

# Bidirectional Self-Normalizing Neural Networks

Thalaiyasingam Ajanthan

Australian National University

August 2020



Australian  
National  
University



# Collaborators



Yao Lu

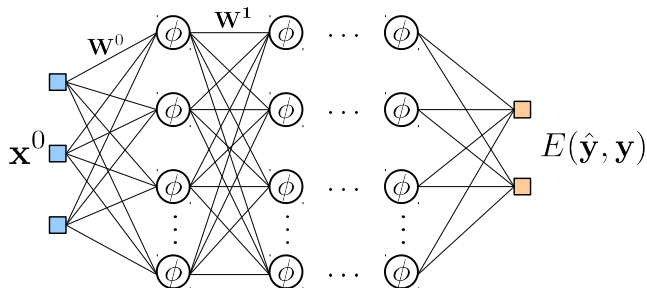


Stephen Gould



Thalaiyasingam Ajanthan

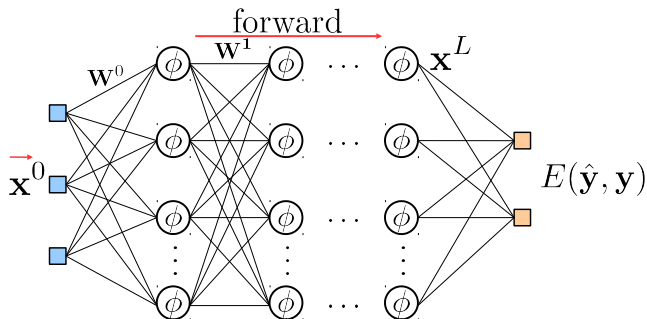
# Signal Propagation in Neural Networks



Forward:  $h^l = \mathbf{W}^l \mathbf{x}^l$ ,  $\mathbf{x}^{l+1} = \sigma(h^l)$

Backward:  $d^l = D^l \left( \prod_{k=l+1}^{L-1} \mathbf{W}^k D^k \right) \frac{\partial E}{\partial \mathbf{x}^l}$ ,  $D^l_{\mathbf{x}} = d'(h^l)$

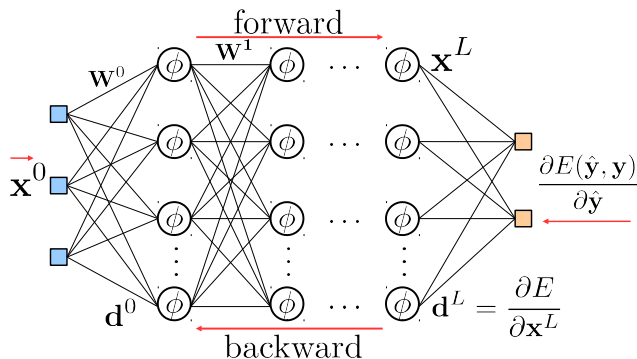
# Signal Propagation in Neural Networks



Forward:  $\mathbf{h}^l = \mathbf{W}^l \mathbf{x}^l$ ,  $\mathbf{x}^{l+1} = \phi(\mathbf{h}^l)$ ,

Backward:  $\mathbf{d}^l = \mathbf{D}^l \left( \prod_{k=l+1}^{L-1} \mathbf{W}^k \mathbf{D}^k \right) \frac{\partial E}{\partial \mathbf{x}^L}$ ,  $D_{ii}^l = \phi'(h_i^l)$ .

# Signal Propagation in Neural Networks



Forward:  $\mathbf{h}^l = \mathbf{W}^l \mathbf{x}^l$ ,  $\mathbf{x}^{l+1} = \phi(\mathbf{h}^l)$ ,

Backward:  $\mathbf{d}^l = \mathbf{D}^l \left( \prod_{k=l+1}^{L-1} \mathbf{W}^k \mathbf{D}^k \right) \frac{\partial E}{\partial \mathbf{x}^L}$ ,  $D_{ii}^l = \phi'(h_i^l)$ .

# Signal Propagation in Neural Networks

## Objective

- ▶ Preserve signal norm in both forward and backward directions.
- ▶ Maintain nonlinear functionality.

## Why?

- ▶ Better trainability: faster convergence and stable training. [Glorot-2010, Klambauer-2017, Pennington-2017]
- ▶ Very deep CNNs and RNNs. [Pennington-2018, Chen-2018]
- ▶ Improved robustness? [Lin-2019]
- ▶ Improved generalization?

## What if signals are not preserved?

- ▶ Signals can saturate  $\Rightarrow$  vanishing/exploding gradients.

# Signal Propagation in Neural Networks

## Objective

- ▶ Preserve signal norm in both forward and backward directions.
- ▶ Maintain nonlinear functionality.

## Why?

- ▶ Better trainability: faster convergence and stable training. [Glorot-2010, Klambauer-2017, Pennington-2017]
- ▶ Very deep CNNs and RNNs. [Pennington-2018, Chen-2018]
- ▶ Improved robustness? [Lin-2019]
- ▶ Improved generalization?

## What if signals are not preserved?

- ▶ Signals can saturate  $\Rightarrow$  vanishing/exploding gradients.

# Signal Propagation in Neural Networks

## Objective

- ▶ Preserve signal norm in both forward and backward directions.
- ▶ Maintain nonlinear functionality.

## Why?

- ▶ Better trainability: faster convergence and stable training. [Glorot-2010, Klambauer-2017, Pennington-2017]
- ▶ Very deep CNNs and RNNs. [Pennington-2018, Chen-2018]
- ▶ Improved robustness? [Lin-2019]
- ▶ Improved generalization?

## What if signals are not preserved?

- ▶ Signals can saturate  $\Rightarrow$  **vanishing/exploding gradients**.



# Existing Approaches to Preserve Signal Propagation

## Initialization

- ▶ Variance scaling initialization. [Glorot-2010, Mishkin-2016]
- ▶ Dynamical isometry and mean-field theory. [Pennington-2017]

## Normalization

- ▶ Self-normalizing neural networks. [Klambauer-2017]
- ▶ Batch normalization and its variants. [Ioffe-2015]
- ▶ {Layer, group, spectral, weight, ...} normalization. [...]

