

United States macroeconomic Analysis with BallMapper

*Ryan Johnson**
University Here

Student: *johnson.ryan1019@gmail.com*
Mentor: *scook25@alaska.edu*

KEYWORDS

Topological Data Analysis; Economics; BallMapper; Macroeconomics;

ABSTRACT

This paper set outs to explore the general use of Ball Mapper (BM) in data analysis as well as add to the small body of current research applying Ball Mapper. We carry out this analysis with economic data from three United States sources: the Bureau of Economic Analysis (BEA), the Federal Reserve Bank (FRB), and Bureau of Labor Statistics. We aim to show the reader how Ball Mapper is used for exploratory data analysis, interpretations of BM graphs, and the stability of Ball Mapper Results.

Reading Key [Removed after Review]

- [?] - Need reference here or as of writing seems like an appropriate place for reference.
- [fn?] - Would a footnote be helpful?
- [any words] - phrasing to keep writing momentum but think it needs changed.

Outline:

- Methods and Procedures
 - Procedures
 - * Literature
 - * Data
 - * Mapper/BallMapper
 - Methods
 - * Descriptive Analysis of BM Graph
 - Groupings
 - Statistical correlation and other known interpretations
 - difference epsilon testing
 - * time travel - coloring by year
 - * robustness check with randomization
 - * other colorings?
 - Results
- Discussion
- Conclusion
- References

INTRODUCTION

In this analysis, we are going to examine the macroeconomic landscape of the United States of America (U.S.) through a lens of an [adjacent mapper-inspired algorithm] called BallMapper (BM). This methodology was developed by Paweł Dłotko in 2019, algorithm is of interest to us because BallMapper significantly reduces the parameters to create a topological map.¹ BallMapper simplifies traditional Mapper by reducing the need for the user to pass data through multiple functions each with their own parameters.² One benefit to this reduction in parameters is that it removes some of the barriers to entry to learning Mapper (or Mapper-like algorithms), and more generally, Topological Data Analysis (TDA). Those who have a mathematical background, introductory course in Topology, or conceptual knowledge of data science methods and algorithms will fare much easier. However, if you are not equipped with any of the aforementioned knowledge, TDA will seem unnecessarily complicated for data analysis.

What follows starts by summarizing our literature review of the traditional Mapper algorithm followed by the BallMapper algorithm, and its differences. We then have a few notes about current research done with BallMapper, and conclude our literature section talking about [...]. We then move to describing our data sources, how we obtained them, and why. We then continue to address any data challenges we faced, calculations needed for analysis, and description of the final data set.

At the end of our procedures we begin by explaining some of the theory behind Mapper and BallMapper, and segue into our methodology: exploratory data analysis, changing algorithm parameters, and provide some tests for robustness of the algorithm.

The [key takeaways] that we want the reader to leave with is, first, to show that BallMapper is another useful tool in the world of Topological Data Analysis. Add to the small body of literature on BallMapper and its applications. Be an accessible paper to interested in TDA, Mapper and adjacent algorithms such as BallMapper. And [...].

¹: <https://arxiv.org/pdf/1901.07410>

¹Need to explain what Topological map is?

²add reference, seems appropriate here.

METHODS AND PROCEDURES

Literature

TDA is a relatively new field and is increasingly being applied to various areas of research such as genetics, financial fraud, biology, and image detection. From the small body of literature we have reviewed, TDABM has removed many of the parameters that are required for more traditional methods of TDA. Generally speaking, traditional TDA contains four broad steps, each with multiple user-defined parameters. Conversely, TDABM has reduced this process to an the user selecting their data, a coloring variable and an epsilon value – details on this follow below. What should be noted, although TDABM reduces the number of steps needed to produce similar outcomes of traditional TDA methods, we lose control of being able fine tune out outputs. We also find that interpreting results becomes more difficult. However, TDABM is still in its infancy, so there is not a large body of research on interpretation. (Paweł Dłotko 2019) One major motivation for this paper is a paper written by Paweł Dłotko, Rudkin, and Qiu (2019) that applied TDABM to a global macroeconomic dataset to compare how countries have evolved over time, their transformation from the Great Depression Era, and various views on wealth and inequality. We could not find other TDA literature found specifically focused on the macroeconomic economy of singular countries. Hence, our topic of choice.

