



复杂关键软件环境全国重点实验室
State Key Laboratory of Complex & Critical Software Environment
北京航空航天大学人工智能研究院
Institute of Artificial Intelligence, Beihang University



ICML
International Conference
On Machine Learning

On the Nonlinearity of Layer Normalization

Yunhao Ni, Yuxin Guo, Junlong Jia, Lei Huang*



Beihang University
Beijing, China.



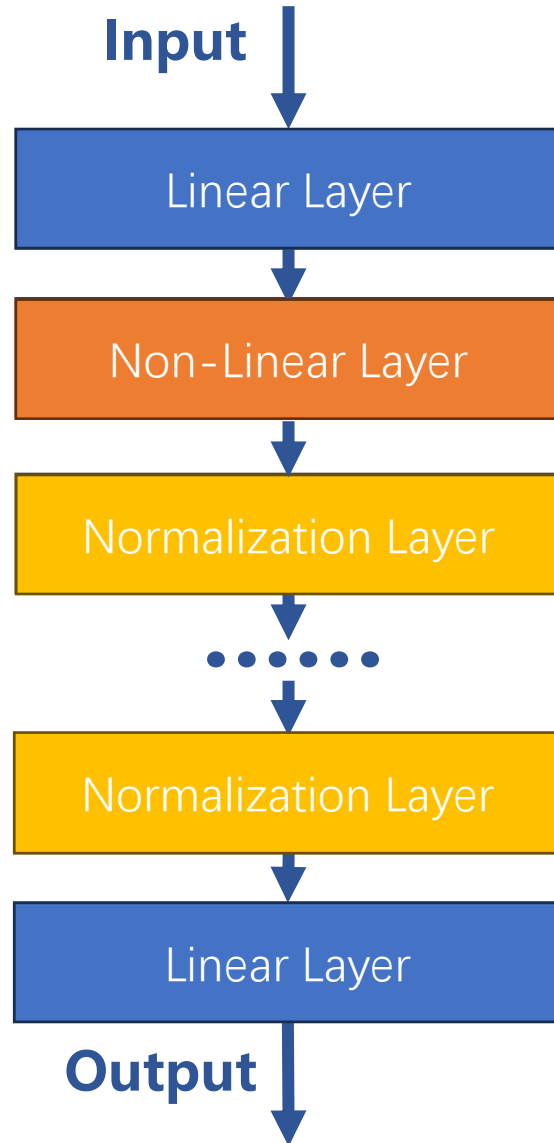
复杂关键软件环境全国重点实验室
State Key Laboratory of Complex & Critical Software Environment



北京航空航天大学人工智能研究院
Institute of Artificial Intelligence, Beihang University

- Background
- The Existence of Nonlinearity in LN
- Capacity of a Network with LN
- Amplify and Exploit the Nonlinearity of LN
- Conclusion

- **Background**
- The Existence of Nonlinearity in LN
- Capacity of a Network with LN
- Amplify and Exploit the Nonlinearity of LN
- Conclusion



→ store the main parameters of the neural network

→ contains the main expressive power

→

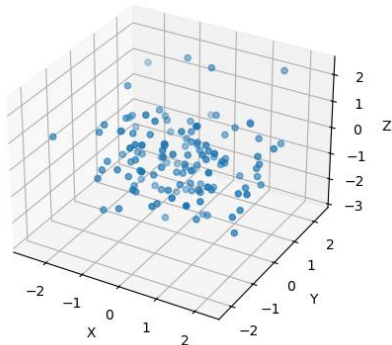
Previous:

→ stablizing training and accelerating optimization

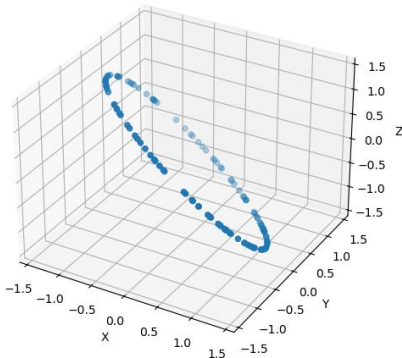
$$\text{BN: } \hat{x} = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad \rightarrow \text{Optimization} > \text{Fitting}$$

But?

