point on a persistence diagram. To combine PH with statistics, a topological summary, called persistence landscapes, was proposed in [26]. This summary is a series of functions in a separable Banach space so that the vector space structure can be used to do statistics, such as computing mean values. Moreover, Robins et al. [13] considered PRF in a Hilbert space and in an affine subspace under reasonable conditions, and they performed functional principal component analysis on experimental data from colloids to study spatial point patterns.

Kernel methods are motivated by searching meaningful measurements to construct kernels used in machine learning models. Persistence scale space kernel (PSSK) was proposed in [27] as a multi-scale kernel by considering a heat diffusion problem. PSSK was proved to be 1-Wasserstein stable, while it does not have higher-degree Wasserstein stability. Kernel embedding of measures was used in [28] to transform a persistence diagram into an element of Hilbert space, and a framework of the k-100 ernel method, persistence weighted Gaussian kernel (PWGK) 101 was developed. Meanwhile, the kernel was proved to be stable 102 with respect to Hausdorff distance. Additionally, another kernel 103 method was proposed in [29], referred to as sliced Wasserstein 104 kernel (SWK), for machine learning tasks using persistence di-105 agrams. Readers can refer to the survey [30] for other recent 106 kernel methods. 107

Production of porous models: Cheah [31, 32] investigat-108 ed and selected various polyhedral shapes suitable for porous 109 structure modeling and created a parametric library of porous 110 structures. Schroeder [33] introduced stochastic geometry theo-111 ry into porous structure modeling to represent porous structures 112

27

28

20

30

31

32

33

34

35

36

44

with density and porosity. Cai and Xi [34] proposed a porous structure modeling method based on a shape function and hexahedral mesh refinement. An approach of irregular porous structure modeling based on subdivision and non-uniform rational B-splines was developed in [35]. You et al. [36] presented an improved method based on centroidal Voronoi tessellation and Bsplines to design porous structures. To overcome the limitation in geometry of pore-making element, easily control the porosity and pore size, and ensure the internal connectivity, several researchers used TPMSs to design porous structures [37, 38, 39]. 10 With the appealing properties of connectivity, smoothness, and 11 geometric representation, TPMS emerged as a significant tool 12 for designing porous structures. 13

14 **3. Vectorization of Barcode**

In this section, we introduce the vectorizing method of bar-15 code based on PRF, as shown in Figure 1. First, the approach of 16 inducing PRF from barcode and the Hilbert space where PRF is 17 considered are summarized and introduced. Second, we present 18 one of our contributions, that is, the method to generate the 19 feature vector by decomposing PRF using Haar basis on the 20 bounded domain. Finally, we mention that, because the feature 21 vector possesses a relatively high dimension, for tasks of clas-22 sification, a layer of dimension reduction is provided based on 23 out-of-sample mapping in supervised manifold learning. The 24 25 approach we employed is introduced in Section 5.



Fig. 1. Pipeline of generating low-dimensional feature vectors from barcodes: (1) Barcodes are computed from a filtration built up on the objects; (2) Extended persistence rank function (PRF) is induced from a barcode in a certain dimension; (3) A finite number of Haar basis functions, of which the first 16 functions are shown in the figure, are used to decompose the extended PRF on a bounded domain; (4) A feature vector with a relatively high dimension is generated by concatenating the coefficients of the Haar basis functions into a vector; (5) A layer of dimensionality reduction is used to produce low-dimensional vectors.

26 3.1. PRF and Functional Space

In real data sets, the number of bars in a barcode is finite, and most of the bars are distributed in a bounded domain. Therefore, it is reasonable to consider PRF within a limited domain. Given a *k*th persistence barcode $\{(b_i, d_i) | 0 \le b_i \le d_i, i \in I, |I| < \infty\}$, the *k*th PRF of a filtration *V* can be induced by the barcode and is defined as:

$$r_{k,V}(s,t) = \sum_{i \in I} r_i(s,t), \tag{1}$$

where

$$r_i(s,t) = \begin{cases} 1, & b_i \le s \le t \le d_i, \\ 0, & \text{otherwise.} \end{cases}$$
(2)

At the coordinate (s, t) where $s \le t$, the integer value is exactly the persistence Betti number, which, intuitively, means the number of *k*-dimensional 'holes' existing in the time interval [s, t]. And the function in Equation (2) draws a right triangle with its hypotenuse coincident with the diagonal.

We denote $r_{k,V}(s, t)$ as r(s, t) if no confusion occurs. The definition indicates that PRF is a non-negative integer-valued function, precisely, a piecewise constant valued function. Further, PRF is defined on the domain $\mathbb{R}^2_{\Delta+} := \{(s, t) \mid 0 \le s \le t\}$, which maps $\mathbb{R}^2_{\Delta+}$ into $\mathbb{Z} \cup 0 \subset \mathbb{R}$.

To evaluate the distance between two PRFs, the space $L^2(\mathbb{R}^2_{\Delta^+},\mu)$ is chosen here with standard inner product, where μ is a measure. To make a PRF be an element of the Hilbert space $L^2(\mathbb{R}^2_{\Delta^+},\mu)$, the measure μ is supposed to be carefully designed because there may be an infinite homology feature $[b,\infty)$ in a persistence barcode, which makes PRF not be finitely supported in $\mathbb{R}^2_{\Delta^+}$. In [13], the measure μ can be obtained by a weighted function $\phi(t-s)$, i.e., $d\mu = \phi(t-s)dsdt$, where t-s has the meaning of persistence of a homology class. It is necessary to show that PRFs are elements of $L^2(\mathbb{R}^2_{\Delta^+},\mu)$. For $r_{k,V}, r_{k,W}$ generated by filtrations V, W, respectively, suppose that V_{∞}, W_{∞} are finite simplicial complexes with $H_k(V_{\infty}) = H_k(W_{\infty})$, and it is proved in [13] that $||r_{k,V} - r_{k,W}||_2 < \infty$ if $\int_0^{+\infty} \phi(x)dx < \infty$. Simply, the weighted function $\phi(s, t)$ in the measure μ can be chosen to be

$$\phi(s,t) = \begin{cases} 1, & 0 \le s, t \le B, \\ 0, & \text{otherwise,} \end{cases}$$
(3)

where *B* is a finite number such that the distance is considered in the bounded domain $\Omega = [0, B] \times [0, B]$. The choice of *B* will be discussed in Section 5.

Finally, to extend the domain of PRF from $\mathbb{R}^2_{\Delta^+}$ to the domain $\mathbb{R}^2_+ = \{(s,t) \mid s,t \ge 0\}$ to make the function decomposed by a basis defined on \mathbb{R}^2_+ , the extended PRF is defined by making PRF be symmetric with respect to the diagonal, that is,

$$\widetilde{r}_{k,V}(s,t) = \sum_{i \in I} \widetilde{r}_i(s,t), \qquad (4)$$

where

$$\widetilde{r}_i(s,t) = \begin{cases} 1, & b_i \le s, t \le d_i, \\ 0, & \text{otherwise.} \end{cases}$$
(5)

Intuitively, an example of extended PRF is given in Figure 1, 40 and it draws several square steps along the diagonal. The extended PRF does not introduce any extra information. In the context, $\tilde{r}_{k,V}(s, t)$ is denoted as $\tilde{r}(s, t)$ if no confusion occurs. 43

3.2. Vector Generation by Haar Decomposition

Here, we present the method to generate vectorizing representations. To extract features of PRF with different scales, one idea is to decompose a PRF on a system of basis that captures both local and global characteristics. Note that extended PRF has a structure of square steps along the diagonal because extended PRF is a piecewise constant-valued bivariant function.

