**ETH**

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

# Persistence of generalized density functions

Master Thesis

Maxim Mikhaylov

January 19, 2038

Supervisors: Prof. Dr. B. Gärtner, Dr. P. Schnider

Department of Computer Science, ETH Zürich

**Abstract**

Topological data analysis (TDA) often uses density-like functions defined on the ambient space $\mathbb{R}^n$ to infer the underlying topological structure of a dataset $P \subseteq \mathbb{R}^n$. The persistent homology of the sublevelset or superlevelset filtrations induced by these functions captures multi-scale topological features. A key desirable property is *stability*: small perturbations in the dataset result in similarly small changes in the persistence diagrams. A classic example is the nearest-neighbor distance function $f(x) := \min_{p \in P} d(x, p)$, for which the bottleneck distance between persistence diagrams $\mathrm{Dgm}_{f(P)}$ and $\mathrm{Dgm}_{f(Q)}$ is bounded by the Hausdorff distance $d_H(P, Q)$ between datasets [1]:

$$d_b(\mathrm{Dgm}_{f(P)}, \mathrm{Dgm}_{f(Q)}) \leq d_H(P, Q).$$

This thesis investigates *generalized density functions* (GDFs), which are functions $f : \mathcal{P}(X) \times X \to \mathbb{R}$, and the conditions under which they satisfy a stability property of the form

$$d_b(\mathrm{Dgm}_{f(P)}, \mathrm{Dgm}_{f(Q)}) \leq c \cdot d_H(P, Q)$$

for some finite constant $c$.

The primary contribution of this work are stability theorems for several classes of generalized density functions. Specifically, we prove stability bounds for:

- Several generalizations of the nearest-neighbor distance function of the form $f(P, x) = \min_{p \in P} h(x, p)$.

- Functons that are Lipschitz continuous with respect to the Hausdorff distance on the space of point clouds.

- Morse functions that satisfy a Lipschitz-like condition.

Beyond these core stability results, we briefly explore the properties of the space of stable functions and investigate how common operations, such as addition and taking minima, affect stability. We identify conditions under which stability is preserved under these operations and provide counterexamples demonstrating cases where it is not.

Our findings unify and extend existing stability results, offering practical guidance for the selection and design of generalized density functions for topological data analysis.

# Contents

Chapter 1

# Introduction

Modern datasets often contain geometric or structural relationships that are poorly captured by traditional statistical methods. Topological data analysis (TDA) addresses this by computing *persistent homology*, a multiscale descriptor of connectivity, holes, and higher-dimensional voids. A critical step in TDA is constructing a filtration of topological spaces from data, often via real-valued functions that reflect data density.

A natural approach is to define a generalized density function (GDF) $f(P, x) : \mathcal{P}(X) \times X \to \mathbb{R}$, where $P$ is a dataset. We can then consider the *sublevelset filtration* of $f$ to analyze the topological features of $P$ at different scales. For example, the nearest-neighbor distance function $f(P, x) = \min_{p \in P} d(x, p)$ is a GDF that captures the distance from a point $x$ to its nearest point in $P$. The sublevelset filtration of this function is precisely the Čech filtration of $P$, which can be intuitvely described as the process of growing balls around each point in $P$. The *persistence diagram* of the Čech filtration is *stable* with respect to the Hausdorff distance between datasets, that is

$$d_b(\mathrm{Dgm}_{f(P)}, \mathrm{Dgm}_{f(Q)}) \leq d_H(P, Q), \tag{1.1}$$

where $d_b$ is the bottleneck distance, a metric on persistence diagrams. This stability property is crucial for TDA, as it ensures that small perturbations in the data lead to small changes in the persistence diagrams, which provides robustness to noise.