Layer Normalization



$$\hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

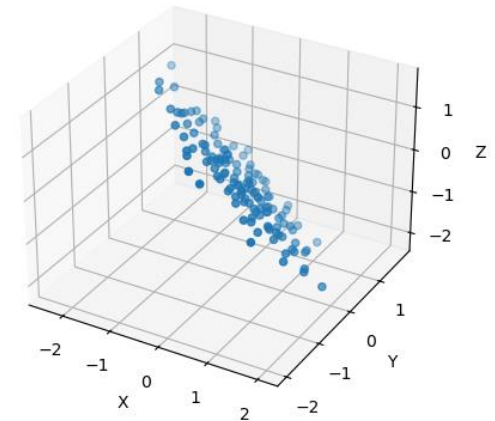


Centering

$$\bar{\mathbf{x}} = \mathbf{x} - \frac{1}{d} (\mathbf{1}_d^\top \mathbf{x}) \cdot \mathbf{1}_d$$

Projected onto the **Hyperplane**

$$\{\mathbf{x} \in \mathbb{R}^d : x_1 + \dots + x_d = 0\}$$

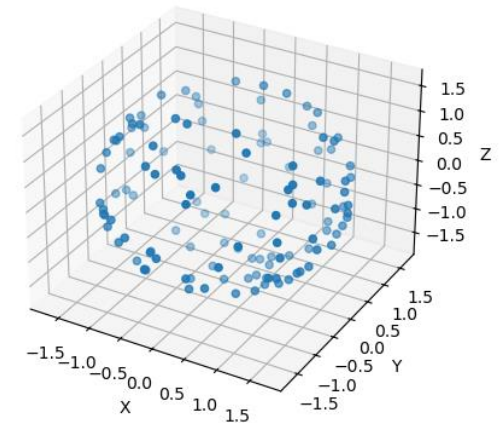


Scaling

$$\hat{\mathbf{x}} = \sqrt{d} \frac{\bar{\mathbf{x}}}{\|\bar{\mathbf{x}}\|_2}$$

Projected onto the **Hypersphere**

$$\{\mathbf{x} \in \mathbb{R}^d : [x_1]^2 + \dots + [x_d]^2 = d\}$$



- Background
- **The Existence of Nonlinearity in LN**
- Capacity of a Network with LN
- Amplify and Exploit the Nonlinearity of LN
- Conclusion

The Existence of Nonlinearity in LN

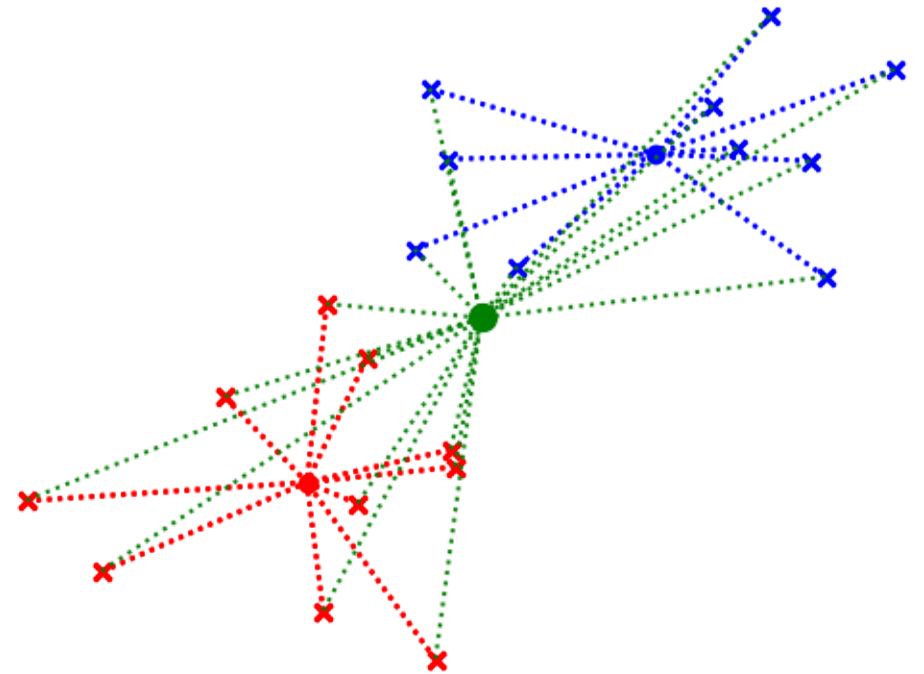
Find a linearly invariant index.