This nature indicates an appropriate decomposition based on a 1

system of non-continuous orthonormal basis, similar to Fouri-2 er decomposition. The system of Haar basis, a series of non-3

continuous orthonormal functions, is exactly a proper alterna-

tive to transform a PRF into a vector consisting of the coeffi-5

cients of the basis. Haar basis is a complete orthonormal basis in $L^{2}[0, 1]$ [40]. The definition of the system of Haar functions is given as follows:

$$\begin{aligned} & \text{har}_{0}(0,t) = 1, \quad 0 \leq t \leq 1, \\ & \text{har}_{n}(k,t) = \begin{cases} & \sqrt{2^{n-1}}, \quad \frac{2k-2}{2^{n}} \leq t < \frac{2k-1}{2^{n}}, \\ & -\sqrt{2^{n-1}}, \quad \frac{2k-1}{2^{n}} \leq t \leq \frac{2k}{2^{n}}, \\ & 0, \quad \text{otherwise}, \end{cases} \end{aligned}$$
(6)

where $n = 1, 2, \cdots, k = 1, 2, 3, \cdots, 2^{n-1}$. The Haar basis can be extended to 2D in the sense of tensor product. Without losing generality, set $\Omega = [0, 1]^2$ and arrange an order on the Haar functions so that $\{har_n(k, t)\}_{n,k}$ is denoted as $\{har_i(t)\}_i$ where i = $0, 1, 2, \dots, N$. The 2D Haar basis in the sense of tensor product is given by

$$har_i(s,t) = har_i(s)har_k(t), \qquad j,k = 1, 2, 3, \dots, N.$$
 (7)

The 2D Haar system is a complete standard orthonormal basis on $L^{2}[0,1]^{2}$. When n in Equation (6) is relatively small, the corresponding Haar functions extract features with a large scale. And the local features are extracted as *n* increases. When *n* is large enough, subtle features can be captured. Therefore, it is reasonable to approximate an extended PRF by a finite number of Haar basis functions denoted as $\{har_i(s, t)\}_{i=1}^{N^2}$ generated sequentially by Equation (7). We truncate extended PRF in the domain Ω , and normalize the domain to be $[0,1]^2$ so that extended PRF can be decomposed by Haar basis. The coefficients λ_i are obtained by

$$\lambda_i = \langle \widetilde{r}(s,t), \operatorname{har}_i(s,t) \rangle, \qquad i = 1, 2, \cdots, N^2.$$
(8)

And then, the coefficients λ_i are concatenated to be a vector in 7 \mathbb{R}^{N^2} , i.e., $v = (\lambda_1, \lambda_2, \cdots, \lambda_{N^2})$, called the *feature vector* based 8 on Haar basis decomposition. The algorithm to generate the 9 feature vector from extended PRF is given in Algorithm 1. In 10 the algorithm, the operator .* means to multiply each entry of 11 a matrix with the corresponding entry of the other matrix. The 12 double sum means to add up all entries of the computed matrix. 13 The time complexity of the algorithm is $O(2^{2(n+N)+1})$, where n 14 is given in Equation (6), and the interval [0, 1] is divided into 15 2^N shares. In practice, we set N = 8 (256 shares) and n = 516 (1024 basis functions) to fast compute the feature vector. 17

Note that it is inevitable to lose some subtle information of a 18 PRF, particularly the information close to the diagonal, because 19 a limited number of Haar basis functions are used to decompose 20 a PRF. Fortunately, this information does not represent promi-21 nent topological features. 22

Practically, feature vectors are sparse vectors with a rather 23 high dimension. It is because, intuitively, high-frequency in-24 formation is mainly concentrated near the diagonal, and a large 25

Algorithm 1: Generation of the Feature Vector

Input: A barcode $\{b_i, d_i\}_{i \in I}$, the bound *B*, an integer *n* for Haar basis functions, and an integer N to split [0, 1]into 2^N shares.

1. Generate the discrete extended PRF according to Equation (4) and (5), truncate it in $[0, B]^2$, and normalize the domain to be $[0, 1]^2$. The discrete extended PRF is stored in the form of matrix with its size $2^N \times 2^N$, denoted as z:

2. Generate 2^n discrete 1D Haar basis functions according to Equation (6), stored in a matrix har1D with its size $2^n \times 2^N$: Initiate $v = \operatorname{zeros}(2^{2n}, 1);$

Initiate index = 1;**for** i = 1 to 2^n **do for** i = 1 to 2^n **do** $har2D = har1D(i, :)^{T} * har1D(j, :);$ $v(\text{index}) = \left(\frac{1}{2^N}\right)^2 * \text{sum}(\text{sum}(\text{har2D}*z));$ index = index + 1;end end **Output**: The feature vector v.

number of high-frequency basis functions do not capture any information. Before feeding data to classifiers, one can use 27 the technique of out-of-sample mapping in supervised manifold 28 learning to obtain low dimensional vectors such that the training 29 time is saved. Therefore, a layer of dimensionality reduction is 30 adopted in the classification framework exhibited in Section 5. 31

32

39

44

4. Theoretical Guarantee: Stability

In this section, we theoretically prove the stability of the gen-33 erated vectorizing representation to tiny noise. Intuitively, the 34 core of the stability of the Haar feature vector shows that the 35 perturbation on the generated feature vector is corresponding-36 ly small if the perturbation on the source data is small. Due 37 to the stability of barcode in [8], in some conditions, it shows 38 that perturbations on the vector can be controlled by those on the barcode. In mathematical literature, for the norm of the d-40 ifference between the original feature vector and the one with 41 perturbation, there is an upper bound given by a distance be-42 tween the original barcode and the perturbed one. 43

4.1. Distance of Barcodes

To evaluate measure similarity between two barcodes, one of 45 alternative distance is *p*-Wasserstein distance, and the stability 46 result is also obtain under the Wasserstein distance for a reason-47 ably large class of function in [9], which guarantees the robust-48 ness of barcode in the presence of noise. Intuitively, to measure 49 the differences between two barcodes, which can be seen as the 50 sets of intervals, a matching is needed between them. There are 51 two cases for the intervals in a barcode: the matched intervals 52 and the others not being matched. 53

Given two barcodes denoted as Bc, Bc' together with index sets I, and I', respectively, let l_i be the interval in a barcode, i.e., $l_i := (b_i, d_i)$ and the differences of matched and non-matched intervals are given by

$$\begin{aligned} \|l_i - l_j\|_{\infty} &:= \max\{|b_i - b_j|, |d_i - d_j|\}, \ i \in I, \ j \in I, \\ \|l_i\|_{\infty} &:= \frac{d_i - b_i}{2}. \end{aligned}$$
(9)

A matching φ is a binary relation between *Bc* and *Bc'* such that any interval in Bc and Bc' matches at most one pair. A matching can be seen as a subset of $I \times I'$. Define M and M' to be the index set of matched intervals in Bc and Bc', respectively. The *p*-Wasserstein distance $W_p(Bc, Bc')$ between the two barcodes Bc and Bc', which intuitively summarizes all the differences under the matching that makes the differences minimal, is defined as 、1/n

$$\inf_{\varphi} \left\{ \sum_{i \in M} \|l_i - l'_{\varphi(i)}\|_{\infty}^p + \sum_{j \in I \setminus M} \|l_j\|_{\infty}^p + \sum_{k \in I' \setminus M'} \|l'_k\|_{\infty}^p \right\}^{1/p}, \quad (10)$$

where *p* is a positive real number.

4.2. 1-Wasserstein Stability 2

At the beginning of the proof of stability, some assumptions are proposed: (1) the number of intervals (b_i, d_i) is limited, i.e., $|I| \leq L$ where L is a positive integer; (2) the domain Ω is 5 bounded, i.e., $\Omega = [0, B] \times [0, B]$, where B is given by a finite constant; (3) the perturbation on the barcode is small.

