TOPOLOGICAL DATA ANALYSIS FOR SPEECH PROCESSING Eduard Tulchinskii¹, Kristian Kuznetsov¹, Laida Kushnareva, Daniil Cherniavskii³, Serguei Barannikov^{1,2}, Irina Piontkovskaya, Sergey Nikolenko⁴, Evgeny Burnaev^{1,3} ¹Skolkovo Institute of Science and Technology (Russia), ²CNRS, Université Paris Cité (France),

³AIRI, Artificial Intelligence Research Institute (Russia),

⁴St. Petersburg Department of the Steklov Institute of Mathematics (Russia)

Introduction

We apply methods of Topological Data Analysis (TDA) to study the inner representation of transformer models to create efficient topologically-backend solutions for various tasks of Speech Processing and explore the inner-works of the attention. Natural Language Processing was revolutionized by transformers – models without any recurrent parts, instead relying on the attention mechanism. They also have shown out-

Our approach

- TDA methods are known to capture well **surface and structural patterns** in different types of data; they are also known to be more robust to noise in data.
- The attention maps generated by the Transformer model can be represented as weighted graphs and further efficiently investigated with TDA.
- We take a pretrained model, extract topological features from its attention maps and embeddings, then use those features to build a linear classifier

standing performance for many tasks in speech processing and other domains. However:

- For some tasks, usage of transformer internal representations yields significantly better performance than the traditional approach (via final embeddings).
- Attention allows model to focus on specific parts of data and learn complex dependencies. It is intended to emulate human cognitive attention, yet attention maps are hard to analyze.

Background: Topological Data Analysis



Fig. 1: Building the persistence barcode from the data.



 Three groups of features: algebraic features of attention matrices (e.g. means of 3 largest diagonals, measure of asymmetry), topological features of attention matrices, and topological features of embedding.

Main results

- Linear classifier built on top of topological features outperforms a fine-tuned classification head
- New state-of-the-art performance result for emotion recognition on CREMA-D
- Topological features are able to reveal functional roles of transformer heads

Model	IEMO CAP	CREMA– D	ASV Spoof	Vox Celeb1	FSDD
	Acc ↑	Acc \uparrow	$EER\downarrow$	$EER\downarrow$	EER ↑
HuBERT	64.92*	71.047	6.649	7.45	99.3
(baseline)		± 0.566			± 1.3
All layer	65.612	71.320	2.706	46.240	96.0
embs, 1st	± 1.050	± 0.479			± 0.7
All layer	69.355	76.260	1.519	8.46	97.7
embs, mean	± 1.801	± 1.148			± 0.5
Attention	69.666	79.200	2.138	30.326	98.7
features	± 1.174	± 1.240			± 0.8
Attn. &	69.955	80.155	1.946	26.443	99.6
non-attn feat	+0.972	+0.680			+0.4

Fig. 2: Our pipeline: from waveform to attention maps via transformer model (e.g. HuBERT, Wav2Vec2), then graph filtration and barcode.



Tab. 1: Experimental results. Accuracy (Acc) and equal-error-rate (EER) are reported

Studying the individual attention heads



Fig. 4: Average length of bar in H_0 barcode for the best HuBERT attention heads for two tasks: (a-c) individual model separation between human (blue) and synthetic (red) speech; (d-e) speaker separation; LA_0069/72/78 – female speakers, LA_0070/76/71 – male speakers.



Fig. 3: H_0 -barcode reflects the hierarchial structure of data. Sample H_0 -barcode and MST on different levels for head (2, 4); speech sample text: "I know it", sample phones: "sil AY1 N OW1 IH1 T sp". Nodes and bars are colored with respect to the phonemes they represent. Black dashed lines show barcode levels corresponding to the trees on the right. Separators inside the bars show levels where nodes from the same phoneme are joined to the bar's component.

Fig. 5: Heads with high correlation between $H_0^{m,pc}$ and MFCC/PLP

More pictures...

... and some additional information can be found on the project's website https://topohubert.github.io/speech-topology-webpages/

