

# Lipschitz Normalization for Self-Attention Layers with Application to Graph Neural Networks

---

George Dasoulas   Kevin Scaman   Aladin Virmaux



# Attention

- A method.
- $\tilde{n}$  focus of structured input.

## Attention mechanism

Let  $\mathbf{P} \in \mathbb{R}^{n \times n}$  be an input matrix and a score function  $s: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ :

$$s(\mathbf{p}, \mathbf{q}) = \mathbf{p} \mathbf{q}^T \mathbf{A} - \frac{1}{2} \mathbf{p} \mathbf{p}^T \mathbf{Q} \mathbf{Q}^T$$

## Standard attention

$$s(\mathbf{p}, \mathbf{q}) = \mathbf{p} \mathbf{q}^T - \frac{1}{2} \mathbf{p} \mathbf{p}^T \mathbf{Q} \mathbf{Q}^T$$

## Transformer

$$s(\mathbf{p}, \mathbf{q}) = \mathbf{p} \mathbf{q}^T - \frac{1}{2} \mathbf{p} \mathbf{p}^T \mathbf{Q} \mathbf{Q}^T \text{ with } \mathbf{p} \mathbf{p}^T \mathbf{Q} \mathbf{Q}^T \ll \mathbf{p} \mathbf{q}^T$$

# Attention

- A method.
- $\tilde{n}$  focus of structured input.

## Attention mechanism

Let  $\mathbf{P} \in \mathbb{R}^{n \times m}$  be an input matrix and a score function  $f: \mathbb{R} \rightarrow \mathbb{R}$ :

$$a_i = \frac{\exp(\mathbf{p}_i \mathbf{q}^T)}{\sum_j \exp(\mathbf{p}_j \mathbf{q}^T)}$$

## Standard attention

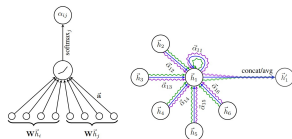
$$a_i = \frac{\exp(\mathbf{p}_i \mathbf{q}^T)}{\sum_j \exp(\mathbf{p}_j \mathbf{q}^T)}$$

## Transformer

$$a_i = \frac{\exp(\mathbf{p}_i \mathbf{q}^T)}{\sum_j \exp(\mathbf{p}_j \mathbf{q}^T)}$$

## Graph Learning and Attention

- Neighbor-wise softmax function.



Graph attention network (Veličković et al, 2018)

# Attention

- A method.
- $\tilde{n}$  focus of structured input.

## Attention mechanism

Let  $\mathbf{P} \in \mathbb{R}^{n \times d}$  be an input matrix and a score function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ :

$$f(\mathbf{p}) = \frac{\exp(\mathbf{p} \cdot \mathbf{q})}{\sum_j \exp(\mathbf{p} \cdot \mathbf{q}_j)}$$

## Standard attention

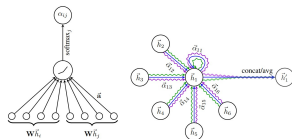
$$f(\mathbf{p}) = \frac{\exp(\mathbf{p} \cdot \mathbf{q})}{\sum_j \exp(\mathbf{p} \cdot \mathbf{q}_j)}$$

## Transformer

$$f(\mathbf{p}) = \frac{\exp(\mathbf{p} \cdot \mathbf{q})}{\sum_j \exp(\mathbf{p} \cdot \mathbf{q}_j)}$$

## Graph Learning and Attention

- Neighbor-wise softmax function.



Graph attention network (Veličković et al, 2018)

## GNN training and model depth

- Deep attention models suffer from **convergence inability** (Alon et al, 2020).
  - Oversmoothing, oversquashing.
  - Gradient explosion/vanishing not explored!
- **Gradient flow and Lipschitz continuity.**

# How to enforce Lipschitz continuity to attention layers?

### Main contributions

1. We derive **general bounds** for the Lipschitz constant of attention layers.
2. We propose a **novel normalization** for attention layers that ensures Lipschitz continuity.
3. We apply this normalization to **graph attention networks**.

