

Higher Methods in Data Science and ML

On the encoding of large, high-dimensional and unorganized datasets

Ulderico Fugacci

30 January 2020





Topological Data Analysis (TDA) and Persistent Homology (PH) allow for extracting the core topological information from *large*, *high-dimensional* and *unorganized datasets*. E.g. *point clouds*, *complex networks*, *(semi-)metric spaces*



Dataset





Topological Information



Topological Data Analysis (TDA) and Persistent Homology (PH) allow for extracting the core topological information from *large*, *high-dimensional* and *unorganized datasets*. E.g. *point clouds*, *complex networks*, *(semi-)metric spaces*





Simplicial Complex:

A family K of subsets (called *simplices*) of a finite set V closed under the operation of taking subsets









Simplicial Complex:

A family K of subsets (called *simplices*) of a finite set V closed under the operation of taking subsets



Development of

compact and efficient data structures

for encoding simplicial complexes

Solution:

0-, 1-, 2-, 3simplices





Outline:

- Which info to be stored?
- Data structures:
 - Simplex-based representations
 - Top-based representations
 - Operator-driven representations
- Issues and solutions in adopting top-based representations
- Comparisons



Out of Scope:

- Data structures for specific classes of complexes
 - E.g., manifolds or complexes of low dimension
- Hierarchical and multi-resolution models
- Construction of a simplicial complex from a dataset

POLITECNICO Di TORINO

Data structure:

The *entities* which a simplicial complex consists of are:

its simplices

 $K = K_0 \cup K_1 \cup \ldots \cup K_d$

where K_i is the collection of the i-simplices of K

the topological relations

$$R_{i,j} \subseteq K_i \times K_j$$

between the simplices of K encoding the (co-)boundary of each simplex



Data structure:

The *entities* which a simplicial complex consists of are:

its simplices

 $K = K_0 \cup K_1 \cup \ldots \cup K_d$

where K_i is the collection of the i-simplices of K

• the *topological relations*

 $R_{i,j} \subseteq K_i \times K_j$

between the simplices of K encoding the (co-)boundary of each simplex

A *data structure* for K has to explicitly *store* a portion of the above information and to (efficiently) *retrieve* the remaining part









Data structures:	
Store all the entities	
	Compactness
	Compaciness
 Simplex-based representations Top-based representations 	
• Operator-driven representations	top-simplices











Simplex-based Representations





Simplex-based Representations





Graph is *not uniquely determined* but it depends on the chosen vertex order

Simplex-based Representations





Graph is *not uniquely determined* but it depends on the chosen vertex order

Top-based Representations





Compact: it explicitly stores just a fraction of the entities of a simplicial complex

Not all the relations between simplices are immediately available

Operator-driven Representations



POLITECNICO DI TORINO

The simplicial complex K is encoded by storing its **1-skeleton** (i.e. the graph consisting of the 0- and the 1-simplices) and a **map** returning, for each 1-simplex σ , the blockers of K containing σ , where:

a simplex τ is a *blocker* if τ does not belong to K but all its faces do



Designed for flag complexes (e.g. VR complexes) and edge contraction

Too specific: *inefficient in any other task*



Top-based representations look like promising data structures for encoding a simplicial complex K

But, how to...

1. Store information associated to each simplex of K (e.g. labels, gradient, ...)?

2. Efficiently perform operators having explicitly stored a fraction of the entities of K?



Top-based representations look like promising data structures for encoding a simplicial complex K

But, how to...

1. Store information associated to each simplex of K (e.g. labels, gradient, ...)?

Attach information to the top simplices only



2. Efficiently perform operators having explicitly stored a fraction of the entities of K?



Top-based representations look like promising data structures for encoding a simplicial complex K

But, how to...

1. Store information associated to each simplex of K (e.g. labels, gradient, ...)?

Attach information to the top simplices only



2. Efficiently perform operators having explicitly stored a fraction of the entities of K?

Re-define the algorithms performing the operators trying to extract the lowest possible amount of non-explicitly stored entities