A natural question arises: which GDFs yield persistence diagrams that are stable? Formally, we say that a GDF $f$ is *c-stable* if for all datasets $P, Q \subseteq X$,

$$d_b(\mathrm{Dgm}_{f(P)}, \mathrm{Dgm}_{f(Q)}) \leq c \cdot d_H(P, Q). \tag{1.2}$$

## 1.1 Overview

TODO:

# Background

This chapter introduces the mathematical concepts and tools used within this thesis, as well as known results. We begin with *persistent homology*, a fundamental tool for measuring topological features at different scales. TODO: finish

## 2.1   Persistent homology

*Simplicial complexes* are a combinatorial structure used to study topological spaces. They are built from simplices, which are the generalization of points, line segments, triangles, and higher-dimensional shapes. A simplicial complex is a set of simplices that satisfy certain conditions TODO: say which ones. A *filtered* simplicial complex is a collection of simplices $(K_\varepsilon)_{\varepsilon \in \mathbb{R}}$ that is indexed by a parameter $\varepsilon$, where $K_\varepsilon \subseteq K_{\varepsilon'}$ for $\varepsilon < \varepsilon'$.

A real-valued function $f : X \to \mathbb{R}$ on a topological space $X$ induces a filtered simplicial complex by considering the sublevel sets TODO: shit definition, talk about triangulation

$$K_\varepsilon = \{x \in X \mid f(x) \leq \varepsilon\}.$$

Such a filtration is called a *sublevelset filtration*.

A filtered simplicial complex can be summarized by its *persistent homology*, which serves as a fundamental tool in TDA that tracks how topological features (connected components, voids) appear and disappear as a parameter is varied [2]. Given a filtered simplicial complex $(K_\varepsilon)_{\varepsilon \in \mathbb{R}}$, persistent homology computes pairs $(b_i, d_j)$ where a feature appears at $\varepsilon = b_i$ (birth) and disappears at $\varepsilon = d_i$ (death).

These birth–death pairs can be visualized in a *persistence diagram*, where each pair is shown as a point in the Euclidean plane. The distance of a point from
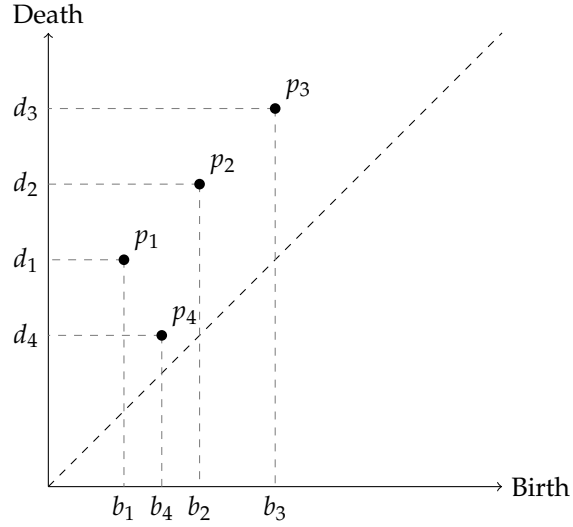
**Figure 2.1:** An example of a persistent diagram. A point $p_i$ is born at $b_i$ and dies at $d_i$.

the diagonal indicates the feature's persistence, which is often interpreted as a measure of its significance. Figure 2.1 shows an example of a persistence diagram.

Two persistence diagrams can be compared using the *bottleneck distance*, which measures the cost of transforming one diagram into another by matching points. Formally, given two persistence diagrams $D_1$ and $D_2$, the bottleneck distance is defined as

$$d_b(D_1, D_2) = \inf_{\pi} \max_{p \in D_1} \|p - \pi(p)\|, \tag{2.1}$$

where $\pi$ is a bijection between the points in $D_1$ and $D_2$, and $\|\cdot\|$ is the Euclidean distance. TODO: wrong definition, talk about allowed diagonal matching An alternative metric is the *Wasserstein distance*, which generalizes the bottleneck distance by computing the $p$-norm of the transport distance TODO: bad way to say this:

$$d_p(D_1, D_2) = \inf_{\pi} \left( \sum_{p \in D_1} \|p - \pi(p)\|^p \right)^{1/p}. \tag{2.2}$$

## 2.2 Stability

Chapter 3

# Writing scientific texts in English

This chapter was originally a separate document written by Reto Spöhel. It is reprinted here so that the template can serve as a quick guide to thesis writing, and to provide some more example material to give you a feeling for good typesetting.