$$SS(\mathbf{X}_c) = \sum_{i=1}^m \|\mathbf{x}_{ci} - \bar{\mathbf{x}}_c\|^2$$

$$SSR(\mathbf{X}_1, \mathbf{X}_2) = \frac{\text{Intra distance } SS(\mathbf{X}_1) + SS(\mathbf{X}_2)}{\text{Total distance } SS([\mathbf{X}_1, \mathbf{X}_2])}$$

$$LSSR(\mathbf{X}_1, \mathbf{X}_2) = \inf_{\varphi \in \mathbb{D}_\varphi(d)} SSR(\varphi(\mathbf{X}_1), \varphi(\mathbf{X}_2))$$

Linear Transformations

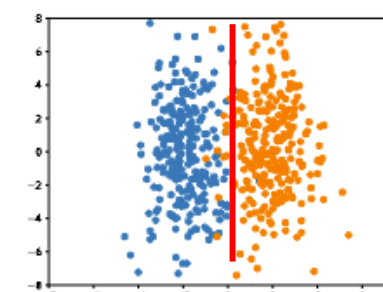
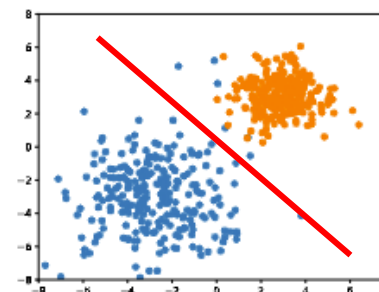
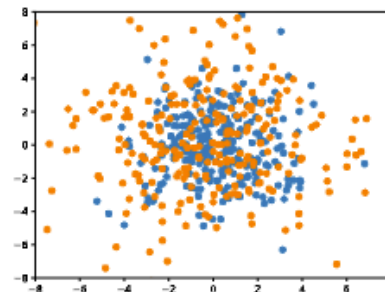
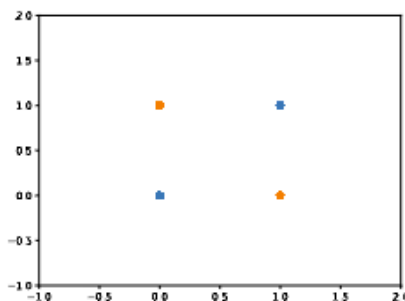


The Existence of Nonlinearity in LN

SSR and LSSR

Data

*XOR data



Linearly separable?

No

No

Yes

Yes

SSR

0.9963

0.9929

0.2304

0.7365

LSSR

0.9929

0.9859

0.1312

0.2157

LSSR is a better index to describe **linear separability** than SSR.

The Existence of Nonlinearity in LN

Proof Method

Taylor's Expansion

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + o(x - x_0)$$

A brief equation

$$\underbrace{SSR(v^T \hat{X}_1, v^T \hat{X}_2)}_{\text{Linear+Scaling}} = \underbrace{LSSR(X_1, X_2)}_{\text{Lower Bound}} - \frac{2(T_1 + T_2)}{T_3} t + o(t)$$

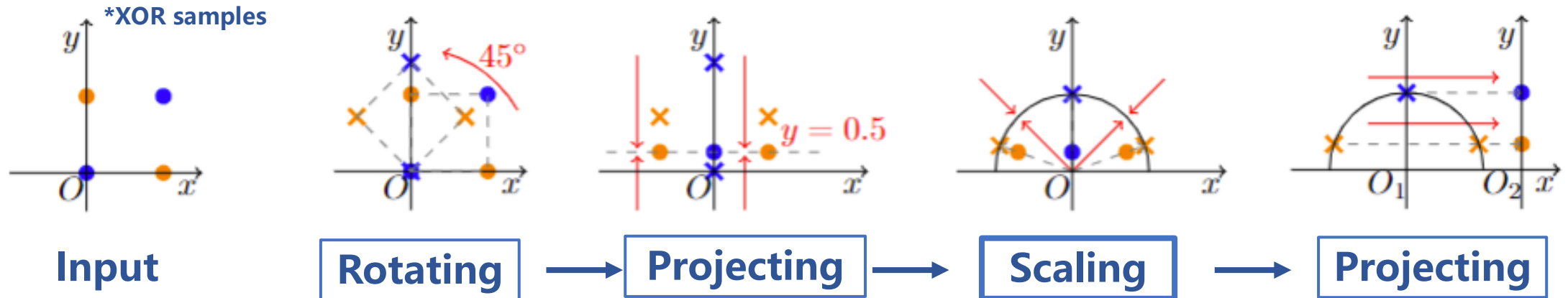
< 0

Linear transformations with LN can break LSSR.