# Is attention Lipschitz continuous ?

- Large scores tend to create large gradients!

## Lemma

o  $\mathbb{P}^R$

$$\sim D, \mathbf{z} \sim_o \alpha \} \text{sb} \mathbb{H} \setminus - \{p \ p \ \text{qq}\}_o \quad ? \bar{2} \}^J \} \text{p} \mathbb{S} ; 2q \sim D \quad \sim_o ; p2 ; \mathbb{S} q :$$

# Is attention Lipschitz continuous ?

- Large scores tend to create large gradients!

## Lemma

o  $\mathbb{P} \mathbb{R}$

$$\sim D, \mathbf{z} \sim_o \alpha \} \text{sb} \mathbb{H} \setminus - \{ \mathbf{p} \mathbf{p} \mathbf{q} \} \circ \quad ? \bar{2} \} \mathbf{J} \} \mathbf{p} \mathbf{8} ; 2 \mathbf{q} \sim D \quad \sim_o ; \mathbf{p} 2 ; \mathbf{8} \mathbf{q} :$$

Impact of a scalar normalization:

Normalize the score function by a scalar function  $\tilde{\mathbf{N}} : \mathbb{R} \rightarrow \mathbb{R} : \mathbf{p} \mathbf{q} \rightarrow \frac{\mathbf{p} \cdot \mathbf{q}}{\|\mathbf{p}\| \|\mathbf{q}\|}$

# Is attention Lipschitz continuous ?

- Large scores tend to create large gradients!

## Lemma

o  $\mathbb{P} \mathbb{R}$

$$\sim D, \mathbf{z} \sim_o \alpha \} \text{sb} \mathbb{H} \setminus - \{ \mathbb{P} \mathbb{P} \mathbb{Q} \mathbb{Q} \}_o \quad ? \bar{2} \}^J \} \mathbb{P} \mathbb{S} ; 2q \sim D \quad \sim_{\alpha; p2; \mathbb{S} q} :$$

## Impact of a scalar normalization:

Normalize the score function by a scalar function  $\mathbb{R} \rightarrow \mathbb{R} : \mathbb{P} \mathbb{Q} \rightarrow \frac{\mathbb{P} \mathbb{Q}}{\mathbb{P} \mathbb{Q}}$

## Theorem

©  $\forall 0 \neq \mathbb{P} \mathbb{R}$

$$(1) \} \sim \mathbb{P} \mathbb{Q} \} \mathbb{S} \propto \mathbb{P} \mathbb{Q}$$

$$(2) \}^J \} \mathbb{P} \mathbb{S} ; 2q \sim D \sim \sim_{\alpha; p2; \mathbb{S} q} \propto \mathbb{P} \mathbb{Q}$$

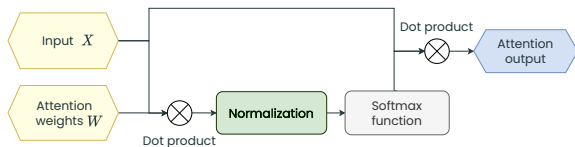
$$(3) \}^J \} \mathbb{P} \mathbb{S} ; 2q \sim D \sim_{\alpha; 1} \} \sim \mathbb{P} \mathbb{Q} \} \mathbb{P} 2; \mathbb{S} q \propto \mathbb{P} \mathbb{Q}^2$$

$$\mathbb{P} \mathbb{Q} \sim \mathbb{P} \mathbb{Q} \{ \mathbb{P} \mathbb{Q} \quad \text{©}$$

$$\text{©}_{\mathbb{P}, \mathbf{z} \mathbb{Q}} \propto \frac{\text{c}}{\text{—}} \quad ? \bar{8} :$$

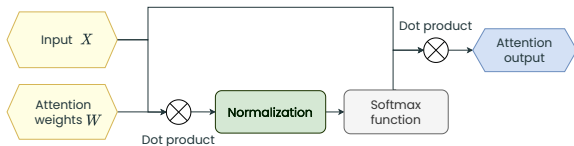


**Core Idea:** Replace original  $\sim \mathbf{p} \mathbf{q}$  with normalized  $\mathbf{p} \mathbf{q}$  in the attention model!