Given a barcode Bc, the corresponding PRF r(s, t) and feature vector $v \in \mathbb{R}^{N \times N}$, a small perturbation is exerted on the barcode to produce the barcode Bc', PRF r'(s, t) and feature 10 vector v'. The *p*-Wasserstein distance of *Bc* and *Bc'* is given 11 by $W_p(Bc, Bc')$. Therefore, suppose that φ is the exact match-12 ing map which reaches the infimum of $W_p(Bc, Bc')$ in Equation 13 (10). Because the perturbation is subtle according to assump-14 tion (3), that is, one can assume $W_1(Bc, Bc') \leq 1$. The following 15 theorem holds. 16

Theorem 4.1. Given a barcode Bc and the corresponding perturbed barcode Bc', r(s,t) and r'(s,t) are the PRFs in $L^{p}(\Omega)$, respectively. If $W_1(Bc, Bc') \leq 1$, then

$$\|r - r'\|_{p}^{p} \le 2^{(p-1)(2L-1)}(2B+1)W_{1}(Bc, Bc'),$$
(11)

where $B < \infty$ is given by the bounded domain $\Omega = [0, B] \times$ [0, B], and L is the largest number of intervals contained in a 18 barcode. 19

For extended PRF, there is a corollary from Theorem 4.1. 21

Corollary 1. In the same conditions as Theorem 4.1 given, $\tilde{r}(s,t)$ and $\tilde{r}(s,t)$ are the extended PRFs in $L^{p}(\Omega)$. Then we have

$$\|\widetilde{r} - \widetilde{r}'\|_{p}^{p} \le 2^{(p-1)(2L-1)+p}(2B+1)W_{1}(Bc, Bc'), \quad (12)$$

- where the notations L and B follow Theorem 4.1.
- Proof. See Appendix.

And then the 1-Wasserstein stability of the feature vector is 24 true by the following theorem. 25

Theorem 4.2 (Stability of Feature Vector). Given a barcode Bc and the corresponding perturbed barcode Bc', v and v' are the feature vectors generated by decomposing the extended PRFs, produced by Bc and Bc', on a finite series of Haar basis $\{har_i(s,t)\}_{i=1}^{N^2}$, respectively. If $W_1(Bc, Bc') \leq 1$, then

$$\|v - v'\|_2^2 \le 2^{2L+1}(2B+1)W_1(Bc, Bc'), \tag{13}$$

where $B < \infty$ is given by the bounded domain $\Omega = [0, B] \times$ 26 [0, B], and L is the largest number of intervals contained in a 27 barcode. 28

Proof. See Appendix. 29

5. Experiments

In this section, we exhibit the technique of dimensionality 31 reduction (DR) through out-of-sample mapping in supervised 32 manifold learning practically used on data sets. And then, we discuss the determination of some parameters in generating the 34 feature vectors, i.e., the bound B in Equation (3), and the appropriate dimension for the DR layer. Then classification experi-36 ments were done in a data set of random images of Brownian 37 motion [41], which can be seen as scalar fields on the grid, a 38 data set of a 2D dynamical system [25] and nine time-series 39 data sets of multi-source signals [42], which are essentially the 40 point clouds embedded in Euclidean space. Accuracy on differ-41 ent classifiers was obtained, and the performance of the feature 42 vectors was compared with state-of-the-art kernel methods and 43 persistence images (PIs). Finally, to clarify the effect of vec-44 torizing process and the layer of DR, the ablation study was 45 conducted. In the experiments, the barcodes of Vietoris-Rips 46 filtration were computed via the Python package Ripser [43]. 47 And Python package dionysus 2 [44] was used to compute the 1-Wasserstein distance between two barcodes. 49

5.1. A Layer of Dimensionality Reduction

To reduce time of training classifiers and preserve or even im-51 prove classification accuracy, a dimensionality reduction technique based on supervised manifold learning is adopted to deal with high dimensional data. We employ the method proposed in [45], which is multi-output kernel ridge regression for outof-sample mapping in supervised manifold learning.

Before choosing the DR method of out-of-sample mapping 57 in supervised manifold learning, we attempted commonly used 58 unsupervised DR methods, such as multidimensional scaling 59 [46], and unsupervised manifold learning methods, like isomet-60 ric feature mapping [47], but found the classification accuracy 61 after DR dropped a lot, compared with that without DR. There-62 fore, we attempted supervised DR method based on manifold 63 learning technique, kernel ridge regression for out-of-sample 64 mapping proposed in [45], in which the category labels are used 65 to improve the embedding of the training points. 66

The supervised approach of manifold learning utilizes the la-67 bels of training data set to compute a projection of training data 68

30

48

50

52

53

54

55

for the subsequent classification task. To increase inter-class 1 dissimilarity and to decrease intra-class dissimilarity, a super-2 vised variant of Isomap based on a hierarchical agglomeration 3 of the components is obtained from training labels. Further-4 more, to evaluate the classification accuracy of a classifier on 5 a test data set, an out-of-sample embedding method based on 6 multi-output kernel ridge regression is used, which projects the 7 test data into the embedding space. Note that before classify-8 ing, the class labels of the test data remain hidden. More details a are introduced in [45]. 10

11 5.2. Evaluation

Determination of Parameters: In practice, the topological 12 features of PH distribute in a bounded domain. PRF is con-13 sidered as an element in L^2 space by choosing the weighted 14 function $\phi(s, t)$ shown in Equation (3). Equivalently, PRF is 15 limited in the domain $\Omega = [0, B]^2$. The problem arises to se-16 lect an appropriate bound value B in the classification tasks. A 17 reasonable alternative to determine the bound B is to find the 18 maximum of death indices of barcodes of training data, i.e., 19 $B = \max_{j \in train} \max_{i \in I_i} (d_i^j)$, and then to truncate bars with B 20 in barcodes of test data. However, for few data sets in which 21 a certain bar emerges with a relatively large death index, re-22 garded as an outlier, the bound B is computed unreasonably 23 large, which may dilute detailed features that Haar basis can 24 capture. To avoid this, one can roughly assume that the se-25 quence of birth and death indices follows a certain distribu-26 tion, such as a normal distribution or an exponential distribu-27 tion. Precisely, let $\mathbf{x} = (b_1, d_1, \dots, b_i, d_i, \dots)$ be a sequence. Fit 28 the sequence **x** by a distribution $\varphi(t)$, and B is determined by 29 $\inf_{y} \{y : \int_{-\infty}^{y} \varphi(t) dt \ge \delta\}$, where δ is a constant in [0, 1]. In this 30 way, one can obtain a reasonable value of bound B. Note that 31 it is an alternative to determine a reasonable value of bound, 32 33 not a rigorous assumption that the sequence follows the distribution. By observing experimental barcode data, we computed 34 the bound value B using exponential distributions with $\delta = 0.99$ 35 on a few data sets, and on most of data sets, B was determined 36 by the maximum of death indices of barcodes. To test the ef-37 fect of the selection of bound B on classification accuracy, we 38 changed the bound in the parameter interval [0.9B, B], where 39 B was determined by the maximum of death indices, and the 40 classification accuracy fluctuates in a small range, about 0.2% 41 to 2.0%. If some of topological features are truncated, the clas-42 sification accuracy will be affected. And when the bound is set 43 in the interval [B, 1.1B], then the classification accuracy has no 44 significant change. 45

Because feature vectors are sparse, the basic idea is to look 46 for a dimension in which the information of the vectors can 47 be preserved as much as possible. To determine the appropri-48 49 ate dimension before the process of dimensionality reduction, principal component analysis (PCA) based on singular value 50 decomposition is used. Concretely, in the implementation of 51 PCA, when singular value decomposition is performed on the 52 covariance matrix, the eigenvalue matrix Λ is obtained. Pro-53 jection error which measures loss of information can be esti-54 mated by $\varepsilon(k) = 1 - \frac{\sum_{i=1}^{k} \lambda_i}{\sum_{i=1}^{n} \lambda_i}$, where $\{\lambda_i\}$ is a sequence of eigenvalues with descending order and *k* is the reduced dimension. 55 56

Let σ be tolerance. The appropriate dimension k_{app} is given by inf $\{k \in \mathbb{Z}^+ : \varepsilon(k) \le \sigma, k \le n\}$. With k_{app} as an estimated reduced dimension, the DR layer introduced in Section 5.1 maps the labeled vectors to a low dimensional vector space. In our experiments, σ was set to be 0.01.



■ DR+Kernel SVM ■ DR+Linear SVM ■ DR+LR ■ DR+kNN (a) Classification accuracy of low-dimensional feature vectors of PRF on classifiers of kernel SVM, linear SVM, LR, and *k*NN, respectively.



■ PRF+SVM ■ PRF+DR+SVM ■ PWGK ■ PSSK ■ SWK ■ PI+SVM (b) The classification accuracy of feature vectors with and without DR on kernel SVM classifier, and classification accuracy of PWGK, PSSK, SWK and PI, each of which is adjusted to have the best performance.