## Architecture

- ▶ Residual connections. [He-2016]

## Drawbacks

- ▶ No rigorous proofs.
- ▶ Do not solve gradient vanishing/explosion. [Philipp-2018]

# Existing Approaches to Preserve Signal Propagation

## Initialization

- ▶ Variance scaling initialization. [Glorot-2010, Mishkin-2016]
- ▶ Dynamical isometry and mean-field theory. [Pennington-2017]

## Normalization

- ▶ Self-normalizing neural networks. [Klambauer-2017]
- ▶ Batch normalization and its variants. [Ioffe-2015]
- ▶ {Layer, group, spectral, weight, ...} normalization. [...]

## Architecture

- ▶ Residual connections. [He-2016]

## Drawbacks

- ▶ No rigorous proofs.
- ▶ Do not solve gradient vanishing/explosion. [Philipp-2018]

# Existing Approaches to Preserve Signal Propagation

## Initialization

- ▶ Variance scaling initialization. [Glorot-2010, Mishkin-2016]
- ▶ Dynamical isometry and mean-field theory. [Pennington-2017]

## Normalization

- ▶ Self-normalizing neural networks. [Klambauer-2017]
- ▶ Batch normalization and its variants. [Ioffe-2015]
- ▶ {Layer, group, spectral, weight, ...} normalization. [...]

## Architecture

- ▶ Residual connections. [He-2016]

### Drawbacks

- ▶ No rigorous proofs.
- ▶ Do not solve gradient vanishing/explosion. [Philipp-2018]

# Existing Approaches to Preserve Signal Propagation

## Initialization

- ▶ Variance scaling initialization. [Glorot-2010, Mishkin-2016]
- ▶ Dynamical isometry and mean-field theory. [Pennington-2017]

## Normalization

- ▶ Self-normalizing neural networks. [Klambauer-2017]
- ▶ Batch normalization and its variants. [Ioffe-2015]
- ▶ {Layer, group, spectral, weight, ...} normalization. [...]

## Architecture

- ▶ Residual connections. [He-2016]

### Drawbacks

- ▶ No rigorous proofs.
- ▶ Do not solve gradient vanishing/explosion. [Philipp-2018]

# Summary of Existing Approaches

	<b>Dynamical Isometry</b>	<b>Self normalization</b>	<b>Batch normalization</b>
Weights ( $\mathbf{W}$ )	Orthogonal	VS-init.	Unconstrained
Activations ( $\phi$ )	Most	SELU	All
Preactivations ( $\mathbf{h}$ )	Linear region	Unconstrained	$\mathbb{E} \approx 0, \mathbb{V} \approx 1$
Forward signal ( $\mathbf{x}$ )	Constrained	$\mathbb{E} \approx 0, \mathbb{V} \approx 1$	Constrained
Backward signal ( $\mathbf{d}$ )	Constrained	Unconstrained	Unconstrained
Functionality	Pseudo-linear	Nonlinear	Nonlinear

[Philipp-2018]

# Summary of Existing Approaches

	<b>Dynamical Isometry</b>	<b>Self normalization</b>	<b>Batch normalization</b>
Weights ( $\mathbf{W}$ )	Orthogonal	VS-init.	Unconstrained
Activations ( $\phi$ )	Most	SELU	All
Preactivations ( $\mathbf{h}$ )	Linear region	Unconstrained	$\mathbb{E} \approx 0, \mathbb{V} \approx 1$
Forward signal ( $\mathbf{x}$ )	Constrained	$\mathbb{E} \approx 0, \mathbb{V} \approx 1$	Constrained
Backward signal ( $\mathbf{d}$ )	Constrained	Unconstrained	Unconstrained
Functionality	Pseudo-linear	Nonlinear	Nonlinear

[Philipp-2018]

## Key missing points

- ▶ Do not preserve signal in both **forward and backward** directions.
- ▶ Do not maintain **nonlinear** functionality.

# Bidirectional Self-Normalizing Neural Networks

Key idea: New class of activation functions: **Gaussian-Poincaré Normalized** activations.

	Dynamical Isometry	Self normalization	BSNN (Ours)
Weights ( $\mathbf{W}$ )	Orthogonal	VS-init.	Orthogonal
Activations ( $\phi$ )	Most	SELU	GPN-versions
Preactivations ( $\mathbf{h}$ )	Linear region	Unconstrained	Unconstrained
Forward signal ( $\mathbf{x}$ )	Constrained	$\mathbb{E} \approx 0, \mathbb{V} \approx 1$	$\mathbb{V} + \mathbb{E}^2 \approx 1$
Backward signal ( $\mathbf{d}$ )	Constrained	Unconstrained	$\mathbb{V} + \mathbb{E}^2 \approx 1$
Functionality	Pseudo-linear	Nonlinear	Nonlinear

# Bidirectional Self-Normalizing Neural Networks

Key idea: New class of activation functions: **Gaussian-Poincaré Normalized** activations.

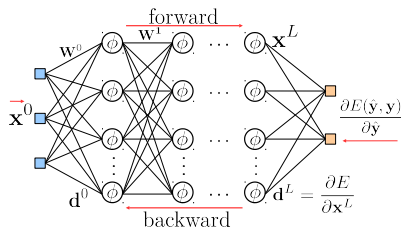
	<b>Dynamical Isometry</b>	<b>Self normalization</b>	<b>BSNN (Ours)</b>
Weights ( $\mathbf{W}$ )	<b>Orthogonal</b>	<b>VS-init.</b>	<b>Orthogonal</b>
Activations ( $\phi$ )	Most	<b>SELU</b>	<b>GPN-versions</b>
Preactivations ( $\mathbf{h}$ )	<b>Linear region</b>	Unconstrained	Unconstrained
Forward signal ( $\mathbf{x}$ )	Constrained	$\mathbb{E} \approx 0, \mathbb{V} \approx 1$	$\mathbb{V} + \mathbb{E}^2 \approx 1$
Backward signal ( $\mathbf{d}$ )	Constrained	Unconstrained	$\mathbb{V} + \mathbb{E}^2 \approx 1$
Functionality	Pseudo-linear	Nonlinear	Nonlinear



# Bidirectional Self-Normalization

**Setting:** Deep fully-connected networks with hidden layers of **same width and no bias**.

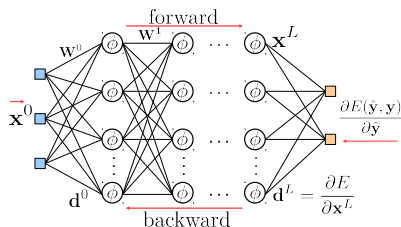
$\mathbf{W}^l \in \mathbb{R}^{n \times n}$ ,  $l = \{1, \dots, L-1\}$   
and  $\phi : \mathbb{R} \rightarrow \mathbb{R}$ .