Pipeline of LipschitzNorm.

**Core Idea:** Replace original  $\sim p q$  with normalized  $p q$  in the attention model!



Pipeline of LipschitzNorm.

Linear score function  $p q$

$$p q = \frac{\sum_{i=1}^J \tilde{o}_i \tilde{w}_i}{\sum_{i=1}^J \tilde{o}_i}$$

Quadratic score function  $p q$

$$p q = \frac{\sum_{i=1}^J \tilde{o}_i^2 \tilde{w}_i}{\sum_{i=1}^J \tilde{o}_i^2}$$

where  $\tilde{o}_i = \sum_{j=1}^J \tilde{w}_j x_{ij}$

$$p q = \frac{\sum_{i=1}^J \tilde{o}_i^2 \tilde{w}_i}{\sum_{i=1}^J \tilde{o}_i^2}$$

# Gradient Explosion

For **deep attention** models, there is a tight connection between:  
**efficient training** and **Lipschitz continuity**.

## Lipschitz constant

Given  $\mu$  attention layers,  $\mathbf{z}, \mathbf{p}, \mathbf{q}$ ,  
 $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_{\mu-1}, \mathbf{z}_\mu$   
is Lipschitz continuous:

$$\|\mathbb{C}_p \mathbf{q}\| \propto \prod_{i=1}^{\mu} \|\mathbb{C}_p \mathbf{z}_i\|$$

- **Multiplicative effect** on the gradient flow.
- Enforcing Lipschitz continuity can **alleviate** gradient explosion.

# Gradient Explosion

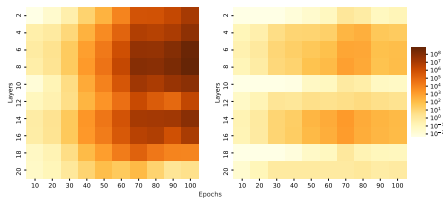
For **deep attention** models, there is a tight connection between:  
**efficient training** and **Lipschitz continuity**.

## Lipschitz constant

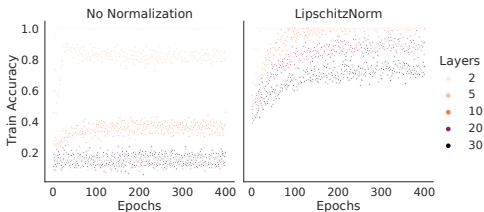
Given  $\mu$  attention layers,  $\mathbf{z}, \mathbf{p}, \mathbf{q}$ ,  
 $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_{\mu-1}, \mathbf{z}_{\mu}$   
is Lipschitz continuous:

$$\|\mathbf{z}_{\mu}\| \propto \prod_{i=1}^{\mu} \|\mathbf{z}_i\|$$

- **Multiplicative effect** on the gradient flow.
- Enforcing Lipschitz continuity can **alleviate** gradient explosion.



Gradient evolution of 25-layer GAT without/with LipschitzNorm

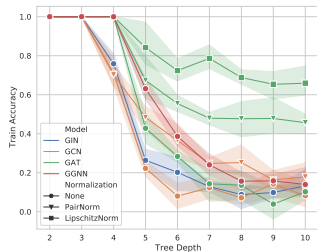


Model convergence w.r.t. model depth

# Synthetic Study

## A. Trees with increasing depth

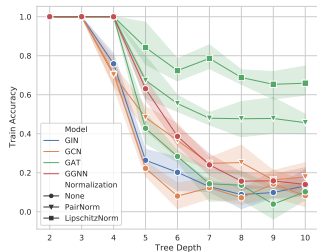
- Synthetic graph benchmark (Alon et al, 2020).
- Simulation of



Train accuracies of four GNN models.

