
Exploring Directed vs. Undirected Graph-Based Topological Data Analysis of Transformer Attention Maps

(Selected Topics in DS 2024 Course)

Kamil Garifullin¹ Ilya Trofimov¹

Abstract

This project investigates the application of directed graph analysis to transformer attention maps, comparing it with traditional undirected graph methods. Using persistent homology, we examine the topological features of attention maps to evaluate their impact on downstream classification tasks. The results of this research show the efficacy of utilizing directed graph representations, providing valuable insights into the dynamic flow of information within transformer models.

Github repo: github.com/TDA-for-Transformers

1. Introduction

In the realm of natural language processing (NLP), transformer models have revolutionized the field by effectively capturing contextual relationships between tokens in textual data. Central to these models are attention mechanisms, which generate attention maps highlighting the interactions between tokens across different layers and heads.

Being rich in information, attention maps often require sophisticated analysis techniques to fully exploit their potential. One such promising direction is the application of persistent homology to extract meaningful signals from these attention maps (Tulchinskii E., 2022), (Cherniavskii D., 2022). By characterizing the topological structure of the attention graphs, persistent homology allows us to understand the underlying relationships between lexemes and their importance in the context of various NLP tasks.

However, existing research works has focused on analyzing attention maps as undirected graphs, ignoring the inherent

directional nature of the relationships encoded in them. This oversight prevents us from obtaining valuable information about the flow and dynamics of information propagation in transformer models.

Fortunately, tools and methodologies exist for computing homologies of directed graphs. This gives us the opportunity to explore the advantages of using attention map representations based on directed graphs.

Therefore, the main goal of this project is to systematically compare the performance and efficiency of persistent homology-based feature extraction from attention maps represented as directed graphs with their non-directed counterparts.

The project plan consists of the following steps:

1. **Dataset and Model Selection:** took the dataset and a pre-trained transformer model (BERT) suitable for the downstream classification task from the article (Cherniavskii D., 2022).
2. **Pipeline Development:** developed a pipeline for extracting topological features from attention maps, implementing two variants: directed and undirected - to cater to different graph representations.
3. **Performance Comparison:** compared the classification performance on the features obtained by the two different methods, drawing conclusions regarding the effectiveness of using features to improve model performance.
4. **Feature Importance Analysis:** analyzed the importance of features obtained in two different ways (directed and undirected).

2. Related works

The study of attention mechanisms in natural language processing (NLP) and their role in encoding linguistic knowledge has attracted considerable attention in recent years. Various studies have focused on understanding how attention heads contribute to different linguistic tasks and phe-

¹Skolkovo Institute of Science and Technology, Moscow, Russia.. Correspondence to: Kamil Garifullin <kamil.garifullin@skoltech.ru>.

nomena. Below we review related works, on the study of the topology of attention maps.

- **Acceptability Judgements via Examining the Topology of Attention Maps.** The work presented by (Cherniavskii D., 2022) introduces one of the first attempt to analyze attention heads in the context of linguistic acceptability (LA) using topological data analysis (TDA), that enabled to examine graph representations of transformers' attention maps. This work demonstrated that features obtained by TDA can be used for further acceptability classification task with results that outperform the established baselines in three Indo-European languages (English, Italian, and Swedish) and confirmed the hypothesis that grammatical phenomena can be encoded through topological properties of the attention map.

- **Topological Data Analysis for Speech Processing.** In the second paper (Tulchinskii E., 2022) the authors explore the application of topological data analysis (TDA) to tackle speech classification problems and understand how a pre-trained speech model, HuBERT, works internally. To do this, authors introduced new features derived from Transformer attention maps. Surprisingly, with the usage of these new features a simple linear classifier performs better than a classifier specifically fine-tuned for the task. Moreover, authors found that on one dataset called CREMA-D, new TDA features set a new state-of-the-art performance with an accuracy of 80.155. Overall, this research shows that TDA could be a really useful approach for analyzing speech, especially in tasks where predicting structure is important.

- **Artificial Text Detection via Examining the Topology of Attention Maps.** In this article authors (Kushnareva L., 2022) introduced three new types of interpretable topological features for natural language processing (NLP) tasks, utilizing Topological Data Analysis (TDA). Authors' empirical findings demonstrate that features derived from the BERT model surpass traditional count- and neural-based approaches by up to 10% across three standard datasets. Moreover, through probing analysis, authors observed that TDA features are sensitive to both surface-level and syntactic properties of language. Overall, results of this research highlight the potential of TDA in enhancing NLP tasks, particularly those involving surface and structural information.