Sources

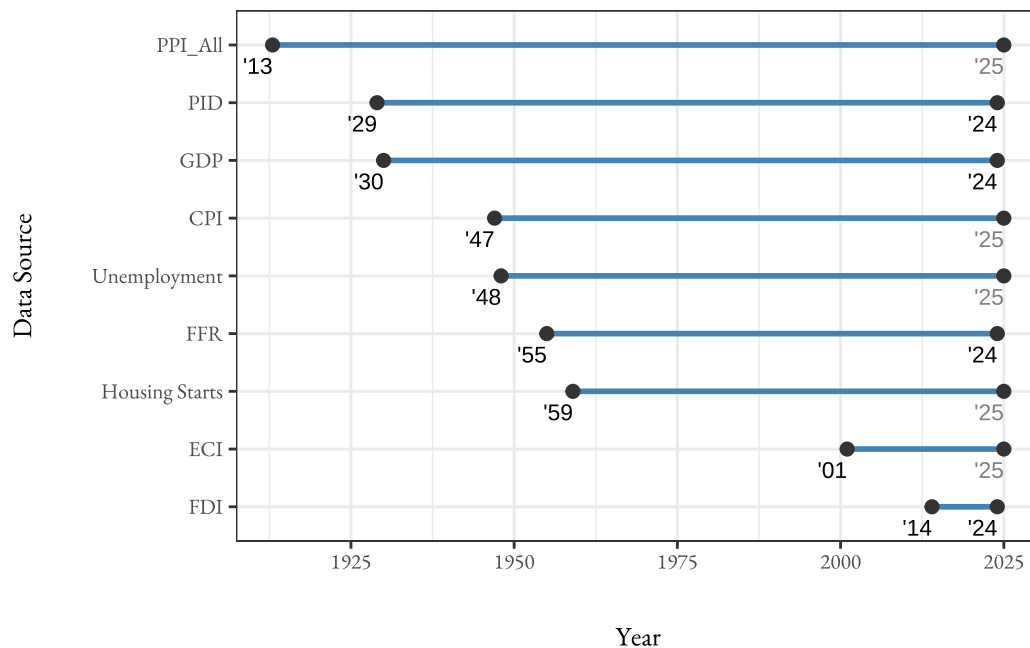


Figure 1: Span of Years by Economic Data Source

This paper relied on three publicly available data from US government sources. The Bureau of Economic Analysis (BEA), The Federal Reserve Board (FRB), and The Bureau of Labor Statistics (BLS). The data were gathered using R using two application programming interfaces (APIs): one for data from the BEA, and the other from the Federal Reserve Economic Data (FRED) API. FRED aggregates data from national and international sources, as well as public and private sources. Relevant to this paper, we used additional data, recession dates, based on business cycle contractions and expansions provided by the National Bureau of Economic Research.

These agencies were selected because of they are authoritative sources for U.S. economic data. Their widespread use in both the private and public sector gives us high confidence in the accuracy and integrity of the data.³

³Hughes-Cromwick, Ellen, and Julia Coronado. 2019. "The Value of US Government Data to US Business Decisions." *Journal of Economic Perspectives* 33 (1): 131201346. DOI: 10.1257/jep.33.1.131

While both quarterly and annual data were collected, the following analysis is focused on an annual timescale. The choice of an annual time frame was to keep consistency with other literature applying BallMapper to economic data on a yearly basis.

Looking at the time ranges in Figure 1, we elected to exclude Employment Cost Index (ECI) and Foreign Direct Investment (FDI) due to their limited available years. Additionally, we took a subset of years from the Producer Price Index - All Commodities (PPI_All), Personal Income and its Disposition (PID), and Gross Domestic Product (GDP). This allowed us to keep as much data as possible for analysis while also maintaining a range of notable economic events.

TDA BallMapper

TDABM, and in general TDA, believes that data contains a “shape” to it. It pulls ideas from the branch of Mathematics, Topology. At a high-level, Topology focuses on the properties of geometric shapes when you deform them without breaking them, e.g. bend, twist, scrunch. The following is a basic overview of TDABM ⁴.