# Bidirectional Self-Normalization

**Setting:** Deep fully-connected networks with hidden layers of **same width and no bias**.

$\mathbf{W}^l \in \mathbb{R}^{n \times n}$ ,  $l = \{1, \dots, L-1\}$   
and  $\phi : \mathbb{R} \rightarrow \mathbb{R}$ .



$$\text{Forward: } \mathbf{h}^l = \mathbf{W}^l \mathbf{x}^l, \quad \mathbf{x}^{l+1} = \phi(\mathbf{h}^l),$$

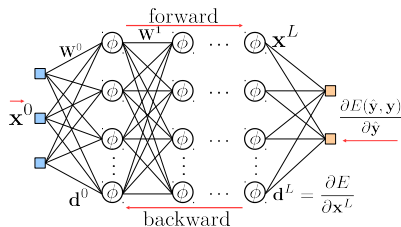
$$\text{Backward: } \mathbf{d}^l = \mathbf{D}^l \left( \prod_{k=l+1}^{L-1} \mathbf{W}^k \mathbf{D}^k \right) \frac{\partial E}{\partial \mathbf{x}^L}, \quad D_{ii}^l = \phi'(h_i^l),$$

$$\text{Gradient: } \frac{\partial E}{\partial \mathbf{W}^l} = \mathbf{d}^l (\mathbf{x}^l)^T.$$

# Bidirectional Self-Normalization

**Setting:** Deep fully-connected networks with hidden layers of **same width and no bias**.

$\mathbf{W}^l \in \mathbb{R}^{n \times n}$ ,  $l = \{1, \dots, L-1\}$   
and  $\phi : \mathbb{R} \rightarrow \mathbb{R}$ .

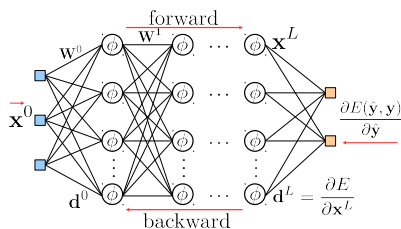


Require:  $\|\mathbf{x}^1\|_2 = \|\mathbf{x}^2\|_2 = \dots = \|\mathbf{x}^L\|_2$ ,  $\{\mathbf{W}^l\}, \phi$  constrained ,  
 $\|\mathbf{d}^1\|_2 = \|\mathbf{d}^2\|_2 = \dots = \|\mathbf{d}^L\|_2$ ,  $\{\mathbf{W}^l\}, \phi'$  constrained .

# Bidirectional Self-Normalization

**Setting:** Deep fully-connected networks with hidden layers of **same width and no bias**.

$\mathbf{W}^l \in \mathbb{R}^{n \times n}$ ,  $l = \{1, \dots, L-1\}$   
and  $\phi : \mathbb{R} \rightarrow \mathbb{R}$ .



Require:  $\|\mathbf{x}^1\|_2 = \|\mathbf{x}^2\|_2 = \dots = \|\mathbf{x}^L\|_2$ ,  $\{\mathbf{W}^l\}, \phi$  constrained ,  
 $\|\mathbf{d}^1\|_2 = \|\mathbf{d}^2\|_2 = \dots = \|\mathbf{d}^L\|_2$ ,  $\{\mathbf{W}^l\}, \phi'$  constrained .

No vanishing/exploding gradients.

# Orthogonal Weight Matrices

$$(\mathbf{W}^l)^T \mathbf{W}^l = \mathbf{W}^l (\mathbf{W}^l)^T = \mathbf{I}_n .$$

## Properties

- ▶ Linear networks: guarantees bidirectional self-normalization. [Saxe-2014]
- ▶ Nonlinear networks: improves trainability with appropriate scaling. [Pennington-2017]
- ▶ Widespread usage in GANs, training sparse networks, quantized networks, *etc.* [Brock-2017, Lee-2020, Lin-2019]

# Orthogonal Weight Matrices

$$(\mathbf{W}^l)^T \mathbf{W}^l = \mathbf{W}^l (\mathbf{W}^l)^T = \mathbf{I}_n .$$

## Properties

- ▶ Linear networks: guarantees bidirectional self-normalization. [Saxe-2014]
- ▶ Nonlinear networks: improves trainability with appropriate scaling. [Pennington-2017]
- ▶ Widespread usage in GANs, training sparse networks, quantized networks, *etc.* [Brock-2017, Lee-2020, Lin-2019]

# GPN Activations

$$\mathbb{E}_{h \sim \mathcal{N}(0,1)} [\phi(h)^2] = \mathbb{E}_{h \sim \mathcal{N}(0,1)} [\phi'(h)^2] = 1 .$$

## Key facts

- ▶ If  $W^l$  is orthogonal,  $h^l$  can be shown to be approximately Gaussian.
- ▶ Function  $\phi$  is GPN and  $\mathbb{E}_{h \sim \mathcal{N}(0,1)}[\phi(h)] = 0$ , if and only if  $\phi$  is linear.
- ▶ A differentiable function  $\phi$  can be transformed into the GPN basis by  $\phi \circ \sigma$ .

# GPN Activations

$$\mathbb{E}_{h \sim \mathcal{N}(0,1)} [\phi(h)^2] = \mathbb{E}_{h \sim \mathcal{N}(0,1)} [\phi'(h)^2] = 1 .$$

## Key facts

- ▶ If  $\mathbf{W}^l$  is orthogonal,  $\mathbf{h}^l$  can be shown to be approximately Gaussian.
- ▶ Function  $\phi$  is GPN and  $\mathbb{E}_{h \sim \mathcal{N}(0,1)}[\phi(h)] = 0$ , if and only if  $\phi$  is **linear**.
- ▶ A differentiable function  $\phi$  can be transformed into its GPN version by  **$a\phi(h) + b$** .



