

More Identifiable yet Equally Performant Transformers for Text Classification

Rishabh Bhardwaj • Navonil Majumder • Soujanya Poria • Eduard Hovy

Key Takeaways

- Attention-based interpretations of Transformer models can be faulty.
- Identifiability issue:** Attention weights are not unique to the Transformer's attention head. They can be made more identifiable!
 - By decreasing the size of key vector.
 - By increasing the size of value vector followed by head addition.

Motivation

Interpretability is an important aspect when it comes to **trust** a model's prediction.

- Famous **Transformer** architecture-based systems are ubiquitous.
- Multi-head self-attention** is its basic building block.

Introduction

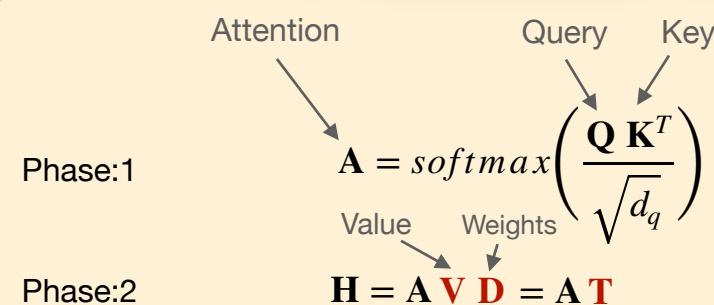
Identifiability of attention weights

"Identifiable if they can be **uniquely** determined from the head's output"
 -On identifiability in transformers [Brunner et al., 2019]

Why is it important?

Explanations based on non-unique attention weights might be **misleading**.

Attention Head



Uniqueness of A

Constraint-1

$$\tilde{\mathbf{A}}\mathbf{T} = \mathbf{0}$$

Sample
LN(T)

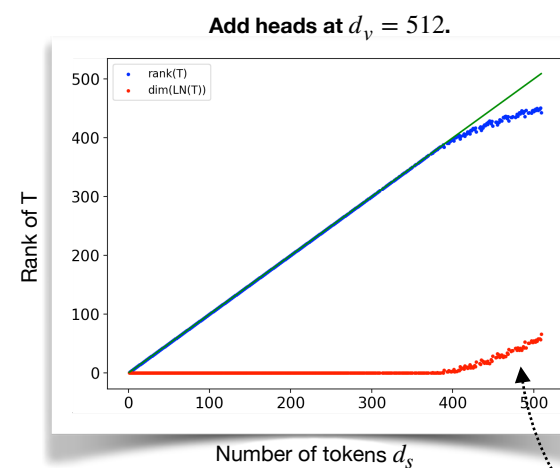
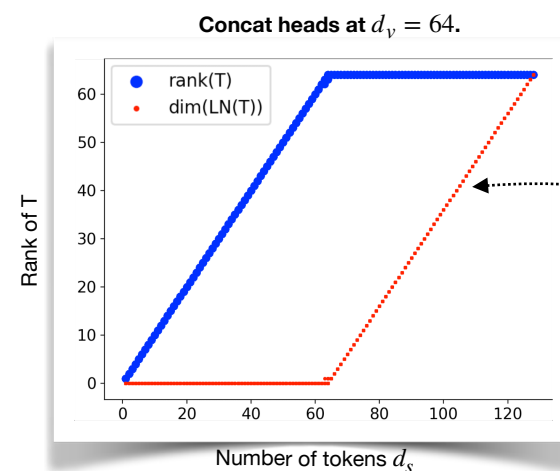
Constraint-2

$$\tilde{\mathbf{A}}\mathbf{1} = \mathbf{0}$$

Probability constraints

Constraint-3

$$(\mathbf{A} + \tilde{\mathbf{A}}) \geq \mathbf{0}$$



Role of key vector size

Let's say $\tilde{\mathbf{A}}$ satisfies above constraints.

$$\text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_q}}\right) \rightarrow (\mathbf{A} + \tilde{\mathbf{A}}). \quad (\text{Phase1 output})$$

Solutions

- We propose the following solutions to identifiability.
- Decrease the size of key vector.
 - Increase the size of value vector and add head outputs.

Classification performance of the model on ten text datasets does not vary significantly.

Paper



Code