Comparison: To evaluate the performance of feature vec-62 tors for extracting topological patterns, classification tasks were 63 performed on eleven data sets. To verify that the vectors can 64 catch the topological patterns of data sets, state-of-the-art ker-65 nel methods for persistence diagrams, including PWGK, PSSK, 66 and SWK, and PIs were used to compare with the feature vec-67 tors. To obtain barcodes, the H_0 barcode was computed based 68 on lower-star filtration of the scalar function on a data set of 69 random images. For a time series of multi-source signals, time-70

58 59 60

61

delay embedding in [48] was adopted to transform a time series into a point cloud in Euclidean space, such that Vietoris-Rips (V-R) filtration was generated to compute H_1 barcodes. In the data set of a 2D dynamical system, V-R filtration with respect to the 2D Euclidean metric was produced to obtain H_1 barcodes.

In the pipeline of generating the feature vectors of PRF, 2^{10} 6 Haar basis functions are used to extract the coefficients, i.e., n = 5 in Equation (6). Meanwhile, because the feature vectors 8 are suitable to be applied on different classifiers, four classifiers, consisting of k nearest neighbors (kNN), logistic regression 10 (LR), kernel support vector machine (kernel SVM) and linear 11 support vector machine (linear SVM), are used to form a layer 12 for classifying the input feature vectors. The parameters of L-13 R are defaulted, while in kNN, linear SVM, and kernel SVM, 14 the hyperparameters are selected by cross-validation to ensure 15 the best performance. Furthermore, in the process of classifica-16 tion using kernel methods, the hyperparameters of kernels are 17 adjusted to reach the best classification performance. For PIs, 18 images are produced with their resolution 10×10 and Gaussian 19 $\sigma = 10^{-3}$. 20

For the performance of feature vectors on the classifiers k-21 ernel SVM, linear SVM, LR, and kNN, the results are shown 22 in Figure 2(a). The low dimensional feature vectors obtained 23 by DR have good performance on the classifiers of kernel SVM 24 for most of data sets (8 in 11), while LR and kNN may not dis-25 tinguish those low dimensional vectors well on some data sets. 26 In Figure 2(b), the comparison results, i.e., the performance of 27 the proposed vectorizing representations with and without DR 28 on kernel SVM classifier, the performance of kernel methods 29 for barcode (PWGK, PSSK, and SWK), and the performance 30 of PIs are shown. On the time-series data sets Beef, CBF [42] 31 and random images, the kernel methods perform best. PIs have 32 the best classification result on the data set of 2D dynamical 33 system. On the rest of data sets (7 in 11), feature vectors of 34 PRF (without DR, shown in orange) perform the best. Overal-35 1, feature vectors can capture prominent topological features of 36 data such that a nice performance is achieved on the classifiers. 37 And, for the comparison of time cost of the proposed method 38 and its competitors, it costs less time to produce feature vectors 39 for classification than kernel methods. The proposed classifica-40 tion pipeline costs about 2 hours, but kernel methods cost about 41 15 hours in total eleven data sets for all procedures. Techniques 42 of parallel computing can be adopted to compute the vectors of 43 a collection of data. 44

Furthermore, we notice that on most of data sets, the perfor-45 mance of the feature vectors without DR on kernel SVM clas-46 sifier is better than that of the vectors with DR. To clarify the 47 effect of DR layer in both aspects of classification accuracy on 48 the rest of classifiers and the time cost of training classifiers 49 with low and high-dimensional vectors, we do ablation study by 50 removing the DR layer. Moreover, the entire framework will be 51 ablated by using the distance matrix with respect to Wasserstein 52 distance equipped on barcodes as the input of kNN classifier in 53 order to show the effect of classification by using the proposed 54 vectoring representations of barcodes. 55



Fig. 3. Classification results for ablation study on eleven data sets: classification accuracies of using feature vectors with DR, feature vectors without DR on LR and *k*NN classifiers, and a distance matrix of 1-Wasserstein distance on *k*NN classifier are illustrated.

5.3. Ablation Study

To investigate the effects of these two components in the process, an ablation study was done. First, the layer of DR is ablated to show its effectiveness for classifying data with different labels. Then, the whole layer of vectorization is ablated to validate the effectiveness of feature vectors for classification tasks. The classification after ablation is implemented by *k*NN with the distance matrix with respect to the 1-Wasserstein distance of barcode.

The purpose of DR is to save time of classification and to preserve and even improve classification accuracy. In the experiments, the appropriate reduced dimension was estimated via PCA. As shown in Figure 3, with the layer of DR, the classification accuracies of eight data sets on kNN and LR classifier are improved. For instance, noticeable improvement occurs on the data sets of ECG200 and of 2D dynamical system. This means that it is feasible to improve the accuracy of classification in some data sets as well as to obtain dense low-dimensional feature vectors.

Table 1. Total time cost of classification using feature vectors with and without DR: Four data sets with a relatively large training and test set are chosen to measure the total time cost of classifying the vectors with reduced dimension and the vectors with 1024 dimensions on four classifiers. The time of selecting optimal hyperparameters is included.

Data set	Training/test data	Reduced dimension	Time cost (s) with or without DR
FordB	3636/810	34	102/891
Chlorine Concentration	467/3840	23	15/143
Distal PhalanxTW	600/276	29	5/24
Random Image	800/160	4	5/89

56

65

66

67

68

69

70

71

72

73



Fig. 4. Pipeline of classifying data set of porous models using feature vectors with the DR technique: an example of nine types of TPMS porous models are exhibited in the column *training data*, and an example of test data sampled from a porous model of bone is shown in the column *test data*.

out DR on kNN classifiers and the method of directly classifying barcodes using 1-Wasserstein distance matrix, Figure 3 also 2 shows the improvement of the proposed method with and with-3 out DR on most of data sets for the extraction of topological 4 patterns. On eight data sets, it improves the classification accu-5 racy by using the proposed vectorizing representation and the 6 DR laver. On five data sets such as ECG200 and Earthquakes. 7 the classification accuracy improves by using vectorizing representation in the classification on kNN classifier. However, 9 there exist a few data sets, such as Beef, on which using 1-10 Wasserstein distance matrix on kNN classifier has a better per-11 formance. However, it spent much more time to compute 1-12 Wasserstein distance even through the efficient method of using 13 geometry of persistence diagrams proposed in [49]. 14

Although, as shown by bars in orange and in red in Figure 15 2(b), there is some loss of accuracy by using DR technique on 16 SVM classifier, the advantage of DR layer is to reduce the train-17 ing time of the classifiers, especially when the data set is large. 18 Apparently, it costs much more time to train the classifiers with 19 high dimensional vectors. As shown in Table 1, four data sets 20 with a large training and test set are chosen to show the time 21 saving after DR process. And the time cost was measured on a 22 PC with Intel(R) Core(TM) i7-4790 CPU@3.60GHz×8. After 23 24 DR process, the total time cost of classification is reduced by 5 to 15 times. 25

6. Application: Classification of Porous Structures

Porous scaffolds are widely used to engineer various human
tissues virtually in tissue engineering and biomaterials [39]. To
extract informative patterns from porous structures for the tasks
of classification is an important and novel problem for material design. In this section, the generating method of porous
structures, the details of training, test data and the approach of
classification are introduced. Finally, the classification accu-

racy and the time cost of the proposed framework and other state-of-the-art methods are analyzed.

34

35

36

37

38

39

40

41

6.1. Porous Data Set

In this application, a data set of 3D porous models was produced by using TPMS to generate nine types of different porous structures. TPMS is a minimal surface with a mean curvature of zero, with periodicity in each direction of 3D space, which is popular for engineering porous models.