BallMapper Algorithm

Before going over the algorithm there are a few For the following algorithm some useful definitions are as follows: - We will use the word pointcloud instead of dataset. - Conventionally the Euclidean distance metric is used; let $d = \text{dist}()$ is used as our distance metric, denote ϵ The core steps of the TDABM algorithm are as follows:

1. From your dataset (pointcloud) select a random point, α_n , and draw circle with radius epsilon, ϵ .
2. For some distance metric where d is the distance from, any points inside the circle will be associated with point α_n .
3. Repeat steps 1 - 3 until there are no more points to select and you have a set $\alpha_1, \dots, \alpha_n$.
4. Draw an edge between α_i and α_j if they share a point(s) for our set $\alpha_1, \dots, \alpha_n$.

Although the above list five steps, it should be noted that this is just the algorithm and is all done with one function call.⁵ From the steps above, the only thing the user of TDABM is responsible for is creation of a pointcloud and choosing an appropriate ϵ .

Procedures

Pointcloud Construction

A primary challenge we encountered when selecting variables for analysis was the variability of time period coverage. Some data spanned from Depression Era 1929 all the way up to 2024. Looking at our final nine variables (see [?@fig-corr](#)) and the range of years available by by source ([?@fig-years](#)span), we limited our years to 1959 - 2024 due to Housing Starts data having the smallest range of years.

Some of the data was only available by month. In these cases we calculated annual changes by finding the difference between January of the current year and January of the previous year. Note, this reduced the number of years in our data set by one. We also standardized the data by converting non-rate based data to percent changes from the previous year to ensure consistency across all variables (Equation 1). One feature of TDABM is that it can be sensitive to data with large scale differences. This not only reduces how precise you can be when selecting a parameter, but can also significantly increase computation time.

Looking at our correlation table ([?@fig-corr](#)), we observe moderate to strong correlations between the pairwise combinations of GDP, PCE, Ig, and Imports. We also see a strong correlation between CPI and Interest Rates. Since these correlations are moderate to high, it's worth noting that the shape of our output may be reflected by these correlations. They also could express themselves with similar coloring of TDABM graphs (Pawel Dłotko, Rudkin, and Qiu 2019).

⁴Some very good examples and explanations can be found in (Dłotko, Qiu, and Rudkin 2022) and (Pawel Dłotko, Qiu, and Rudkin 2019) which touch on the finer details of TDABM algorithm. A very good illustrative example can be seen in the pre-print (P. Dłotko, Qiu, and Rudkin 2021).

⁵There are packages available both in R and Python under the names BallMapper and pyBallMapper, respectively.

Cloud Building

In our Data section, we noted that some calculations were needed before our analysis. Additionally, for the Federal Funds Rate, the Unemployment Rate, and CPI we had to take extra steps and considerations while constructing out pointcloud. Since the Federal Funds Rate and the Unemployment Rate are already a rate-based measurements, instead of taking a percent change from the previous year, we instead took the standard non-weight average rate by each year. For CPI, this data was given on a monthly-basis, we calculated it using our [General Percent Change Formula](#) with $k = 12$.

The analysis reportes here are guided by Paweł Dłotko, Rudkin, and Qiu (2019)'s approach. We will be referencing our TDABM output in Figure 2 for the remainder of this section.

Results

Interpretation of TDA BallMapper Graphs, An Abridged Version

In Figure 2 we are presented with what looks like a graph, in the mathematical sense; a display of points (nodes) and edges –some connected while others are not connected at all (satellites). If we do not change the poincloud or ϵ of our TDABM graph, we will get the same “shape” of the graph regardless of the coloration we choose.

Nodes are colored depending on the users interest. In (?), we are focusing color by year. Coloration calculated by taking the average of the data that is “inside” each node, in our case the average value of the years. Further, if there is only one data point in a node (i.e. a satellite points), the node will appear small and the coloration will reflect as seen on the legend. Conversely, if the node is large, this indicates more data is “inside” the node, and the average value could be reflecting the mean values of a large or small variance.

General Figure 2 Remarks

We initially notice a large connected component⁶ (component A) on the right side of our graph, as well as a smaller connected component (component B) on the lower left side. Noting the coloration, there are three main sections: the right, bottom, and upper-left. On the bottom and sweeping upward to the left we see there are multiple satellite points. These can be of interest because they may indicate outliers in our data.⁷ Focusing back towards component A, we notice that the two largest nodes, 28 and 22, which seem to represent the late 2010s and the dot-Com boom, respectively. We also note that there is a arm coming off the upper-right portion of component A, as well as smaller, lollipop, features emerging at the bottom and top of component A.⁸

Time Travel

Looking at Figure 2, we see that our graph is colored by year.⁹ In this section we will only highlight a few observations in detail due to the nature of this paper. Additional areas for investigation can be found in the Conclusion.