# GPN Activations

$$\mathbb{E}_{h \sim \mathcal{N}(0,1)} [\phi(h)^2] = \mathbb{E}_{h \sim \mathcal{N}(0,1)} [\phi'(h)^2] = 1 .$$

## Key facts

- ▶ If  $\mathbf{W}^l$  is orthogonal,  $\mathbf{h}^l$  can be shown to be approximately Gaussian.
- ▶ Function  $\phi$  is GPN and  $\mathbb{E}_{h \sim \mathcal{N}(0,1)}[\phi(h)] = 0$ , if and only if  $\phi$  is **linear**.
- ▶ A differentiable function  $\phi$  can be transformed into its GPN version by  $a\phi(h) + b$ .

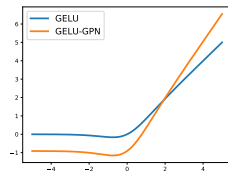
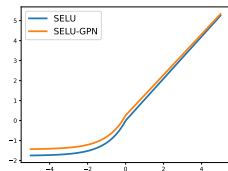
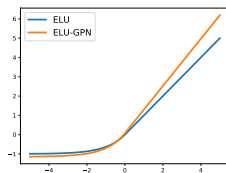
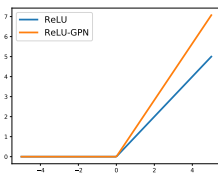
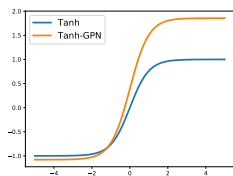
# GPN Activations

$$\mathbb{E}_{h \sim \mathcal{N}(0,1)} [\phi(h)^2] = \mathbb{E}_{h \sim \mathcal{N}(0,1)} [\phi'(h)^2] = 1 .$$

## Key facts

- ▶ If  $\mathbf{W}^l$  is orthogonal,  $\mathbf{h}^l$  can be shown to be approximately Gaussian.
- ▶ Function  $\phi$  is GPN and  $\mathbb{E}_{h \sim \mathcal{N}(0,1)}[\phi(h)] = 0$ , if and only if  $\phi$  is **linear**.
- ▶ A differentiable function  $\phi$  can be transformed into its GPN version by  **$a\phi(h) + b$** .

# GPN Activations



Common activation functions and their GPN versions.

# Thin-Shell Concentration

- ▶ Probabilistic version of the vector norm constraint.

A random vector  $\mathbf{x} \in \mathbb{R}^n$  is **TSC** if for any  $\epsilon > 0$

$$\mathbb{P} \left\{ \left| \frac{1}{n} \|\mathbf{x}\|_2^2 - 1 \right| \geq \epsilon \right\} \rightarrow 0, \quad \text{as } n \rightarrow \infty. \text{ [Bobkov-2003]}$$

Examples: Multivariate Gaussian

Any distribution on  $n$ -unit-sphere scaled by  $\sqrt{n}$ .

Intuitive visualization of Gaussian. [Vershynin-2018]

# Thin-Shell Concentration

- ▶ Probabilistic version of the vector norm constraint.

A random vector  $\mathbf{x} \in \mathbb{R}^n$  is **TSC** if for any  $\epsilon > 0$

$$\mathbb{P} \left\{ \left| \frac{1}{n} \|\mathbf{x}\|_2^2 - 1 \right| \geq \epsilon \right\} \rightarrow 0, \quad \text{as } n \rightarrow \infty. \text{ [Bobkov-2003]}$$

Examples: Multivariate Gaussian

Any distribution on  $n$ -unit-sphere scaled by  $\sqrt{n}$ .

Intuitive visualization of Gaussian. [Vershynin-2018]

# Thin-Shell Concentration

- ▶ Probabilistic version of the vector norm constraint.

A random vector  $\mathbf{x} \in \mathbb{R}^n$  is **TSC** if for any  $\epsilon > 0$

$$\mathbb{P} \left\{ \left| \frac{1}{n} \|\mathbf{x}\|_2^2 - 1 \right| \geq \epsilon \right\} \rightarrow 0, \quad \text{as } n \rightarrow \infty. \text{ [Bobkov-2003]}$$

**Examples:** Multivariate Gaussian

Any distribution on  $n$ -unit-sphere scaled by  $\sqrt{n}$ .

Intuitive visualization of Gaussian. [Vershynin-2018]

# Thin-Shell Concentration

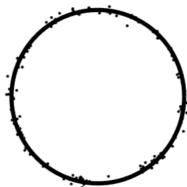
- ▶ Probabilistic version of the vector norm constraint.

A random vector  $\mathbf{x} \in \mathbb{R}^n$  is **TSC** if for any  $\epsilon > 0$

$$\mathbb{P} \left\{ \left| \frac{1}{n} \|\mathbf{x}\|_2^2 - 1 \right| \geq \epsilon \right\} \rightarrow 0, \quad \text{as } n \rightarrow \infty. \text{ [Bobkov-2003]}$$

**Examples:** Multivariate Gaussian

Any distribution on  $n$ -unit-sphere scaled by  $\sqrt{n}$ .



Intuitive visualization of Gaussian. [Vershynin-2018]

# Norm Preservation Theorems

## Assumptions

1. Random vector  $\mathbf{x} \in \mathbb{R}^n$  is **TSC**.
    - ▶ Normalize the input vector.
  2. Random orthogonal weight matrix  $\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n)^T$  is **uniformly distributed**.
  3. Activation function  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  is **GPN**.
  4. Activation function  $\phi$  and its derivative are **Lipschitz continuous**.
- ▶ Most common activation functions satisfy 3 and 4.



# Norm Preservation Theorems

## Assumptions

1. Random vector  $\mathbf{x} \in \mathbb{R}^n$  is **TSC**.
  - ▶ Normalize the input vector.
2. Random orthogonal weight matrix  $\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n)^T$  is **uniformly distributed**.
3. Activation function  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  is **GPN**.
4. Activation function  $\phi$  and its derivative are **Lipschitz continuous**.