## 3.1 Basic writing rules

The following rules need little further explanation; they are best understood by looking at the example in the booklet by Knuth et al., §2–§3.

**Rule 3.1** Write texts, not chains of formulas.

More specifically, write full sentences that are logically interconnected by phrases like 'Therefore', 'However', 'On the other hand', etc. where appropriate.

**Rule 3.2** Displayed formulas should be embedded in your text and punctuated with it.

In other words, your writing should not be divided into 'text parts' and 'formula parts'; instead the formulas should be tied together by your prose such that there is a natural flow to your writing.

## 3.2 Being nice to the reader

Try to write your text in such a way that a reader enjoys reading it. That's of course a lofty goal, but nevertheless one you should aspire to!

**Rule 3.3** Be nice to the reader.

Give some intuition or easy example for definitions and theorems which might be hard to digest. Remind the reader of notations you introduced

many pages ago – chances are he has forgotten them. Illustrate your writing with diagrams and pictures where this helps the reader. Etc.

**Rule 3.4** Organize your writing.

Think carefully about how you subdivide your thesis into chapters, sections, and possibly subsections. Give overviews at the beginning of your thesis and of each chapter, so the reader knows what to expect. In proofs, outline the main ideas before going into technical details. Give the reader the opportunity to 'catch up with you' by summing up your findings periodically.

*Useful phrases:* 'So far we have shown that . . . ', 'It remains to show that . . . ', 'Recall that we want to prove inequality (7), as this will allow us to deduce that . . . ', 'Thus we can conclude that . . . . Next, we would like to find out whether . . . ', etc.

**Rule 3.5** Don't say the same thing twice without telling the reader that you are saying it twice.

Repetition of key ideas is important and helpful. However, if you present the same idea, definition or observation twice (in the same or different words) without telling the reader, he will be looking for something new where there is nothing new.

*Useful phrases:* 'Recall that [we have seen in Chapter 5 that] . . . ', 'As argued before / in the proof of Lemma 3, . . . ', 'As mentioned in the introduction, . . . ', 'In other words, . . . ', etc.

**Rule 3.6** Don't make statements that you will justify later without telling the reader that you will justify them later.

This rule also applies when the justification is coming right in the next sentence! The reasoning should be clear: if you violate it, the reader will lose valuable time trying to figure out on his own what you were going to explain to him anyway.

*Useful phrases:* 'Next we argue that . . . ', 'As we shall see, . . . ', 'We will see in the next section that . . . , etc.

## 3.3   A few important grammar rules

**Rule 3.7** There is (almost) *never* a comma before 'that'.

It's really that simple. Examples:

> We assume that . . .
> *Wir nehmen an, dass . . .*

> It follows that . . .
> *Daraus folgt, dass . . .*
>
> 'thrice' is a word that is seldom used.
> *'thrice' ist ein Wort, das selten verwendet wird.*

Exceptions to this rule are rare and usually pretty obvious. For example, you may end up with a comma before 'that' because 'i.e.' is spelled out as 'that is':

> For $p(n) = \log n / n$ we have . . . However, if we choose $p$ a little bit higher, that is $p(n) = (1 + \varepsilon) \log n / n$ for some $\varepsilon > 0$, we obtain that. . .

Or you may get a comma before 'that' because there is some additional information inserted in the middle of your sentence:

> Thus we found a number, namely $n_0$, that satisfies equation (13).

If the additional information is left out, the sentence has no comma:

> Thus we found a number that satisfies equation (13).

(For 'that' as a relative pronoun, see also Rules 3.9 and 3.10 below.)

**Rule 3.8** There is usually no comma before 'if'.

Example:

> A graph is not 3-colorable if it contains a 4-clique.
> *Ein Graph ist nicht 3-färbbar, wenn er eine 4-Clique enthält.*

However, if the 'if' clause comes first, it is usually separated from the main clause by a comma:

> If a graph contains a 4-clique, it is not 3-colorable .
> *Wenn ein Graph eine 4-Clique enthält, ist er nicht 3-färbbar.*

There are more exceptions to these rules than to Rule 3.7, which is why we are not discussing them here. Just keep in mind: don't put a comma before 'if' without good reason.