- Background
- The Existence of Nonlinearity in LN
- **Capacity of a Network with LN**
- Amplify and Exploit the Nonlinearity of LN
- Conclusion

Capacity of a Network with LN

XOR



Classify XOR samples with **linear transformations** and **scaling** only.

※ **Hint:** **Scaling** can be represented by **LN** and **linear transformations** only.

Capacity of a Network with LN

Binary Classification

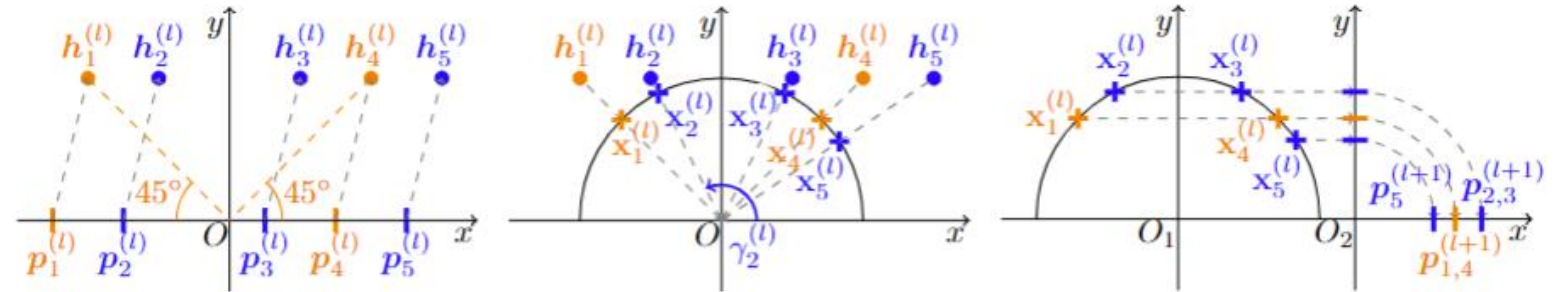
Algorithm 1 Projection Merge Algorithm

input The initial input $P^{(1)}$.

output The final output $P^{(L)}$.

```

1:  $l \leftarrow 1$ ;
2:  $\mathbb{P} \leftarrow \{p_1^{(l)}, p_2^{(l)}, \dots, p_m^{(l)}\}$ ;
3: while  $\mathbb{P} \neq \emptyset$  do
4:    $i \leftarrow \arg \min_k \{p_k^{(l)} : p_k^{(l)} \in \mathbb{P}\}$ ;
5:    $\mathbb{J}_i \leftarrow \{p_j^{(l)} \in \mathbb{P} : p_j^{(l)} \neq p_i^{(l)}, y_j = y_i\}$ ;
6:   if  $\mathbb{J}_i \neq \emptyset$  then
7:      $j \leftarrow \arg \min_k \{p_k^{(l)} : p_k^{(l)} \in \mathbb{J}_i\}$ ;
8:     for  $k \leftarrow 1$  to  $m$  do
9:        $h_k^{(l)} \leftarrow p_k^{(l)} - \begin{bmatrix} p_i^{(l)} + p_j^{(l)} \\ p_i^{(l)} - p_j^{(l)} \end{bmatrix} / 2$ ;
10:       $x_k^{(l)} \leftarrow h_k^{(l)} / \|h_k^{(l)}\|$ ;
11:       $p_k^{(l+1)} \leftarrow \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} x_k^{(l)}$ ;
12:     end for
13:      $l \leftarrow l + 1$ ;
14:      $\mathbb{P} \leftarrow \{p_1^{(l)}, p_2^{(l)}, \dots, p_m^{(l)}\}$ ;
15:   else
16:     remove  $p_j^{(l)}$  from  $\mathbb{P}$ , as long as  $p_j^{(l)} = p_i^{(l)}$ ;
17:   end if
18: end while
19: return  $P^{(l)}$ ;
    
```



Projecting



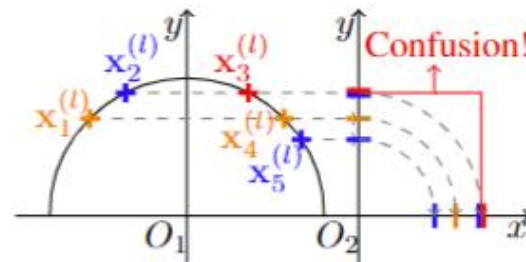
Scaling



Projecting

Binarily classify any m samples with **linear transformations** and **scaling** only.