The Generation of TPMS in B-Spline Solid: We approximated the TPMS using a periodic nodal surface defined by a Fourier series [50],

$$\psi(\mathbf{r}) = \sum_{k} A_{k} cos[2\pi (\mathbf{h}_{k} \cdot \mathbf{r})/\lambda_{k} - P_{k}] = C, \qquad (14)$$

where **r** is the location vector in the Euclidean space, A_k is the 42 amplitude, \mathbf{h}_k is the k^{th} lattice vector in the reciprocal space, 43 λ_k is the wavelength of the period, P_k is the phase shift, and C 44 is the threshold constant. And we set C = 0 in this applica-45 tion. In this application, marching tetrahedra (MT) algorithm 46 is employed to extract the TPMS. In consideration of accura-47 cy and storage of TPMS, the physical domain is divided into 48 $100 \times 100 \times 100$ hexahedrons, and each hexahedron is further 49 divided into 6 tetrahedrons. And then all of intersection trian-50 gles constitute a mesh approximating the iso-surface. 51

To generate a TPMS in the 3D mesh model, the B-spline sol-52 id needs to be approximated by hexahedral mesh model. We 53 use a set of sampling points along each parametric direction ac-54 cording to the preset resolution. All these sampling points are 55 calculated and mapped into the Cartesian space by the trivari-56 ate B-spline function. These sampling evaluation points con-57 stitute numerous hexahedron elements to approximately repre-58 sent the B-spline solid. Every direction of parametric domain 59 is evenly sampled to generate parametric coordinates of eval-60 uation points. And then, the Cartesian coordinates of the 3D 61 mesh vertices can be calculated by the trivariate B-spline func-62 tion. The B-spline solid can be approximately represented by

hexahedron elements. Therefore, each hexahedron is further 1 divided into 6 tetrahedrons, and the TPMS can be extracted by

the MT algorithm mentioned above.

Training Data: The training data are produced by the following operations. First, the physical domain is set to $[x, x + 2] \times$ $[y, y + 2] \times [z, z + 2]$, and the variables x, y, z are sampled from [-1, 1] using random sampling algorithm. Then, with different parameters in Equation (14), the discrete hexahedral physical domain and the iso-surface $\psi(x, y, z) = 0$ are determined, and the nine types of TPMS can be extracted by the MT algorithm. 10 Finally, 100 models with its size $2 \times 2 \times 2$ are generated for each 11 TPMS type. For each model, 1000 points are randomly sam-12 pled twice to produce a point cloud. In total, 1800 point clouds 13 are generated as training data with labels. 14

Test Data: As for the test data, five porous models with nine 15 types of TPMS are produced, namely, bone (balljoint), isis, 16 moai, tooth, and venus (45 models in total). From each of the 17 models, three cubes with their size $2 \times 2 \times 2$ are sampled. The 18 samples do not contain the boundary of the models. And for 19 each cube, three point clouds are randomly sampled with 1000 20 points. Therefore, 405 point clouds are obtained as test data. 21 The labels (TPMS types) keep hidden before classification. 22

6.2. Classification and Results 23

The pipeline of classifying porous structures is given in Fig-24 ure 4. For each point cloud in both training and test data set, 25 the V-R filtration is built on each point cloud in 3D Euclidean 26 space to compute a 1-dimensional persistence barcode via PH. 27 Then, the vectorizing representation based on PRF is comput-28 ed to obtain the feature vector. And alternatively, the layer of 29 DR can be used to obtain the low-dimensional vector. Finally, 30 the classifiers are trained with cross-validation and classifica-31 tion accuracy is obtained. 32

Table 2. Classification accuracy and time cost on the data set of porous models: 9-dimensional vectors were obtained through DR, and the hyperparameter C of linear SVM was selected to be 0.01 via cross validation. PIs was obtained by setting the resolution of 10×10 and $\sigma = 10^{-3}$ to achieve the best performance. Time cost for our methods and PIs consists of vectorization, training classifier, and testing. Time cost for 1-Wasserstein+kNN consists of computing distance matrix, training classifier, and testing. And time cost of kernel methods consists of training and testing.

Methods	Accuracy	Time Cost
PRF+DR+Linear SVM	81.6%	5.2min/3.0s/0.1s
PRF+Linear SVM	79.4%	5.2min/12.2s/0.2s
PWGK	64.9%	1.5h/0.6h
PSSK	71.6%	2.7h/1.1h
SWK	54.7%	2.5h/1.0h
1-Wasserstein+kNN	72.6%	12.0h/13.5s/0.2s
PIs+kernel SVM	66.4%	4.9min/5.6s/0.1s

As shown in Table 2, feature vectors with DR technique on 33 the linear SVM classifier have the best performance. The layer 34 of DR helps improve classification accuracy, and 9-dimensional 35 vectors are obtained to represent topological patterns of porous 36 models. Kernel methods and PIs do not perform well because 37

it is likely to be unable to capture detailed topological features 38 but to be disturbed by subtle topological noises in the test data. For time cost, the calculation of the proposed feature vectors is 40 nearly as efficient as the generation of PIs. And we mention 41 that it costs 1.0 min in DR procedure in our method. The kernel 42 methods cost a large mount of time to train the classifier, and 43 it is unacceptable to use the traditional 1-Wasserstein distance 44 on kNN classifier for this classification task. To eliminate the 45 concern that the nice performance of feature vectors with DR 46 occurs because the decomposition based on a few number of 47 Haar basis removes the topological noise of data, 2¹² Haar basis 48 functions (n = 6 in Equation 6) are adopted, compared with 2^{10} 49 functions (n = 5). The classification accuracy is 80.6% on the 50 classifier of linear SVM, which remains high. In general, the 51 results on data sets of porous models show the potential to use 52 PRF to extract topological patterns on data sets with prominent 53 topological features for classification tasks. 54

7. Discussion and Conclusion

This paper introduces a novel vectorizing representation of PRF based on Haar basis decomposition on a bounded domain to extract topological patterns for the classification tasks of point clouds, which is one of the issues of interest to the computer graphics community. The generated vectorizing represen-60 tation is proved to have 1-Wasserstein stability, which provides theoretical guarantee of the proposed method to deal with data with noise. Classification experiments in different data sets show its effectiveness. Meanwhile, it is shown that on the data set of porous models with topological noises, feature vectors of PRF perform the best. In the classification pipeline, we employed the DR technique of out-of-sample mapping in supervised manifold learning to reduce time of training classfiers and preserve or even improve classification accuracy in practice. It is more effective than the unsupervised DR methods we attempted. One can also attempt other DR methods that match the practical needs for specific data sets.

In essence, the proposed vectorizing representation of bar-73 codes is used for classification tasks of point clouds, and the 74 classification of porous models exhibited in Section 6 is essen-75 tially a task of 3D point cloud classification. In deep learning, 76 the designed networks, such as PointNet [51], usually have s-77 tunning performance for point cloud classification, especially 78 in the scenarios involving model semantics. However, the in-79 terpretability of the networks is still an urgent issue. In ma-80 chine learning, kernel methods based on persistence diagram-81 s, such as those we used to compare with the proposed meth-82 ods, transform the classification of persistence diagrams into a 83 classic non-linear classification problem by designing kernels. 84 These methods have strong generalization ability. But finite-85 dimensional vectors can not be explicitly given, and it costs a large amount of time in practice to train the classifiers. As a 87 hand-crafted feature extraction, the proposed method captures 88 the informative features of PRF, and can be clearly explained. 89 Meanwhile, the proposed method is shown to have good per-90 formance on the data sets with significant topological features, 91 such as porous model data set. Unlike kernel methods, the fea-92

55

56

57

58

59

62

63

64

65

66

67

68

69

70

71

/Computers & Graphics (2020)

ture vectors computed by the proposed method are directly in-1 put into classifiers, and it may reduce the generalizing power of 2

the classification pipeline. 3

For future work, the pore categories on real-world porous material data sets are usually unknown. The proposed method 5 transforms a point cloud of porous material into a finite-6 dimensional vector containing its topological information, such 7 as connectivity and loops. The generated vectors can be used 8 as a topological descriptor for material retrieval and classifica-9 tion. It will help researchers find porous materials with similar 10 topology in application scenarios. 11

Acknowledgments 12

This work is supported by the National Natural Science 13 Foundation of China under Grant Nos. 61872316, 61932018. 14