Mentioned above, we observed the satellite points in the lower half of the graph. If we look at the upper-right quadrant, we see that it consists of nodes {1, 16, 25, 27}. Looking at [?@tbl-covidplus](#) we see that it covers the 2020 Pandemic and the Financial Crisis and its aftermath. Seeing a general coloration from the recent 10 - 15 years we would expect to see some of the recent outlier economic events here. Indeed, node 25 consists of two of the most recent economic downturns, 2009 and 2020. However, what is interesting is that node 27 consists of three years closely following the 2008-2009 Financial crisis. However, the years following the 2020 COVID Pandemic can be found scattered across nodes {6, 8, 24, 28, 29}. Looking at node 27 we find that the common thread is high unemployment. In contrast, the years following the Pandemic are not the same year to year.¹⁰

Shifting our focus towards component A, it consists of data most similar to the economic years between the 1990s - 2010s. Meanwhile, looking at component B and the general lower area of our graph, we see that it represents the Great Inflation period of the 1970s - early

⁶In graph theory, a connected component is one in which there exists a path from a node to every other node for a set of nodes and edges.

⁷We have not found any literature yet on whether this observation is empirically true.

⁸The observations mentioned above good starting places for interpreting BallMapper graphs.

⁹Supplementary graphs displaying coloration based on different variables can be found in [Supplement_1.pdf](#).

¹⁰2021: An outlier, we see higher levels of unemployment and inflation. We see high PCE, Ig, and imports; 2022: Nodes {8, 28} we see that there is higher than normal inflation, export, and imports; 2023: Nodes {6, 24}, we see low inflation and low unemployment.

1980s [investopedia article]. Looking at nodes 12 – 20, we get exactly ten years of data¹¹ where there was know high inflation, high unemployment, and in general know to be bad time economically for the US (?@tbl-greatinflation). Something of interest to note is that nodes 14 and 15 consist of 1976 - 1978. This three year period based on data seems to show more normal economic conditions. On the other hand, we see that all the other years in this grouping (nodes 12 – 20) show signs of a struggling economy in some way.

Discussion

This paper serves as a proof-of-concept of the usefulness of TDABM as a methodology for exploratory data analysis. By no means is this a totally comprehensive method to gain insight into data, but instead TDA BallMapper proves itself to be useful tool when employed in addition to other traditional analyses.

In our Results section we briefly went over two ways to interpret TDABM graphs. As seen in many other, longer papers written by Dlotko and colleagues, there is a lot of room to expound our TDABM graph here, though that is beyond the scope of this paper. In general, though, we can see that TDBM can give insight into our data that we might not see otherwise. Seeing that there is data that is not like the others gives an indication that there is something to investigate. Consider our comparison of the 08' - 09' financial crisis and the 2020 COVID Pandemic: We see learned that the years following the financial crisis were similar economically. Another way to view this is, clustering reveals that the financial crisis was a singular problem that spurred the collapse of the financial system. Conversely, the COVID pandemic was a singular issue but exposed many different weaknesses in our current world.

For the future development of TDABM, and TDA in general, there are many possibilities of future research and one glaring downside. The downside to TDA is that it is a high barrier to entry to understand the methodologies inner workings. Those with a technical background will have an easier time, but nonetheless a higher barrier than methods such as linear regressions. Subject wise, TDA has an advantage in some of the social sciences because it can be somewhat of a bridge between qualitative and quantitative analyses. TDABM has been used to analyze the Brexit Vote data, seeking to understand quantitative and qualitative motivations behind its outcome Rudkin et al. (2023). TDA also has the advantages of dimension reduction which could be useful to the life-sciences. Fields such as genetics, chemistry, and biology have many fields of data, and being able to consolidate it into an understandable form could benefit the body of knowledge greatly.