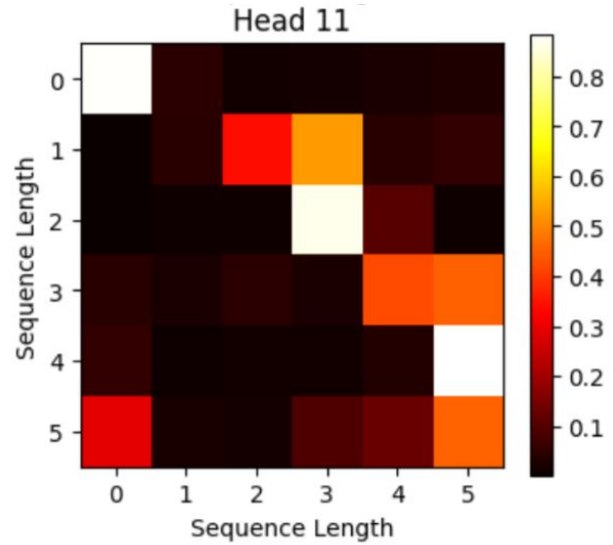


Figure 1. An attention map obtained from the 1st layer and 11th head

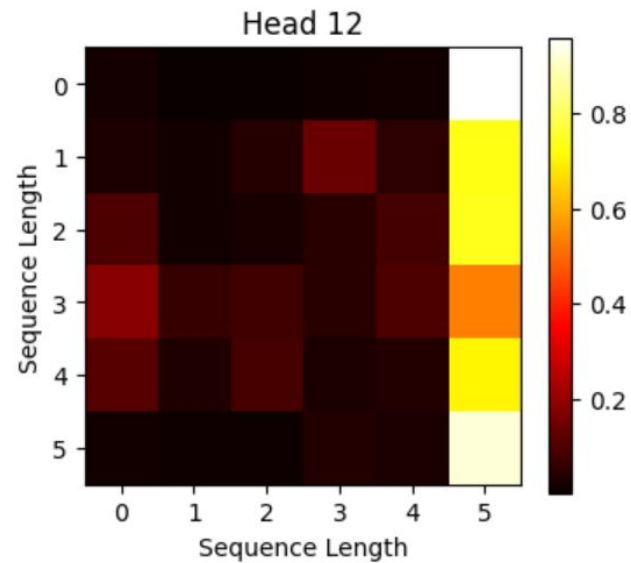


Figure 2. An attention map obtained from the 9th layer and 12th head

3. Methodology. Barcode features

3.1. Attention maps

Attention maps represent a fundamental concept in the realm of neural network architectures, particularly within models like Transformers, where they play a pivotal role in capturing dependencies and relationships between different elements of input sequences. In tasks like machine translation, sentiment analysis, and text generation, attention maps provide valuable context for understanding model predictions

	En-BERT + linear layer(baseline)	En-BERT + undirected TDA	En-BERT + directed TDA
Accuracy	75.0	?	?
Precision	74.6	?	?
Recall	75.6	?	?

Table 1. Results of experiments (classification) for various models: En-BERT with linear layer(baseline), En-BERT with undirected TDA and En-BERT with directed TDA

and decision-making processes. By visualizing attention patterns, researchers can identify linguistic structures, semantic relationships, and syntactic dependencies captured by the model.

Attention maps are generated dynamically as the model processes input sequences. These maps depict the distribution of attentional weights or probabilities assigned to each token in the input sequence relative to every other token. The process involves computing pairwise attention scores between tokens, often followed by a softmax normalization to obtain attention probabilities. These probabilities form the basis of the attention map, which is typically visualized as a heatmap, with intensity indicating the strength of attention.

For example, figure 1 and figure 2 show attention maps obtained for various layers and heads for the sentence: “It was raining yesterday”.

3.2. Barcode features

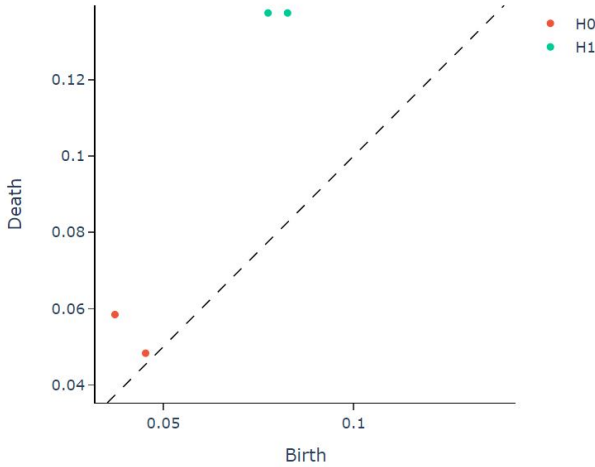


Figure 3. TDA features obtained for the sentence: “It was raining yesterday”

In this research we study barcode features (features that were extracted from barcodes) inferred from directed graphs and undirected graphs.

To calculate this features we used a high-performance topological machine learning toolbox in Python built on top of scikit-learn [giotto-ai](#). Using this tool one can compute

topological summaries, called persistence diagrams, from collections of point clouds or weighted graphs.

In the analysis of each text sample, we perform calculations to obtain the barcodes associated with the first two persistent homology groups, designated as H_0 and H_1 , across every attention head of the BERT model. Our research covers several key characteristics of these barcodes, including:

- Total aggregate length of bars
- Mean length of bars
- Bars lengths variance
- Birth and death times of the longest bar
- Overall count of bars present
- Barcode’s entropy assessment

For example, the TDA features obtained for the sentence: “It was raining yesterday”, are shown in figure 3.

4. Experiments and Results

4.1. Dataset

In our analysis, we selected the IMDb Movie Reviews dataset as the basis for our investigation. This dataset is widely known for its extensive collection of film reviews. IMDb contains movie reviews labeled as positive or negative sentiment. It’s commonly used for sentiment analysis tasks. For our task, we intend to employ this dataset for binary classification purposes, specifically categorizing instances into two classes: positive and negative.

4.2. Baseline

For the baseline in this work, a combination of BERT-base-uncased and one linear layer was used. We leveraged the BERT model to generate representations for our textual data. Specifically, we utilize the pooled output, often associated with the special [CLS] token, which serves as a condensed representation of the entire input sequence. This pooled output encapsulates the semantic information learned by the BERT model from the input text. And after obtaining the pooled outputs from the pre-trained BERT model, the next

step in our methodology involved training a classifier using these representations as features. As a result, the accuracy of our baseline solution was 75% on test set.

4.3. Undirected TDA features

5. Conclusion

References

- Cherniavskii D., e. a. Acceptability judgements via examining the topology of attention maps, 2022.
- Kushnareva L., e. a. Artificial text detection via examining the topology of attention maps, 2022.
- Tulchinskii E., e. a. Topological data analysis for speech processing., 2022.