References 15

33

34

35

36

37

38

- [1] Edelsbrunner, H, Letscher, D, Zomorodian, A. Topological persistence 16 and simplification. Foundations of Computer Science 2000;28(4):454-17 463 18
- Robins, V, Wood, PJ, Sheppard, A. Theory and algorithms for construct-19 [2] ing discrete morse complexes from grayscale digital images. IEEE Trans-20 21 actions on Pattern Analysis and Machine Intelligence 2011;33(8):1646-1658. 22
- [3] Rieck, B, Fugacci, U, Lukasczyk, J, Leitte, H. Clique commu-23 nity persistence: A topological visual analysis approach for complex 24 networks. IEEE Transactions on Visualization and Computer Graphics 25 2018:24(1):822-831. 26
- Carlsson, GE. Topological pattern recognition for point cloud data. Acta 27 [4] Numerica 2014;23:289-368. 28
- Gabrielsson, RB, Ganapathisubramanian, V, Skraba, P, Guibas, LJ. [5] 29 Topology-aware surface reconstruction for point clouds. arXiv: Compu-30 31 tational Geometry 2018:.
- Collins, AD, Zomorodian, A, Carlsson, GE, Guibas, LJ. A bar-[6] 32 code shape descriptor for curve point cloud data. Computers & Graphics 2004:28(6):881-894.
 - Ghrist, R. Barcodes: The persistent topology of data. Bulletin of the [7] American Mathematical Society 2008;45(1):61-75.
 - [8] Cohen-Steiner, D, Edelsbrunner, H, Harer, J. Stability of persistence diagrams. Discrete & Computational Geometry 2007;37(1):103-120.
- Edelsbrunner, H, Harer, J. Computational Topology: An Introduction. [9] 39 40 American Mathematical Soc.; 2010.
- Cerri, A. Fabio, BD. Ferri, M. Frosini, P. Landi, C. Betti numbers in [10] 41 multidimensional persistent homology are stable functions. Mathematical 42 Methods in The Applied Sciences 2013;36(12):1543-1557. 43
- 44 [11] Frosini, P, Landi, C. Persistent betti numbers for a noise toleran-45 t shape-based approach to image retrieval. Pattern Recognition Letters 2013;34(8):863-872 46
- 47 [12] Zomorodian, A, Carlsson, GE. Computing persistent homology. Discrete & Computational Geometry 2005;33(2):249-274. 48
- Robins, V, Turner, K. Principal component analysis of persistent ho-49 [13] 50 mology rank functions with case studies of spatial point patterns, sphere packing and colloids. Physica D: Nonlinear Phenomena 2016;334:99-51 117. 52
- [14] Simon, CM, Kim, J, Lin, L, Martin, RL, Haranczyk, M, Smit, B. 53 Optimizing nanoporous materials for gas storage. Physical Chemistry 54 Chemical Physics 2014;16(12):5499-5513. 55
- Kalil, T, Wadia, C. Materials genome initiative for global competitive-56 ness. https://www.mgi.gov/;2011. 57
- [16] Lee, Y, Barthel, S, Dlotko, P, Moosavi, SM, Hess, K, Smit, B. Quanti-58 fying similarity of pore-geometry in nanoporous materials. Nature Com-59 munications 2017;8(1):15396. 60
- [17] Donatini, P, Frosini, P, Lovato, A. Size functions for signature recogni-61 tion. Proceedings of Spie the International Society for Optical Engineer-62 63 ing 1998;3454.

- [18] Frosini, P, Landi, C. Size functions and formal series. Applicable Algebra in Engineering, Communication and Computing 2001;12(4):327-349.
- MK, Johnson, SC, Singh, [19] Pachauri. D, Hinrichs, C, Chung, Topology-based kernels with application to inference problems in Alzheimer's disease. IEEE Transactions on Medical Imaging 2011;30(10):1760-1770.
- [20] Carriere, M, Oudot, SY, Ovsjanikov, M. Stable topological signatures for points on 3D shapes. In: Computer Graphics Forum; vol. 34. Wiley Online Library; 2015, p. 1-12.
- [21] Di Fabio, B, Ferri, M. Comparing persistence diagrams through complex vectors. In: International Conference on Image Analysis and Processing. Springer; 2015, p. 294-305.
- [22] Adcock, A, Carlsson, E, Carlsson, G. The ring of algebraic functions on persistence bar codes. Homology, Homotopy and Applications 2016:18(1):381-402
- [23] Kališnik, S. Tropical coordinates on the space of persistence barcodes. Foundations of Computational Mathematics 2018;:1-29.
- [24] Dong, Z, Lin, H, Zhou, C. Persistence B-spline grids: Stable vector representation of persistence diagrams based on data fitting. arXiv preprint arXiv:190908417 2019:.
- [25] Adams, H, Emerson, T, Kirby, M, Neville, R, Peterson, C, Shipman, P, et al. Persistence images: A stable vector representation of persistent homology. Journal of Machine Learning Research 2017;18(1):218-252.
- [26] Bubenik, P. Statistical topological data analysis using persistence landscapes. Journal of Machine Learning Research 2015;16(1):77-102.
- [27] Reininghaus, J, Huber, S, Bauer, U, Kwitt, R. A stable multi-scale kernel for topological machine learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2015, p. 4741-4748
- [28] Kusano, G, Fukumizu, K, Hiraoka, Y. Kernel method for persistence diagrams via kernel embedding and weight factor. Journal of Machine Learning Research 2018;18(189):1-41.
- [29] Carriere, M, Cuturi, M, Oudot, S. Sliced wasserstein kernel for persistence diagrams. In: Proceedings of the 34th International Conference on Machine Learning-Volume 70. JMLR. org; 2017, p. 664-673.
- [30] Pun, CS, Xia, K, Lee, SX. Persistent-homology-based machine learning and its applications - a survey. arXiv: Algebraic Topology 2018;.
- [31] Cheah, CM, Chua, C, Leong, K, Chua, S. Development of a tissue engineering scaffold structure library for rapid prototyping. part 1: Investigation and classification. International Journal of Advanced Manufacturing Technology 2003;21(4):291-301.
- [32] Cheah, CM, Chua, C, Leong, K, Chua, S. Development of a tissue engineering scaffold structure library for rapid prototyping. part 2: Parametric library and assembly program. International Journal of Advanced Manufacturing Technology 2003;21(4):302-312.
- [33] Schroeder, C, Regli, WC, Shokoufandeh, A, Sun, W. Computer-aided design of porous artifacts 2005;37(3):339-353.
- [34] Cai, S, Xi, J. A control approach for pore size distribution in the bone scaffold based on the hexahedral mesh refinement. Computer-Aided Design 2008;40(10-11):1040-1050.
- Kou, ST, Tan, ST. An approach of irregular porous structure modeling [35] based on subdivision and nurbs. Computer-Aided Design and Applications 2013;10(2):355-369.
- [36] You, Y, Kou, S, Tan, S. A new approach for irregular porous structure modeling based on centroidal voronoi tessellation and B-spline. Computer-Aided Design and Applications 2016;13(4):484-489
- [37] Feng, J, Fu, J, Shang, C, Lin, Z, Li, B. Porous scaffold design by solid T-splines and triply periodic minimal surfaces. Computer Methods in Applied Mechanics & Engineering 2018;336:333-352.
- [38] Yang, N, Zhou, K. Effective method for multi-scale gradient porous scaffold design and fabrication. Materials Science & Engineering C 2014:43:502-505.
- [39] Yoo, D. Porous scaffold design using the distance field and triply periodic minimal surface models. Biomaterials 2011;32(31):7741-7754.
- Lagarias, JC, Wang, Y. Haar type orthonormal wavelet bases in R2. [40] Journal of Fourier Analysis and Applications 1995;2(1):1-14.
- Kimura, M, Obayashi, I, Takeichi, Y, Murao, R, Hiraoka, Y. Non-[41] empirical identification of trigger sites in heterogeneous processes using persistent homology. Scientific Reports 2018;8(1):3553-3553.
- Dau, HA, Keogh, E, Kamgar, K, Yeh, CCM, Zhu, Y, Gharghabi, S, [42] et al. The ucr time series classification archive. 2018. https://www. cs.ucr.edu/~eamonn/time_series_data_2018/.

