



SENIOR THESIS IN MATHEMATICS

Victor's Senior Thesis

Author:

Victor de Fontnouvelle

Advisor:

Dr. Vin De Silva

Submitted to Pomona College in Partial Fulfillment
of the Degree of Bachelor of Arts

May 1, 2019

Abstract

We will explore various methods of analyzing high-dimensional data. We'll investigate techniques that make inferences about the underlying structure of the data, cluster the data, or reduce the dimension of the data. We'll apply these methods to real-world datasets, compare the methods, and suggest improvements to the methods.

Contents

1	Introduction	1
1.1	Inference About the Structure of the Data	1
1.2	Clustering	1
1.3	Dimensionality Reduction	1
2	Review of Literature	3
3	Helmholtz Decomposition	5
3.1	Explanation of Helmholtz Decomposition	5
3.1.1	Intuition	5
3.1.2	Computation of G , C , and H	6
3.2	Application of Helmholtz Decomposition	7
4	Cohomology Analysis	8
4.1	Explanation of Cohomology Analysis	8
4.2	Application of Cohomology Analysis	8
5	Mapper	9
5.1	Explanation of Mapper	9
5.2	Explanation of Mapper	9
6	Laplacian Eigenvector Analysis	10
6.1	Explanation of Laplacian Eigenvector Analysis	10
6.2	Application of Laplacian Eigenvector Analysis	10
7	Discussion	11
8	Conclusion	12

Chapter 1

Introduction

1.1 Inference About the Structure of the Data

Cohomology analysis and Helmholtz decomposition provide insight on the structure of the data. Cohomology analysis provides a broader framework for detecting clusters, holes, and higher-order features within the data. Helmholtz decomposition is a specific application of cohomology analysis which seeks to find an ordered list which best accounts for a weighted graph.

1.2 Clustering

Mapper clusters the data, by mapping clusters of points onto intervals on the real line using a filter function, and connecting overlapping clusters.

1.3 Dimensionality Reduction

Calculating the eigenvalues of the laplacian matrix reduces the dimension of the data. The laplacian is a symmetric matrix encoding the pairwise distances between points, normalized by row. This matrix can be thought of as a linear operator encoding heat flows. Given an input of initial temperatures, it outputs the changes in temperatures after one time step. Eigenvectors corresponding to low eigenvalues thus correspond to stable temperature configurations. The eigenvectors are perpendicular, and thus eigenvectors corresponding to slightly higher eigenvalues often capture geometric structure

existing in the data. Additionally, nearby points will have similar values in the eigenvectors, and thus the eigenvectors are also a useful tool for reducing the dimension of the data.

Chapter 2

Review of Literature

Papers by Gunnar Carlsson [1] and Vin de Silva [3] explain the intuition behind cohomology analysis, and provides several examples. Carlsson [2] and deSilva [3] also describe the algorithm used to compute cohomology, which will be useful if I choose to implement it.

Curto [4] and Ulmer [5] both provide various examples of the uses of homology analysis, both in detecting underlying structure, and in providing fingerprints that identify different phenomena. Curto analyzes a dataset representing connection strengths between neurons in rats responsible for spatial recognition. Curto first keeps a certain proportion of edges ρ s.t. $0 < \rho < 1$, keeping those edges which are the strongest. Curto then runs cohomology analysis on this graph to determine the Betti curve. Curto determined that the Betti curve obtained from spatially organized neurons is different than the Betti curve that would be obtained from neurons with random structure. Cohomology analysis thus provides a method for detecting spatial neurons. Ulmer uses cohomology analysis to evaluate two different models of social interaction for aphids roaming in a dish. Standard measures that compare the models to real data include angular momentum, and average distance to closest neighbor. Ulmer found that cohomology analysis provided an equally strong measurement for assessing model accuracy.

A paper by Jiang [6] explains the Helmholtz decomposition that converts ranked data into ordinal data. It provides both the algorithm for the decomposition, as well as three examples of its use. Jiang uses Helmholtz decomposition to rank movies. Most users rate several movies, so each time a user rated two movies, this introduced an edge from one movie to the other indicating the user's preference. Jiang used Helmholtz decomposition

to create an absolute index of currency values based off trading rates. Jiang also used Helmholtz decomposition to rank websites based off of connection strengths. Jiang found that Helmholtz decomposition performed as well as some of the standard methods for website ranking, indicating that it is a useful tool.

Carlsson [1] also explains the Mapper tool, and provides examples of its use. I have already used the algorithm described in this paper to analyze several datasets.

Singh [7] describes the intuition behind the laplacian analysis, which I have also implemented.

Chapter 3

Helmholtz Decomposition

3.1 Explanation of Helmholtz Decomposition

3.1.1 Intuition

The Hodge Decomposition Theorem provides that \bar{Y} can be expressed as the sum of three orthogonal components:

1. The gradient flow G
2. The curl flow C
3. The harmonic flow H

The gradient flow G represents a comparison matrix that corresponds directly to a vector v where each item is assigned a real value. G is the gradient of v , thus each entry is computed as $G[ij] = v[j] - v[i]$.

The harmonic flow H represents a ranking that is curl-free and divergence-free. This means that any three items in H will have pairwise rankings that are logically consistent. Specifically, $H[ij] + H[jk] + H[ki] = 0, \forall i, j, k$. Because $\text{div}(H) = 0$, H contains plausible values in the sense that it could have been produced by real-world data. A harmonic flow thus indicates that there are cycles in the graph with more than three edges, where the edge weights don't sum to zero.