▶ Most common activation functions satisfy 3 and 4.

# Norm Preservation Theorems

## Assumptions

1. Random vector  $\mathbf{x} \in \mathbb{R}^n$  is **TSC**.
  - ▶ Normalize the input vector.
2. Random orthogonal weight matrix  $\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n)^T$  is **uniformly distributed**.
3. Activation function  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  is **GPN**.
4. Activation function  $\phi$  and its derivative are **Lipschitz continuous**.

▶ These assumptions are satisfied by ReLU and GPN.

# Norm Preservation Theorems

## Assumptions

1. Random vector  $\mathbf{x} \in \mathbb{R}^n$  is **TSC**.
  - ▶ Normalize the input vector.
2. Random orthogonal weight matrix  $\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n)^T$  is **uniformly distributed**.
3. Activation function  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  is **GPN**.
4. Activation function  $\phi$  and its derivative are **Lipschitz continuous**.

▶ Most common activation functions satisfy 3 and 4.

# Norm Preservation Theorems

## Assumptions

1. Random vector  $\mathbf{x} \in \mathbb{R}^n$  is **TSC**.
  - ▶ Normalize the input vector.
2. Random orthogonal weight matrix  $\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n)^T$  is **uniformly distributed**.
3. Activation function  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  is **GPN**.
4. Activation function  $\phi$  and its derivative are **Lipschitz continuous**.
  - ▶ Most common activation functions satisfy 3 and 4.

# Norm Preservation Theorems

## Assumptions

1. Random vector  $\mathbf{x} \in \mathbb{R}^n$  is **TSC**.
  - ▶ Normalize the input vector.
2. Random orthogonal weight matrix  $\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n)^T$  is **uniformly distributed**.
3. Activation function  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  is **GPN**.
4. Activation function  $\phi$  and its derivative are **Lipschitz continuous**.
  - ▶ Most common activation functions satisfy 3 and 4.

# Norm Preservation Theorems

## Theorem 1: Forward norm preservation

Random vector

$$\left( \phi(\mathbf{w}_1^T \mathbf{x}), \phi(\mathbf{w}_2^T \mathbf{x}), \dots, \phi(\mathbf{w}_n^T \mathbf{x}) \right)^T,$$

is **TSC**.

- ▶ Multiplication by  $\mathbf{W}$  followed by  $\phi$  preserves the norm with high probability.

## Theorem 2: Backward norm preservation

Let  $\mathbf{y} \in \mathbb{R}^n$  with  $\|\mathbf{y}\|_\infty$  and  $D_{ii} = \phi'(\mathbf{w}_i^T \mathbf{x})$ . Then for any  $\epsilon > 0$

$$\mathbb{P} \left\{ \left| \frac{1}{n} \|\mathbf{D}\mathbf{y}\|_2^2 - \|\mathbf{y}\|_2^2 \right| \geq \epsilon \right\} \rightarrow 0,$$

as  $n \rightarrow \infty$ .

- ▶ Multiplication by  $\mathbf{D}$  preserves the norm with high probability.

# Norm Preservation Theorems

## Theorem 1: Forward norm preservation

Random vector

$$\left( \phi(\mathbf{w}_1^T \mathbf{x}), \phi(\mathbf{w}_2^T \mathbf{x}), \dots, \phi(\mathbf{w}_n^T \mathbf{x}) \right)^T,$$

is **TSC**.

- ▶ Multiplication by  $\mathbf{W}$  followed by  $\phi$  **preserves the norm with high probability**.

## Theorem 2: Backward norm preservation

Let  $\mathbf{y} \in \mathbb{R}^n$  with  $\|\mathbf{y}\|_\infty$  and  $D_{ii} = \phi'(\mathbf{w}_i^T \mathbf{x})$ . Then for any  $\epsilon > 0$

$$\mathbb{P} \left\{ \left| \frac{1}{n} \|\mathbf{D}\mathbf{y}\|_2^2 - \|\mathbf{y}\|_2^2 \right| \geq \epsilon \right\} \rightarrow 0,$$

as  $n \rightarrow \infty$ .

- ▶ Multiplication by  $\mathbf{D}$  **preserves the norm with high probability**.

# Norm Preservation Theorems

## Theorem 1: Forward norm preservation

Random vector

$$\left( \phi(\mathbf{w}_1^T \mathbf{x}), \phi(\mathbf{w}_2^T \mathbf{x}), \dots, \phi(\mathbf{w}_n^T \mathbf{x}) \right)^T,$$

is **TSC**.

- ▶ Multiplication by  $\mathbf{W}$  followed by  $\phi$  **preserves the norm with high probability**.

## Theorem 2: Backward norm preservation

Let  $\mathbf{y} \in \mathbb{R}^n$  with  $\|\mathbf{y}\|_\infty$  and  $D_{ii} = \phi'(\mathbf{w}_i^T \mathbf{x})$ . Then for any  $\epsilon > 0$

$$\mathbb{P} \left\{ \left| \frac{1}{n} \|\mathbf{D}\mathbf{y}\|_2^2 - \|\mathbf{y}\|_2^2 \right| \geq \epsilon \right\} \rightarrow 0,$$

as  $n \rightarrow \infty$ .

- ▶ Multiplication by  $\mathbf{D}$  **preserves the norm with high probability**.



# Norm Preservation Theorems

## Theorem 1: Forward norm preservation

Random vector

$$\left( \phi(\mathbf{w}_1^T \mathbf{x}), \phi(\mathbf{w}_2^T \mathbf{x}), \dots, \phi(\mathbf{w}_n^T \mathbf{x}) \right)^T,$$

is **TSC**.

- ▶ Multiplication by  $\mathbf{W}$  followed by  $\phi$  **preserves the norm with high probability**.

## Theorem 2: Backward norm preservation

Let  $\mathbf{y} \in \mathbb{R}^n$  with  $\|\mathbf{y}\|_\infty$  and  $D_{ii} = \phi'(\mathbf{w}_i^T \mathbf{x})$ . Then for any  $\epsilon > 0$

$$\mathbb{P} \left\{ \left| \frac{1}{n} \|\mathbf{D}\mathbf{y}\|_2^2 - \|\mathbf{y}\|_2^2 \right| \geq \epsilon \right\} \rightarrow 0,$$

as  $n \rightarrow \infty$ .