**Rule 3.9** Non-defining relative clauses have commas.

**Rule 3.10** Defining relative clauses have no commas.

In English, it is very important to distinguish between two types of relative clauses: defining and non-defining ones. This is a distinction you absolutely need to understand to write scientific texts, because mistakes in this area actually distort the meaning of your text!

It's probably easier to explain first what a *non-defining* relative clause is. A non-defining relative clauses simply gives additional information *that could also be left out* (or given in a separate sentence). For example, the sentence

> The WeirdSort algorithm, which was found by the famous mathematician John Doe, is theoretically best possible but difficult to implement in practice.

would be fully understandable if the relative clause were left out completely. It could also be rephrased as two separate sentences:

> The WeirdSort algorithm is theoretically best possible but difficult to implement in practice. [By the way,] WeirdSort was found by the famous mathematician John Doe.

This is what a non-defining relative clause is. *Non-defining relative clauses are always written with commas.* As a corollary we obtain that you cannot use 'that' in non-defining relative clauses (see Rule 3.7!). It would be wrong to write

> ~~The WeirdSort algorithm, that was found by the famous mathematician John Doe, is theoretically best possible but difficult to implement in practice.~~

A special case that warrants its own example is when 'which' is referring to the entire preceding sentence:

> Thus inequality (7) is true, which implies that the Riemann hypothesis holds.

As before, this is a non-defining relative sentence (it could be left out) and therefore needs a comma.

So let's discuss *defining* relative clauses next. A defining relative clause tells the reader *which specific item the main clause is talking about*. Leaving it out either changes the meaning of the sentence or renders it incomprehensible altogether. Consider the following example:

> The WeirdSort algorithm is difficult to implement in practice. In contrast, the algorithm that we suggest is very simple.

Here the relative clause 'that we suggest' cannot be left out – the remaining sentence would make no sense since the reader would not know which algorithm it is talking about. This is what a defining relative clause is. *Defining relative clauses are never written with commas.* Usually, you can use both 'that' and 'which' in defining relative clauses, although in many cases 'that' sounds better.

As a final example, consider the following sentence:

> For the elements in $\mathcal{B}$ which satisfy property (A), we know that equation (37) holds.

This sentence does not make a statement about all elements in $\mathcal{B}$, only about those satisfying property (A). The relative clause is *defining*. (Thus we could also use 'that' in place of 'which'.)

**Table 3.1:** Things you (usually) don't say

| It holds (that) ... | We have ... | *Es gilt ...* |
|---|---|---|
| ('Equation (5) holds.' is fine, though.) | | |
| x fulfills property $\mathcal{P}$. | $x$ satisfies property $\mathcal{P}$. | *x erfüllt Eigenschaft $\mathcal{P}$.* |
| in average | on average | *im Durchschnitt* |
| estimation | estimate | *Abschätzung* |
| composed number | composite number | *zusammengesetzte Zahl* |
| with the help of | using | *mit Hilfe von* |
| surely | clearly | *sicher, bestimmt* |
| monotonously increasing | monotonically incr. | *monoton steigend* |
| (Actually, in most cases 'increasing' is just fine.) | | |

In contrast, if we add a comma the sentence reads

> For the elements in $\mathcal{B}$, which satisfy property (A), we know that equation (37) holds.

Now the relative clause is *non-defining* – it just mentions in passing that all elements in $\mathcal{B}$ satisfy property (A). The main clause states that equation (37) holds for *all* elements in $\mathcal{B}$. See the difference?

## 3.4 Things you (usually) don't say in English – and what to say instead

Table 3.1 lists some common mistakes and alternatives. The entries should not be taken as gospel – they don't necessarily mean that a given word or formulation is wrong under all circumstances (obviously, this depends a lot on the context). However, in nine out of ten instances the suggested alternative is the better word to use.