The curl flow C is the image of curl^* , the adjoint of the curl. Nonzero values of C thus indicate cycles of length three whose edge weights don't sum to zero.

The gradient flow provides the ranking that minimizes the least-squared residual, while the harmonic and curl flows characterize the residual. Specifically, $\bar{Y} - H = G + C$. A large curl flow indicates that the ordering of items ranked closely together is unreliable, while a large harmonic flow indicates that the ordering of items ranked further apart is unreliable. For example, if the curl flow is large but the harmonic flow is small, this indicates that the ranking is valid at a larger scale, but that the specific ordering of closely-ranked alternatives isn't very precise.

3.1.2 Computation of Gradient Flow

We seek to find $S \in \mathbb{R}^n$ that minimizes $\|\delta_0 s - \bar{Y}\|$. Such an s must satisfy $\bar{Y} - \delta_0 s \in \delta_0^\perp = \ker(\delta_0^*)$. Thus $\delta_0^*(\bar{Y} - \delta_0 s) = 0 \implies \delta_0^* \delta_0 s = \delta_0^* \bar{Y}$. Note that $\Delta_0 = \delta_0^* \delta_0$ and $\delta_0^* s = -\text{div}$, thus $s = -\delta_0^+ \text{div} Y$. We compute the gradient flow as $G = \delta_0 s$.

3.1.3 Computation of Curl Flow

We seek a three-dimensional matrix M that minimizes $\|\delta_1 * M - \bar{Y}\|$. By similar logic, $\text{curl} \text{curl}^* M = \text{curl} Y \implies M = (\text{curl} \text{curl}^*)^{-1} \text{curl} Y$. Thus the curl flow $C = \text{curl}^* (\text{curl} \text{curl}^*)^{-1} \text{curl} Y = \text{curl}^+ \text{curl} Y$. We represent the curl matrix as a matrix with width $\binom{n}{2}$ and height $\binom{n}{3}$, which transforms a vector holding all entries $\bar{Y}[ij]$ s.t. $j > i$ and outputs a vector corresponding to all triples (i, j, k) s.t. $i < j < k$. A triple (i, j, k) will be assigned the value $\bar{Y}[ij] + \bar{Y}[jk] + \bar{Y}[ki]$. For example, in the case of four variables, we have

$$\begin{bmatrix} 1 & -1 & 0 & 1 & 0 & 0 \\ 1 & 0 & -1 & 0 & 1 & 0 \\ 0 & 1 & -1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & -1 & 1 \end{bmatrix} \begin{bmatrix} ab \\ ac \\ ad \\ bc \\ bd \\ cd \end{bmatrix} = \begin{bmatrix} abc \\ abd \\ acd \\ bcd \end{bmatrix}$$

To check correctness, we examine the triple (a, b, c) . This triple is assigned the value $\bar{Y}[ab] - \bar{Y}[ac] + \bar{Y}[bc]$. This equals $\bar{Y}[ac] + \bar{Y}[bc] + \bar{Y}[ca]$, as $-\bar{Y}[ac] = \bar{Y}[ca]$ since \bar{Y} is symmetric. This is the curl, as we expect.

3.1.4 Computation of Harmonic Flow

Harmonic flow: The Hodge decomposition theorem provides that $\bar{Y} = G + H + C$, thus we compute H as $\bar{Y} - G - H$.

3.2 Application of Helmholtz Decomposition

Chapter 4

Cohomology Analysis

4.1 Explanation of Cohomology Analysis

4.2 Application of Cohomology Analysis

Chapter 5

Mapper

5.1 Explanation of Mapper

5.2 Explanation of Mapper

Chapter 6

Laplacian Eigenvector Analysis

- 6.1 Explanation of Laplacian Eigenvector Analysis
- 6.2 Application of Laplacian Eigenvector Analysis

Chapter 7

Discussion

Chapter 8

Conclusion

Bibliography

- [1] Carlsson, Gunnar. "Topology and data." *Bulletin of the American Mathematical Society* 46.2 (2009): 255-308.
- [2] Zomorodian, Afra, and Gunnar Carlsson. "Computing persistent homology." *Discrete & Computational Geometry* 33.2 (2005): 249-274.
- [3] De Silva, Vin, Dmitriy Morozov, and Mikael Vejdemo-Johansson. "Persistent cohomology and circular coordinates." *Discrete & Computational Geometry* 45.4 (2011): 737-759.
- [4] Giusti, Chad, et al. "Clique topology reveals intrinsic geometric structure in neural correlations." *Proceedings of the National Academy of Sciences* 112.44 (2015): 13455-13460.
- [5] Ulmer, M., Lori Ziegelmeier, and Chad M. Topaz. "Assessing biological models using topological data analysis." *arXiv preprint arXiv:1811.04827* (2018).
- [6] Jiang, Xiaoye, et al. "Statistical ranking and combinatorial Hodge theory." *Mathematical Programming* 127.1 (2011): 203-244.
- [7] Singh, Gurjeet, Facundo Mmoli, and Gunnar E. Carlsson. "Topological methods for the analysis of high dimensional data sets and 3d object recognition." *SPBG*. 2007.