## A. Trees with increasing depth

- Synthetic graph benchmark (Alon et al, 2020).
- Simulation of

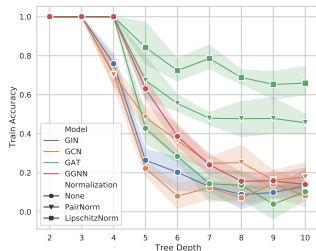


Train accuracies of four GNN models.

- \* Normalization always helps 😊
- \* LipschitzNorm for the win!

## A. Trees with increasing depth

- Synthetic graph benchmark (Alon et al, 2020).
- Simulation of



Train accuracies of four GNN models.

- \* Normalization always helps 😊
- \* LipschitzNorm for the win!

## B. Node Classification with Missing Features

(Zhao and Akoglou, 2019)

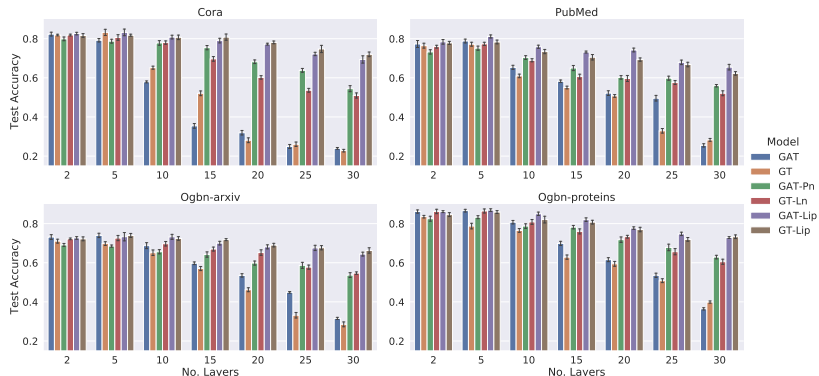
	Cora		CiteSeer		Pubmed							
	0%	100%	0%	100%	0%	100%						
GCN	82.5	1.2 (2)	58.8	3.5 (2)	69.5	2.1 (2)	31.3	2.7 (2)	77.9	1.4 (2)	44.9	4.4 (2)
GGNN	81.8	2.0 (2)	68.2	2.5 (6)	68.5	1.9 (3)	40.5	1.4 (5)	78.4	2.1 (4)	56.6	1.9 (4)
GAT	82.3	2.3 (2)	65.3	2.1 (4)	69.3	1.6 (2)	42.8	1.6 (4)	77.4	0.5 (6)	63.1	0.7 (4)
GAT-Pn	78.8	0.6 (4)	73.8	1.2 (12)	67.2	0.8 (4)	<b>51.7</b>	<b>1.1 (10)</b>	77.6	1.6 (8)	70.4	1.1 (12)
GAT-Lip	<b>83.1</b>	<b>0.5 (5)</b>	<b>75.3</b>	<b>0.9 (11)</b>	69.1	1.5 (3)	50.9	1.9 (9)	<b>78.9</b>	<b>1.3 (5)</b>	<b>73.3</b>	<b>1.4 (15)</b>

Accuracy of GAT for 100%-PubMed

- \* GAT-Lip: strong across **shallow** and **deep** settings!

# Model depth in real-world datasets

- We evaluate the performance w.r.t. **increasing** number of layers.






Test accuracies of a GAT and a Graph Transformer (GT). By '-Lip' we denote the application of  $\text{Lip}$ , by '-Ln' the  $\text{Ln}$  and by '-Pn' the  $\text{Pn}$ .



Thank you for your **normalized**  
attention!

## References

---

- P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, Y. Bengio, "Graph Attention Networks"  $\ddagger$  , 2018.
- L.Zhao & L.Akoglou, "PairNorm: Tackling Oversmoothing in GNNs"  $\ddagger$  , 2019.
- U. Alon & E. Yahav, "On the Bottleneck of Graph Neural Networks and its Practical Implications,"  $\ddagger$  , 2021.
- W. Hu, M. Fey, M. Zitnik, Y. Dong, H. Ren, B. Liu, M. Catasta & J. Leskovec, "Open Graph Benchmark: Datasets for Machine Learning on Graphs," :  $\emptyset$ , 2020.