- ▶ Multiplication by  $\mathbf{D}$  **preserves the norm with high probability**.

# Sketch of the Proofs



Key theory: Concentration of measure [Vershynin-2018]

- ▶ Most mass of some high-dimensional probability distributions is concentrated around a certain range.
- ▶ Lipschitz functions do not affect this property.

## Sketch of the Proof: Theorem 1

$$\frac{1}{n} \sum_i \phi(\mathbf{w}_i^T \mathbf{x})^2$$

Rows  $\{\mathbf{w}_i\}$  of a random orthogonal matrix are approximately **independent for large  $n$** .

# Sketch of the Proof: Theorem 1

$$\frac{1}{n} \sum_i \phi(\mathbf{w}_i^T \mathbf{x})^2$$



$$\frac{1}{n} \sum_i \phi(\boldsymbol{\theta}_i^T \mathbf{x})^2$$

Rows  $\{\mathbf{w}_i\}$  of a random orthogonal matrix are approximately **independent for large  $n$** .

$\{\boldsymbol{\theta}_i^T \mathbf{x}\}$  is **approximately Gaussian** when  $\{\boldsymbol{\theta}_i\}$  are independent and  $\mathbf{x}$  is TSC.

# Sketch of the Proof: Theorem 1

$$\frac{1}{n} \sum_i \phi(\mathbf{w}_i^T \mathbf{x})^2$$

$$\frac{1}{n} \sum_i \phi(\boldsymbol{\theta}_i^T \mathbf{x})^2$$

$$\frac{1}{n} \sum_i \phi(z_i)^2$$

Rows  $\{\mathbf{w}_i\}$  of a random orthogonal matrix are approximately **independent for large  $n$** .

$\{\boldsymbol{\theta}_i^T \mathbf{x}\}$  is **approximately Gaussian** when  $\{\boldsymbol{\theta}_i\}$  are independent and  $\mathbf{x}$  is TSC.

Proof follows from Lipschitz continuous and GPN function  $\phi$ .

# Sketch of the Proof: Theorem 1

$$\frac{1}{n} \sum_i \phi(\mathbf{w}_i^T \mathbf{x})^2$$

$$\frac{1}{n} \sum_i \phi(\boldsymbol{\theta}_i^T \mathbf{x})^2$$

$$\frac{1}{n} \sum_i \phi(z_i)^2$$

**1**

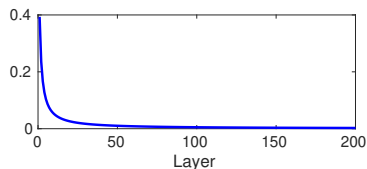
Rows  $\{\mathbf{w}_i\}$  of a random orthogonal matrix are approximately **independent for large  $n$** .

$\{\boldsymbol{\theta}_i^T \mathbf{x}\}$  is **approximately Gaussian** when  $\{\boldsymbol{\theta}_i\}$  are independent and  $\mathbf{x}$  is TSC.

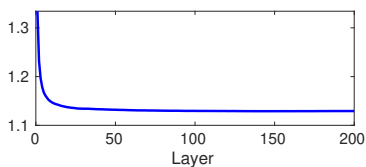
Proof follows from Lipschitz continuous and GPN function  $\phi$ .

Each of these steps are **rigorously proved**.

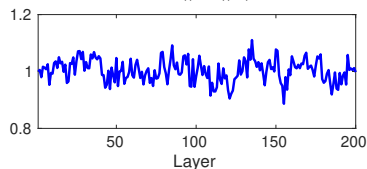
# Synthetic Experiments



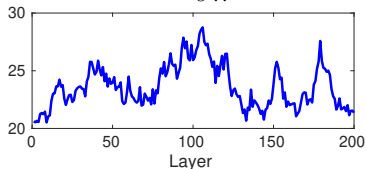
Tanh,  $\|\mathbf{x}^l\|_2^2/n$



Tanh,  $\|\frac{\partial E}{\partial \mathbf{W}^l}\|_F$



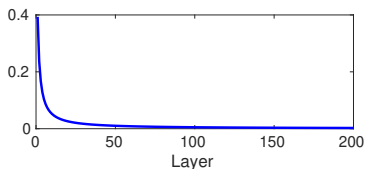
Tanh-GPN,  $\|\mathbf{x}^l\|_2^2/n$



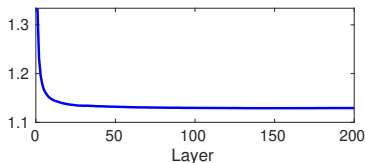
Tanh-GPN,  $\|\frac{\partial E}{\partial \mathbf{W}^l}\|_F$

Forward-backward signal propagation,  $L = 200$ ,  $n = 500$ .

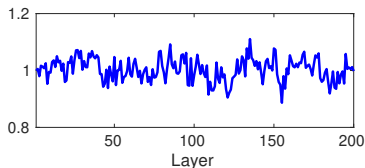
# Synthetic Experiments



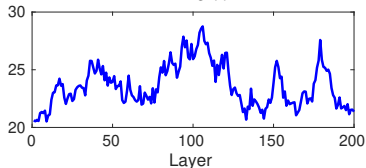
Tanh,  $\|\mathbf{x}^l\|_2^2/n$



Tanh,  $\|\frac{\partial E}{\partial \mathbf{W}^l}\|_F$



Tanh-GPN,  $\|\mathbf{x}^l\|_2^2/n$



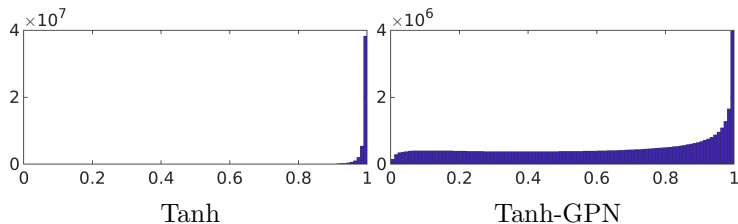
Tanh-GPN,  $\|\frac{\partial E}{\partial \mathbf{W}^l}\|_F$

Forward-backward signal propagation,  $L = 200$ ,  $n = 500$ .