64

65

66

67

68

69

70

71

112

113

114

115

116

117

118

119

120

121

122

123

124

129

130

131

132

133

134

- [43] Tralie, C, Saul, N, Bar-On, R. Ripser.py: A lean persistent homology library for python. The Journal of Open Source Software 2018;3(29):925. URL: https://doi.org/10.21105/ joss.00925.doi:10.21105/joss.00925.
- [44] Morozov, D. Dionysus 2 documentation. https://www.mrzv.org/ software/dionysus2/;2019.
- [45] Orsenigo, C, Vercellis, C. Kernel ridge regression for out-of-sample mapping in supervised manifold learning. Expert Systems With Applications 2012;39(9):7757-7762.
- [46] Borg, I, Groenen, PJF. Modern multidimensional scaling: Theory and applications. Journal of Educational Measurement 2003;40(3):277-280.
- [47] Tenenbaum, JB, De Silva, V, Langford, J. A global geometric framework for nonlinear dimensionality reduction. 13 Science 2000;290(5500):2319-2323.
- [48] Seversky, LM, Davis, S, Berger, M. On time-series topological data 15 analysis: New data and opportunities. In: Computer Vision & Pattern 16 Recognition Workshops. 2016,. 17
- [49] Kerber, M, Morozov, D, Nigmetov, A. Geometry helps to com-18 pare persistence diagrams. Journal of Experimental Algorithmics (JEA) 19 $2017 \cdot 22 \cdot 1 - 4$ 20
- [50] Gandy, PJF, Bardhan, S, Mackay, AL, Klinowski, J. Nodal surface 21 approximations to the ja:math and I-WP triply periodic minimal surfaces. 22 Chemical Physics Letters 2001;336(3):187-195. 23
- Charles, RQ, Su, H, Kaichun, M, Guibas, LJ. Pointnet: Deep learn-24 [51] ing on point sets for 3D classification and segmentation. In: 2017 IEEE 25 Conference on Computer Vision and Pattern Recognition (CVPR). 2017, 26 p. 77-85 27
- [52] Rudin, W, et al. Principles of mathematical analysis; vol. 3. McGraw-hill 28 New York: 1976. 29

Appendix A: Proof of Theorem 4.1

At first, two inequations in $L^p(\mathbb{R}^2)$ are proved in Lemma 1. 31

Lemma 1. 1. Given $f, g \in L^p(\mathbb{R}^2)$, we have

$$\|f + g\|_p^p \le 2^{p-1} \left(\|f\|_p^p + \|g\|_p^p \right), \text{ where } p \ge 1.$$
(15)

2. Given $f_1, f_2, \dots, f_n \in L^p(\mathbb{R}^2)$, where $n < \infty$, it follows that

$$\left|\sum_{i=1}^{n} f_{i}\right|_{p}^{p} \le 2^{(p-1)(n-1)} \sum_{i=1}^{n} ||f_{i}||_{p}^{p}, \text{ where } p \ge 1.$$
 (16)

32

2

3

4

5

8

9

10

11

12

14

Proof. To show Equation (15), that is,

$$\int_{\mathbb{R}^2} |f+g|^p \ dxdy \le 2^{p-1} \left(\int_{\mathbb{R}^2} |f|^p \ dxdy + \int_{\mathbb{R}^2} |g|^p \ dxdy \right)$$
(17)

holds, we show that for $p \ge 1$,

$$|f+g|^{p} \le 2^{p-1}(|f|^{p}+|g|^{p}).$$
(18)

It is obvious when p = 1. Then, it is easy to show that $y(x) = x^p$ is convex over \mathbb{R}^+ , i.e., for $0 < x_1 \le x_2$, we have $y((x_1 + x_2)/2) \le y(x_1 + x_2)/2$ $(y(x_1) + y(x_2))/2$. Therefore, for p > 1, we have

$$\left|\frac{1}{2}f + \frac{1}{2}g\right|^{p} \le \left|\frac{1}{2}|f| + \frac{1}{2}|g|\right|^{p} \le \frac{1}{2}\left(|f|^{p} + |g|^{p}\right).$$
(19)

It shows that Equation (18) holds. Since $f, g \in L^p(\mathbb{R}^2)$, Equa-33 tion (15) holds by integrating on both sides, which suggests that 34 $f + g \in L^p(\mathbb{R}^2).$ 35

Moreover, Equation (16) follows by repeatedly using the E-36 quation (15). 37

Proof of Theorem 4.1: 38



Fig. 5. Illustration of two cases of the relationship between S_i and S'_i .

Proof. For $r(s,t), r'(s,t) \in L^p(\Omega)$, the *p*-norm of the difference of r(s, t) and r'(s, t) is to be evaluated. It follows that

$$||r - r'||_p^p = \int_{\Omega} \left| \sum_{i \in I} r_i - \sum_{j \in I'} r'_j \right|^p \, ds dt.$$
(20)

 r_i 's and r'_i 's are paired according to the matching φ . We have

$$\begin{aligned} \|r - r'\|_{p}^{p} &= \int_{\Omega} \left| \sum_{i \in M} r_{i} - r'_{\varphi(i)} + \sum_{j \in I \setminus M} r_{j} + \sum_{k \in I' \setminus M'} r'_{k} \right|^{p} ds dt \\ &\leq \int_{\Omega} \left(\sum_{i \in M} |r_{i} - r'_{\varphi(i)}| + \sum_{j \in I \setminus M} |r_{j}| + \sum_{k \in I' \setminus M'} |r'_{k}| \right)^{p} ds dt. \end{aligned}$$

$$(21)$$

According to the assumption (2) that the number of intervals in Bc and Bc' is less than L, the total number of pairs and single elements determined by φ is less than 2L, that is, $|M| + |I \setminus M| +$ $|I' \setminus M'| \le 2L$. Hence, it follows by Equation (16) in Lemma 1 that

$$\begin{split} \|r - r'\|_{p}^{p} &\leq 2^{(p-1)(2L-1)} \left(\sum_{i \in M} \int_{\Omega} |r_{i} - r'_{\varphi(i)}|^{p} \, ds dt \right. \\ &+ \sum_{j \in I \setminus M} \int_{\Omega} |r_{j}|^{p} \, ds dt + \sum_{k \in I' \setminus M'} \int_{\Omega} |r'_{k}|^{p} \, ds dt \right) \\ &= 2^{(p-1)(2L-1)} \left(\sum_{i \in M} \|r_{i} - r'_{\varphi(i)}\|_{p}^{p} + \sum_{j \in I \setminus M} \|r_{j}\|_{p}^{p} + \sum_{k \in I' \setminus M'} \|r'_{k}\|_{p}^{p} \right). \end{split}$$
(22)

Let $\varepsilon_i = \max\{|b_i - b'_{\varphi(i)}|, |d_i - d'_{\varphi(i)}|\}$ if the interval is paired by the matching φ to achieve the infimum, $\varepsilon_j = ||l_j||_{\infty}$, and $\varepsilon_k =$ $||l'_k||_{\infty}$. Hence, the 1-Wasserstein distance is able to be expressed by

$$W_1(Bc, Bc') = \left\{ \sum_{i \in M} \varepsilon_i + \sum_{j \in I \setminus M} \varepsilon_j + \sum_{k \in I' \setminus M'} \varepsilon_k \right\}.$$
 (23)

The persistence of (b_i, d_i) is $d_i - b_i \ge 0$. We denote it as *per_i*. Since $W_1(Bc, Bc') \le 1$, it follows by Equation (11) in the paper that $\varepsilon_* \leq 1$, where * represents subscripts *i*, *j* and *k* given above.