Chapter 4

# Typography

## 4.1 Punctuation

**Rule 4.1** Use opening (') and closing (') quotation marks correctly.

In LaTeX, the closing quotation mark is typed like a normal apostrophe, while the opening quotation mark is typed using the French *accent grave* on your keyboard (the *accent grave* is the one going down, as in *frère*).

Note that any punctuation that *semantically* follows quoted speech goes inside the quotes in American English, but outside in Britain. Also, Americans use double quotes first. Oppose

> "Using 'lasers,' we punch a hole in ... the Ozone Layer," Dr. Evil said.

to

> 'Using "lasers", we punch a hole in ... the Ozone Layer', Dr. Evil said.

**Rule 4.2** Use hyphens (-), en-dashes (–) and em-dashes (—) correctly.

A hyphen is only used in words like 'well-known', '3-colorable' etc., or to separate words that continue in the next line (which is known as hyphenation). It is entered as a single ASCII hyphen character (-).

To denote ranges of numbers, chapters, etc., use an en-dash (entered as two ASCII hyphens --) with no spaces on either side. For example, using Equations (1)–(3), we see...

As the equivalent of the German *Gedankenstrich*, use an en-dash with spaces on both sides – in the title of Section 3.4, it would be wrong to use a hyphen instead of the dash. (Some English authors use the even longer emdash (—)

instead, which is typed as three subsequent hyphens in LATEX. This emdash is used without spaces around it—like so.)

## 4.2 Spacing

**Rule 4.3** Do not add spacing manually.

You should never use the commands \\ (except within tabulars and arrays), \␣ (except to prevent a sentence-ending space after Dr. and such), \vspace, \hspace, etc. The choices programmed into LATEX and this style should cover almost all cases. Doing it manually quickly leads to inconsistent spacing, which looks terrible. Note that this list of commands is by no means conclusive.

**Rule 4.4** Judiciously insert spacing in maths where it helps.

This directly contradicts Rule 4.3, but in some cases TEX fails to correctly decide how much spacing is required. For example, consider

$$f(a, b) = f(a + b, a - b).$$

In such cases, inserting a thin math space \, greatly increases readability:

$$f(a, b) = f(a + b, \, a - b).$$

Along similar lines, there are variations of some symbols with different spacing. For example, Lagrange's Theorem states that $|G| = [G : H]|H|$, but the proof uses a bijection $f : aH \to bH$. (Note how the first colon is symmetrically spaced, but the second is not.)

**Rule 4.5** Learn when to use \␣ and \@.

Unless you use 'french spacing', the space at the end of a sentence is slightly larger than the normal interword space.

The rule used by TEX is that any space following a period, exclamation mark or question mark is sentence-ending, except for periods preceded by an upper-case letter. Inserting \ before a space turns it into an interword space, and inserting \@ before a period makes it sentence-ending. This means you should write

```
Prof.\ Dr.\ A. Steger is a member of CADMO\@.
If you want to write a thesis with her, you
should use this template.
```

which turns into

Prof. Dr. A. Steger is a member of CADMO. If you want to write a thesis with her, you should use this template.

The effect becomes more dramatic in lines that are stretched slightly during justification:

Prof. Dr. A. Steger is a member of CADMO. If you

**Rule 4.6** Place a non-breaking space (˜) right before references.

This is actually a slight simplification of the real rule, which should invoke common sense. Place non-breaking spaces where a line break would look 'funny' because it occurs right in the middle of a construction, especially between a reference type (Chapter) and its number.

## 4.3   Choice of 'fonts'

Professional typography distinguishes many font attributes, such as family, size, shape, and weight. The choice for sectional divisions and layout elements has been made, but you will still occasionally want to switch to something else to get the reader's attention. The most important rule is very simple.

**Rule 4.7** When emphasising a short bit of text, use \emph.

In particular, *never* use bold text (\textbf). Italics (or Roman type if used within italics) avoids distracting the eye with the huge blobs of ink in the middle of the text that bold text so quickly introduces.

Occasionally you will need more notation, for example, a consistent typeface used to identify algorithms.