Tanh shows **pseudo-linearity** while Tanh-GPN is nonlinear.



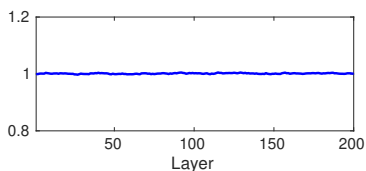
# Synthetic Experiments



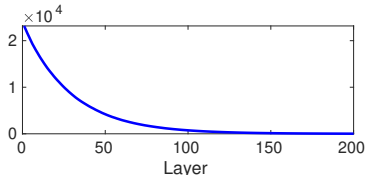
Histogram of  $\phi'(h_i^l)$ , *i.e.*, singular value distribution,  $L = 200$ ,  $n = 500$ .

Tanh shows **pseudo-linearity** while Tanh-GPN is nonlinear.

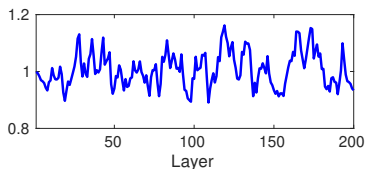
# Synthetic Experiments



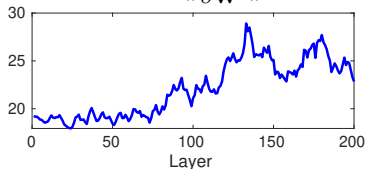
SELU,  $\|\mathbf{x}^l\|_2^2/n$



SELU,  $\|\frac{\partial E}{\partial \mathbf{W}^l}\|_F$



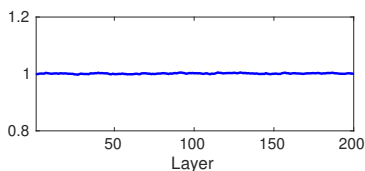
SELU-GPN,  $\|\mathbf{x}^l\|_2^2/n$



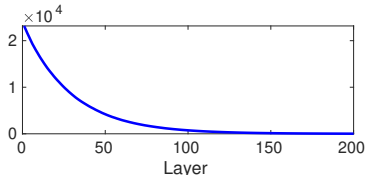
SELU-GPN,  $\|\frac{\partial E}{\partial \mathbf{W}^l}\|_F$

Forward-backward signal propagation,  $L = 200$ ,  $n = 500$ .

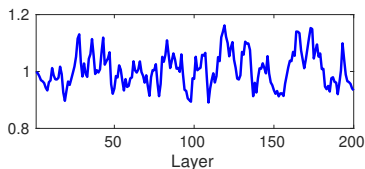
# Synthetic Experiments



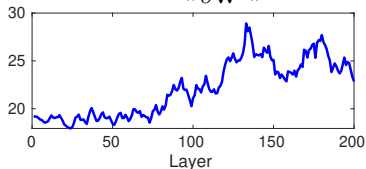
SELU,  $\|\mathbf{x}^l\|_2^2/n$



SELU,  $\|\frac{\partial E}{\partial \mathbf{W}^l}\|_F$



SELU-GPN,  $\|\mathbf{x}^l\|_2^2/n$

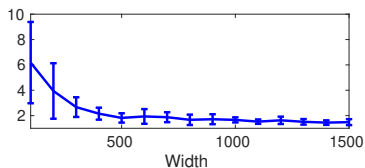


SELU-GPN,  $\|\frac{\partial E}{\partial \mathbf{W}^l}\|_F$

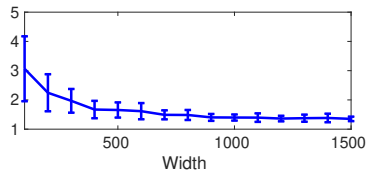
Forward-backward signal propagation,  $L = 200$ ,  $n = 500$ .

SELU suffers **gradient explosion** while SELU-GPN is stable.

# Synthetic Experiments



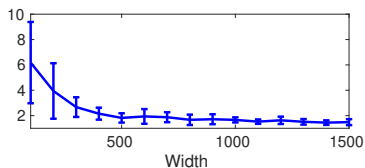
Tanh-GPN



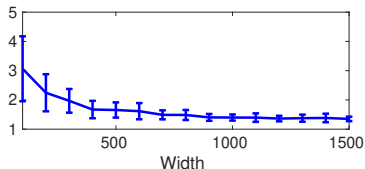
SELU-GPN

Gradient norm ratio, *i.e.*,  $\max_l \left\| \frac{\partial E}{\partial \mathbf{W}^l} \right\|_F / \min_l \left\| \frac{\partial E}{\partial \mathbf{W}^l} \right\|_F$ ,  $L = 200$ .

# Synthetic Experiments



Tanh-GPN



SELU-GPN

Gradient norm ratio, *i.e.*,  $\max_l \left\| \frac{\partial E}{\partial \mathbf{W}^l} \right\|_F / \min_l \left\| \frac{\partial E}{\partial \mathbf{W}^l} \right\|_F$ ,  $L = 200$ .

Larger width leads to more stable gradients.

# Experiments on Real-World Data

	MNIST		CIFAR-10	
	Non-GPN	GPN (Ours)	Non-GPN	GPN (Ours)
Tanh	99.05 (87.39)	<b>99.81</b> (84.93)	80.84 (27.90)	<b>96.39</b> (25.13)
ReLU	11.24 (11.24)	<b>33.28</b> (11.42)	10.00 (10.00)	<b>46.60</b> (10.09)
LReLU	11.24 (11.24)	<b>43.17</b> (11.19)	10.00 (10.21)	<b>51.85</b> (09.89)
ELU	99.06 (98.24)	<b>100.0</b> (97.86)	80.73 (42.39)	<b>99.37</b> (43.35)
SELU	99.86 (97.82)	<b>99.92</b> (97.91)	29.23 (46.47)	<b>98.24</b> (47.74)
GELU	11.24 (12.70)	<b>97.67</b> (11.22)	10.00 (10.43)	<b>90.51</b> (10.00)

Training accuracy with various activation functions,  $L = 200$ ,  $n = 500$ .