Conclusion

Endnotes

General Percent Change Formula

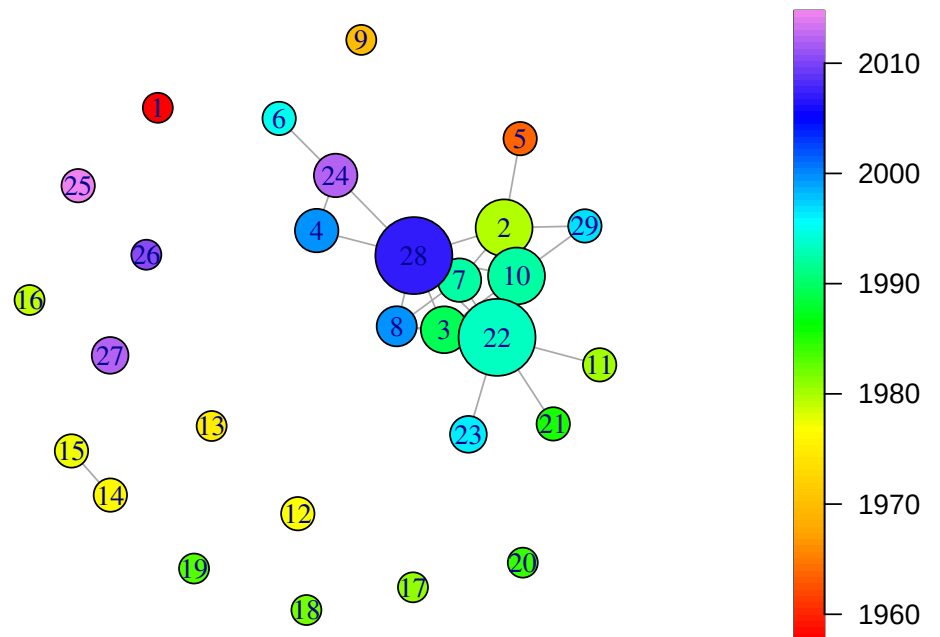
$$\Delta p_n = \frac{p_n - p_{n-k}}{p_{n-k}} \text{ for some } p \in P \text{ and } n > 0 \quad (1)$$

where P is a column of our data with x_1, \dots, x_n observations and k is a lag integer such that $k > 0$.

Appendix

¹¹1974 - 1984

BallMapper output Colored by Year



Epsilon: 0.45

Figure 2: Data is normalized on a [0,1] scale due to all variables not being normally distributed See Supplement_2.pdf.

References

- Dłotko, Paweł, Wanling Qiu, and Simon Rudkin. 2022. “Topological Data Analysis Ball Mapper for Finance.” <https://arxiv.org/abs/2206.03622>.
- Dłotko, Paweł. 2019. “Ball Mapper: A Shape Summary for Topological Data Analysis.” <https://arxiv.org/abs/1901.07410>.
- Dłotko, Paweł, Wanling Qiu, and Simon Rudkin. 2019. “Financial Ratios and Stock Returns Reappraised Through a Topological Data Analysis Lens.” <https://arxiv.org/abs/1911.10297>.
- Dłotko, Paweł, Simon Rudkin, and Wanling Qiu. 2019. “Topologically Mapping the Macroeconomy.” <https://arxiv.org/abs/1911.10476>.
- Dłotko, P., W. Qiu, and S. T. Rudkin. 2021. “Financial Ratios and Stock Returns Reappraised Through a Topological Data Analysis Lens.” *The European Journal of Finance* 30 (1): 53–77. <https://doi.org/10.1080/1351847x.2021.2009892>.
- Rudkin, Simon, Lucy Barros, Paweł Dłotko, and Wanling Qiu. 2023. “An Economic Topology of the Brexit Vote.” *Regional Studies* 58 (3): 601–18. <https://doi.org/10.1080/00343404.2023.2204123>.

Data Sources

- Bureau of Economic Analysis
 - [Real Exports and Imports](#)
 - [New Foreign Direct Investment in the United States](#)
 - [Real GDP](#)
 - [Personal Income and Its Disposition](#)
- Federal Reserve Bank
 - [Selected Interests Rates](#)
- Bureau of Labor Statistics
 - [CPI](#)
 - [Employment Cost Index](#)
 - [PPI](#)
 - [Housing Starts](#)
 - [Unemployment Rate](#)

Old Intro Material

Dłotko demonstrates that, mathematically, TDABM produces a near similar output to that of traditional TDA methods via a factor map (Paweł Dłotko 2019). TDA BallMapper and closely follows other papers that have applied TDABM. We will begin with a general overview of TDABM: current published literature and useful concepts to know. We then will describe our data and reasoning behind their selection. Our Methodology and Results will touch on more theoretical concepts behind TDABM and then adress a couple of findings, respectively.