**Rule 4.8** Vary one attribute at a time.

For example, for WEIRDSORT we only changed the shape to small caps. Changing two attributes, say, to bold small caps would be excessive (LaTeX does not even have this particular variation). The same holds for mathematical notation: the reader can easily distinguish $g_n$, $G(x)$, $\mathcal{G}$ and $\mathsf{G}$.

**Rule 4.9** Never underline or uppercase.

No exceptions to this one, unless you are writing your thesis on a typewriter. Manually. Uphill both ways. In a blizzard.

## 4.4   Displayed equations

**Rule 4.10** Insert paragraph breaks *after* displays only where they belong. Never insert paragraph breaks *before* displays.

LATEX translates sequences of more than one linebreak (i.e., what looks like an empty line in the source code) into a paragraph break in almost all contexts. This also happens before and after displays, where extra spacing is inserted to give a visual indication of the structure. Adding a blank line in these places may look nice in the sources, but compare the resulting display

$$a = b$$

to the following:

$$a = b$$

The first display is surrounded by blank lines, but the second is not. It is bad style to start a paragraph with a display (you should always tell the reader what the display means first), so the rule follows.

**Rule 4.11** Never use `eqnarray`.

It is at the root of most ill-spaced multiline displays. The *amsmath* package provides better alternatives, such as the `align` family

$$f(x) = \sin x,$$
$$g(x) = \cos x,$$

and `multline` which copes with excessively long equations:

$$P[X_{t_0} \in (z_0, z_0 + dz_0], \ldots, X_{t_n} \in (z_n, z_n + dz_n]]$$
$$= \nu(dz_0) K_{t_1}(z_0, dz_1) K_{t_2 - t_1}(z_1, dz_2) \cdots K_{t_n - t_{n-1}}(z_{n-1}, dz_n).$$

## 4.5 Floats

By default this style provides floating environments for tables and figures. The general structure should be as follows:

```
\begin{figure}
  \centering
  % content goes here
  \caption{A short caption}
  \label{some-short-label}
\end{figure}
```

Note that the label must follow the caption, otherwise the label will refer to the surrounding section instead. Also note that figures should be captioned at the bottom, and tables at the top.

The whole point of floats is that they, well, *float* to a place where they fit without interrupting the text body. This is a frequent source of confusion and changes; please leave it as is.

**Rule 4.12** Do not restrict float movement to only 'here' (h).

If you are still tempted, you should avoid the float altogether and just show the figure or table inline, similar to a displayed equation.

Chapter 5

# Example Chapter

Dummy text.

## 5.1 Example Section

Dummy text.

### 5.1.1 Example Subsection

Dummy text.

#### Example Subsubsection

Dummy text.

**Example Paragraph** Dummy text.

*Example Subparagraph* Dummy text.

Appendix A

# Dummy Appendix

You can defer lengthy calculations that would otherwise only interrupt the flow of your thesis to an appendix.

# Bibliography

[1] Frederic Chazal, Vin de Silva, and Steve Oudot. Persistence stability for geometric complexes, 2013.

[2] H. Edelsbrunner and J. Harer. *Computational Topology: An Introduction*. Applied Mathematics. American Mathematical Society, 2010.

**ETH**

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

# Declaration of originality

The signed declaration of originality is a component of every semester paper, Bachelor's thesis, Master's thesis and any other degree paper undertaken during the course of studies, including the respective electronic versions.

Lecturers may also require a declaration of originality for other written papers compiled for their courses.

---

I hereby confirm that I am the sole author of the written work here enclosed and that I have compiled it in my own words. Parts excepted are corrections of form and content by the supervisor.

**Title of work** (in block letters):

**Authored by** (in block letters):
*For papers written by groups the names of all authors are required.*

**Name(s):**                                **First name(s):**

With my signature I confirm that
- I have committed none of the forms of plagiarism described in the 'Citation etiquette' information sheet.
- I have documented all methods, data and processes truthfully.
- I have not manipulated any data.
- I have mentioned all persons who were significant facilitators of the work.

I am aware that the work may be screened electronically for plagiarism.

**Place, date**                             **Signature(s)**

*For papers written by groups the names of all authors are required. Their signatures collectively guarantee the entire content of the written paper.*