# Experiments on Real-World Data

	MNIST		CIFAR-10	
	Non-GPN	GPN (Ours)	Non-GPN	GPN (Ours)
Tanh	99.05 (87.39)	<b>99.81</b> (84.93)	80.84 (27.90)	<b>96.39</b> (25.13)
ReLU	11.24 (11.24)	<b>33.28</b> (11.42)	10.00 (10.00)	<b>46.60</b> (10.09)
LReLU	11.24 (11.24)	<b>43.17</b> (11.19)	10.00 (10.21)	<b>51.85</b> (09.89)
ELU	99.06 (98.24)	<b>100.0</b> (97.86)	80.73 (42.39)	<b>99.37</b> (43.35)
SELU	99.86 (97.82)	<b>99.92</b> (97.91)	29.23 (46.47)	<b>98.24</b> (47.74)
GELU	11.24 (12.70)	<b>97.67</b> (11.22)	10.00 (10.43)	<b>90.51</b> (10.00)

Training accuracy with various activation functions,  $L = 200$ ,  $n = 500$ .

GPN yields **significant improvements** in training accuracy.

# Experiments on Real-World Data

	MNIST		CIFAR-10	
	Non-GPN	GPN (Ours)	Non-GPN	GPN (Ours)
Tanh	<b>96.57</b> (89.32)	95.54 (87.11)	<b>42.71</b> (29.32)	40.95 (26.58)
ReLU	11.35 (11.42)	<b>28.13</b> (11.34)	10.00 (10.00)	<b>34.96</b> (09.96)
LReLU	11.35 (11.63)	<b>49.28</b> (11.66)	10.00 (10.06)	<b>39.38</b> (10.00)
ELU	95.41 ( <b>97.48</b> )	96.56 (96.69)	<b>45.76</b> (44.16)	43.12 (44.36)
SELU	97.33 (97.38)	96.97 ( <b>97.39</b> )	29.55 (45.88)	<b>45.90</b> (45.52)
GELU	11.35 (10.28)	<b>95.82</b> (09.74)	10.00 (10.00)	<b>36.94</b> (10.00)

Test accuracy with various activation functions,  $L = 200$ ,  $n = 500$ .



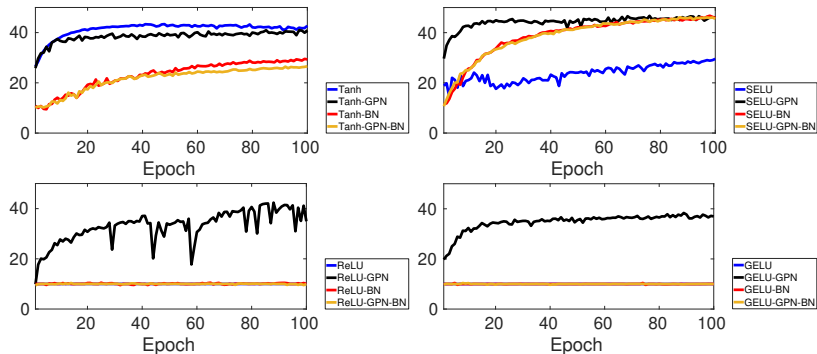
# Experiments on Real-World Data

	MNIST		CIFAR-10	
	Non-GPN	GPN (Ours)	Non-GPN	GPN (Ours)
Tanh	<b>96.57</b> (89.32)	95.54 (87.11)	<b>42.71</b> (29.32)	40.95 (26.58)
ReLU	11.35 (11.42)	<b>28.13</b> (11.34)	10.00 (10.00)	<b>34.96</b> (09.96)
LReLU	11.35 (11.63)	<b>49.28</b> (11.66)	10.00 (10.06)	<b>39.38</b> (10.00)
ELU	95.41 ( <b>97.48</b> )	96.56 (96.69)	<b>45.76</b> (44.16)	43.12 (44.36)
SELU	97.33 (97.38)	96.97 ( <b>97.39</b> )	29.55 (45.88)	<b>45.90</b> (45.52)
GELU	11.35 (10.28)	<b>95.82</b> (09.74)	10.00 (10.00)	<b>36.94</b> (10.00)

Test accuracy with various activation functions,  $L = 200$ ,  $n = 500$ .

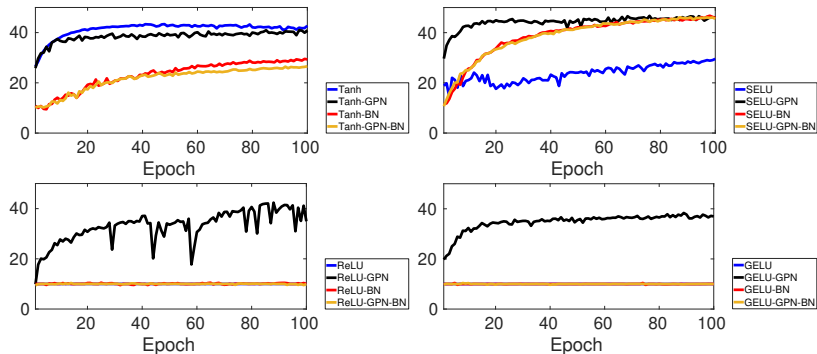
GPN yields **improvements** in many cases.

# Experiments on Real-World Data



Test accuracy curves on CIFAR-10,  $L = 200$ ,  $n = 500$ .

# Experiments on Real-World Data



Test accuracy curves on CIFAR-10,  $L = 200$ ,  $n = 500$ .

GPN **accelerates training** in many cases.

# Summary

- ▶ We introduced BSNN which constrains signal norm in both directions in nonlinear networks via orthogonal weights and GPN activation functions.
- ▶ Many common activations functions can be transformed into their respective GPN versions.
- ▶ Rigorously proved that gradient vanishing/exploding problem disappears with high probability if the width is sufficiently large.

# Current Limitations

- ▶ Theoretical analysis is limited to same width, fully-connected networks.
- ▶ Generalization capabilities are unclear.
- ▶ Universality of BSNN is an open question.

Questions?

Thank you!