Comparisons



			Top-k	Top-based vs Simplex-based					
	Dataset	d	$ \Sigma_0 $	$ \Sigma_{i} $	$ \Sigma $	Storage Cos		ost	
	Davasev		20	<i>2</i> top		IA^*	IG	ST	
	DTI-SCAN	3	0.9M	$5.5\mathrm{M}$	24M	0.97	11.9	2.4	
	VISMALE	3	4.6M	26M	118M	4.7	-	9.7	
	Ackley4	4	$1.5\mathrm{M}$	$32\mathrm{M}$	$204\mathrm{M}$	6.8	-	12.8	
	Amazon01	6	0.2M	0.4M	2.2M	0.12	1.6	0.3	
	Amazon02	7	0.4M	1.0M	$18.4\mathrm{M}$	0.28	9.8	1.5	
	Roadnet	3	$1.9\mathrm{M}$	$2.5\mathrm{M}$	$4.8\mathrm{M}$	0.8	3.3	1.0	
	Sphere-1.0	16	100	224	0.6M	0.003	0.9	0.04	
	Sphere-1.2	21	100	285	$26\mathrm{M}$	0.0032	-	1.5	
	Sphere-1.3	23	100	382	$197 \mathrm{M}$	0.0034	-	11.01	

Comparisons





Comparisons



Top-based vs Operator-driven

data	ω		contr. timings			memory peak		
uutu			edges	check	contr.	tot	gen.	simpl.
	28	weak	6.38K	9.15h	2.27 <i>m</i>	9.19h	5.6	57.2K
0		top		0.01s	0.02s	0.09s		7.6
CAC		Skel.		0.00s	0.15s	0.15s	7.8	7.8
CHIC		weak		ou	t-of-memo	6.2	_	
Ũ	56	top	7.99K	0.04 <i>s</i>	0.06s	0.23s	14.1	10.8
		Skel.		0.00s	0.71 <i>s</i>	0.71 <i>s</i>		14.1
		weak		out-of-memory			11.6	_
s	63	top	27.9K	0.08s	0.11s	0.38s	11.0	14.9
HEN		Skel.		0.00s	0.74 <i>s</i>	0.75 <i>s</i>	26.4	26.8
ATF		weak		out-of-memory			10.0	_
	126	top	31.2K	0.40s	0.49s	1.36s	10.0	25.9
		Skel.		0.01s	7.73s	7.74 <i>s</i>	66.1	66.7
LE	3.5	weak		34.3 <i>m</i>	1.28 <i>m</i>	40.4 <i>m</i>	1.01/	2.0K
MA		top	4.23M	4.34m	0.89m	7.20m	1.0K	2.0K
VISI		Skel.		0.76 <i>m</i>	3.34h	3.35h	8.0K	8.0K
E		weak		killed	after 25	7 5V	_	
00	4.5	top	4.69M	2.89h	26.0m	3.32h	/. 5 K	10.7K
ц		Skel.		killed after 25 hours			19.4K	-
Υ		weak		killed after 25 hours			7 5V	_
'nc	1.5	top	14.0M	11.9m	14.8 <i>m</i>	32.0m	7.5K	15.4K
Г		Skel.		23.19 <i>s</i>	14.6h	14.6h	50.9K	52.1K

In Summary



We have briefly overviewed the *most common data structures* proposed in the literature

Future Directions:

- Express the most frequently adopted operators in terms of top simplices
- Face the bottleneck concerning the *construction of a simplicial complex from a dataset* (maybe proposing an approximated construction?)
- Investigate with your help the connections between simplicial complexes (and the data structures to store them) and *itemsets in association rule learning*



Thank you!

Main References:



Simplex-based Representations:

- Boissonnat, Maria. The Simplex Tree: an Efficient Data Structure for General Simplicial Complexes. In Algorithmica, 2014

Top-based Representations:

- Canino, De Floriani, Weiss. IA*: An Adjacency-Based Representation for Non-manifold Simplicial Shapes in Arbitrary Dimensions. In Computers & Graphics, 2011
- Fellegara, Weiss, De Floriani. The Stellar Tree: a Compact Representation for Simplicial Complexes and Beyond. arXiv preprint, 2017

• Operator-driven Representations:

- Attali, Lieutier, Salinas. Efficient Data Structure for Representing and Simplifying Simplicial Complexes in High Dimensions. In International Journal of Computational Geometry & Applications, 2012

• Comparisons:

- Fugacci, Iuricich, De Floriani. Computing Discrete Morse Complexes from Simplicial Complexes.
 In Graphical Models, 2019
- Fellegara, Iuricich, De Floriani, Fugacci. Efficient Homology-Preserving Simplification of High-Dimensional Simplicial Shapes. In Computer Graphics Forum, 2019