For the intervals paired according to the matching φ , the 43 nonzero valued region S_i of $r_i(s, t)$, induced by $(b_i, d_i) \in Bc$, has 44 different relevant relation with the region S'_{i} of r'(s, t), induced by $(\dot{b}_{\varphi(i)}, \dot{d}_{\varphi(i)})$. In general, three cases are taken into considera-46 tion. 47

Case 1: S_i and S'_i overlap but one does not contain the other. That is, we have

$$\|r_i - r'_{\varphi(i)}\|_p^p = \int_{T_1 \cup T_2} |1|^p \, ds dt = \operatorname{Area}(T_1 \cup T_2), \tag{24}$$

39

40

41

42

45

where T_1 and T_2 are two trapezoids formed by $(S_i \cup S'_i) \setminus (S_i \cap S'_i) \setminus (S_i \cap S'_i)$ S'). Without loss of generality, Area $(T_1 \cup T_2)$ is computed according to the case shown in Figure 5(a). It is obtained by directly computing the areas that

Area
$$(T_1) = \frac{1}{2} \left[per'_{\varphi(i)} + per_i - (d'_{\varphi(i)} - d_i) \right] (b'_{\varphi(i)} - b_i),$$
 (25)

and

Area
$$(T_2) = \frac{1}{2} \left[per'_{\varphi(i)} + per_i - (b'_{\varphi(i)} - b_i) \right] (d'_{\varphi(i)} - d_i).$$
 (26)

And then we have

$$Area(T_{1} \cup T_{2}) = Area(T_{1}) + Area(T_{2})$$

$$\leq \frac{1}{2}\varepsilon_{i} \left[2per'_{\varphi(i)} + 2per_{i} - (d'_{\varphi(i)} - d_{i}) - (b'_{\varphi(i)} - b_{i}) \right]$$

$$\leq \frac{1}{2}\varepsilon_{i} \left| 2per'_{\varphi(i)} + 2per_{i} + d_{i} - d'_{\varphi(i)} + b_{i} - b'_{\varphi(i)} \right|$$

$$\leq \frac{1}{2}\varepsilon_{i} \left(2per'_{\varphi(i)} + 2per_{i} + |d_{i} - d'_{\varphi(i)}| + |b_{i} - b'_{\varphi(i)}| \right)$$

$$\leq (2per + \varepsilon_{i})\varepsilon_{i},$$

$$(27)$$

where *per* represents the largest persistence in Bc and Bc', i.e., $per = \max\{\max_{i \in I} (d_i - b_i), \max_{i \in I'} (d'_i - b'_i)\}$. Since $\varepsilon_i \le 1$ mentioned above,

$$\|r_i - r'_{\varphi(i)}\|_p^p \le 2per \cdot \varepsilon_i + \varepsilon_i^2 \le (2per + 1)\varepsilon_i.$$
(28)

Case 2: For S_i and S'_i , one contains the other.

Without loss of generality, the case is considered as shown in Figure 5(b).

$$\begin{aligned} \|r_{i} - r_{\varphi(i)}^{'}\|_{p}^{p} &= \int_{S_{i} \setminus S_{i}^{'}} |1|^{p} \, dsdt \\ &= \operatorname{Area}(S_{i} \setminus S_{i}^{'}) = \frac{1}{2}(per_{i}^{2} - per_{\varphi(i)}^{'2}) \\ &= \frac{1}{2}(per_{i} + per_{\varphi(i)}^{'})(per_{i} - per_{\varphi(i)}) \\ &\leq \frac{1}{2}(per_{i} + per_{\varphi(i)}^{'}) \left| d_{i} - d_{\varphi(i)}^{'} + b_{\varphi(i)}^{'} - b_{i} \right| \\ &\leq \frac{1}{2}(per_{i} + per_{\varphi(i)}^{'}) \left(|d_{i} - d_{\varphi(i)}^{'}| + |b_{\varphi(i)}^{'} - b_{i}| \right) \\ &\leq 2per \cdot \varepsilon_{i} \leq (2per + 1)\varepsilon_{i}. \end{aligned}$$
(29)

- where per has the same meaning as Equation (27).
- **Case 3**: S_i and S'_i do not overlap. 3 Notice that in this case we have

$$\max\{|b_{i} - b'_{\varphi(i)}|, |d_{i} - d'_{\varphi(i)}|\} \ge \max\{per_{i}, per'_{\varphi(i)}\}$$

$$\ge \frac{d_{i} - b_{i}}{2} + \frac{d'_{\varphi(i)} - b'_{\varphi(i)}}{2},$$
(30)

i.e.,

6

$$\|l_i - l'_{\varphi(i)}\|_{\infty} \ge \|l_i\|_{\infty} + \|l'_{\varphi(i)}\|_{\infty}, \tag{31}$$

which means that it is better if l_i and $l'_{\omega(i)}$ are not paired. There-4

fore, case 3 does not happen if the matching φ is the one to 5 reach the infimum.

For the intervals not paired via the matching φ , it follows that

$$\int_{\Omega} |r_*|^p \, dsdt = \int_{S_*} |1|^p \, dsdt = \frac{1}{2} per_*^2, \tag{32}$$

As defined in Equation (9), it follows by $per_* = 2\varepsilon_*$ that

$$\int_{\Omega} |r_*|^p \, dsdt = per_* \cdot \varepsilon_* \le (2per+1)\varepsilon_*,\tag{33}$$

where * represents the subscripts $j \in I \setminus M$ and $k \in I' \setminus M'$ in the corresponding situation, and *per* is the largest persistence in the barcodes.

Eventually, it follows by Equation (22) that

$$\begin{split} ||r - r'||_{p}^{p} &\leq 2^{(p-1)(2L-1)} \left(\sum_{i \in M} ||r_{i} - r_{\varphi(i)}^{'}||_{p}^{p} + \sum_{j \in I \setminus M} ||r_{j}||_{p}^{p} + \sum_{k \in I' \setminus M'} ||r_{k}^{'}||_{p}^{p} \right) \\ &\leq 2^{(p-1)(2L-1)} (2per + 1) \left(\sum_{i \in M} \varepsilon_{i} + \sum_{j \in I \setminus M} \varepsilon_{j} + \sum_{k \in I' \setminus M'} \varepsilon_{k} \right) \\ &= 2^{(p-1)(2L-1)} (2per + 1) W_{1}(Bc, Bc'). \end{split}$$
(34)

Because the domain is restricted on $[0, B] \times [0, B]$, in a barcode, it is equivalent to consider the intervals satisfying $0 \le b_i \le d_i \le$ B, to truncate the intervals in which $b_i < B$ and $d_i > B$ into (b_i, B) , and to neglect the intervals in which $b_i > B$. Whence, it is equivalent to assume that $per = \max_{i \in I} (d_i - b_i) \le B$ for any 14 barcode on the domain, which makes Equation (11) hold. 15

10

11

12

13

16

17

19

20

23

Appendix B: Proof of Corollary 1

Proof. Let r_{upper} represent the function restricting \tilde{r} on the domain above the diagonal and r_{down} represent the function below the diagonal. Then the equation $\tilde{r} = r_{upper} + r_{down}$ holds. Because of the definition of extended PRF given in Equation (4) and Equation (5), we have

$$||r - r'||_p = ||r_{upper} - r'_{upper}||_p = ||r_{down} - r'_{down}||_p.$$
(35)

Therefore, it follows by Equation (15) that

$$\|\widetilde{r} - \widetilde{r}'\|_{p}^{p} = \|(r_{upper} - r_{upper}') + (r_{down} - r_{down}')\|_{p}^{p} \le 2^{p} \|r - r'\|_{p}^{p}.$$
 (36)

By using the result in Theorem 4.1, Equation (12) holds.

Appendix C: Proof of Theorem 4.2

Before proving the final stability theorem of the feature vectors, a lemma is provided as follows.

Lemma 2. For 2D Haar basis $\{har_i(s, t)\}_{i=1}^{\infty}$ and $f \in L^2(\Omega)$, the Parseval identity [52] holds, that is,

$$\|f\|_{2}^{2} = \sum_{i=1}^{\infty} \langle f, har_{i} \rangle^{2} = \sum_{i=1}^{\infty} \lambda_{i}^{2}.$$
 (37)

Proof. It follows the fact that Haar basis is a complete and or-21 thonormal basis in $L^2(\Omega)$ [40]. 22

Eventually, we give the proof of Theorem 4.2.

Proof. It is left to show that $||v - v'||_2^2 \le ||\tilde{r} - \tilde{r}'||_2^2$. According to the definition of feature vector given in Equation (8) and Equation 37 in Lemma 2, we have

$$\|v - v'\|_{2}^{2} = \sum_{i=1}^{N^{2}} (\lambda_{i} - \lambda_{i}')^{2} \le \|\widetilde{r} - \widetilde{r}'\|_{2}^{2}.$$
 (38)

Then, it follows by Equation (12) that when p = 2, Equation 24 (13) holds. 25