Multiclass Classification



Difference: **Confusion** is possible.

Solution: Breaking Parallelization.

Capacity of a Network with LN

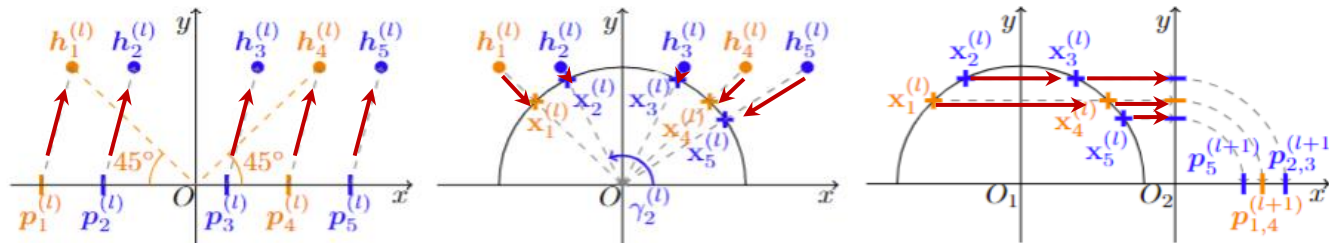
- ▶ **The first** to discuss about the expressive power of layer normalization.

Universal Approximation Theory of LN-Net

An infinitely deep LN-Net can classify **any given m samples** correctly.

Computation Operation \Rightarrow **Merging Operation**

Universal Approximation \Rightarrow **Universal Classification**

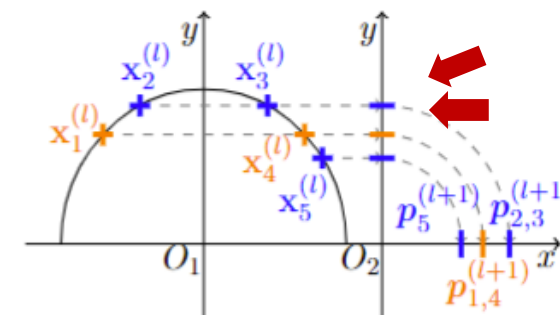


Projecting

Scaling

Projecting

VC dimension



Samples with the same label: **at least two** are merged with each LN.
Samples with different labels: **none** are merged.
 Width: at least **3**.

- ▶ Given an LN-Net $f_{\theta}(\cdot)$ with **width 3** and **depth L** its VC dimension $VCdim(f_{\theta}(\cdot))$ is **lower bounded by $L + 2$** .

- Background
- The Existence of Nonlinearity in LN
- Capacity of a Network with LN
- **Amplify and Exploit the Nonlinearity of LN**
- Conclusion

Group based LN (LN-G) has stronger nonlinearity than LN

Measurement of Nonlinearity

$$\text{Hessian: } \mathcal{H}(f; x) = \sum_{i=1}^d \left\| \frac{\partial^2 y_i}{\partial x^2} \right\|_F^2$$

Noted that $\mathcal{H}(f; x) \geq 0$, and $\mathcal{H}(f; x) = 0$ if and only if f is linear.

→ We *assume* that the larger $\mathcal{H}(f; x)$, the more nonlinearity f contains.

Amplifying Nonlinearity by Group

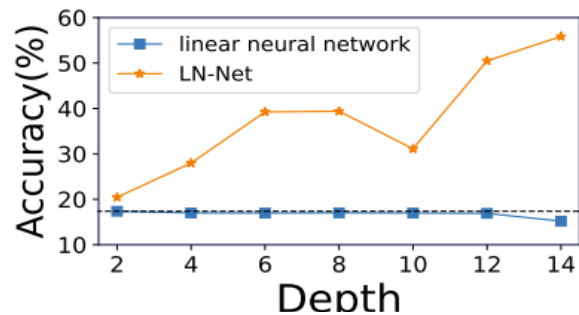
Proposition: Given $g \leq d/3$, we have $\frac{\mathcal{H}(\psi_G(g; \cdot); x)}{\mathcal{H}(\psi_L(\cdot); x)} \geq 1$. When $g = d/4$, $\frac{\mathcal{H}(\psi_G(g; \cdot); x)}{\mathcal{H}(\psi_L(\cdot); x)} \geq \frac{d}{8}$.

※ $\psi_G(g; \cdot)$ denotes LN-G on \mathbb{R}^d with group number g , $\psi_L(\cdot)$ denotes LN on \mathbb{R}^d .

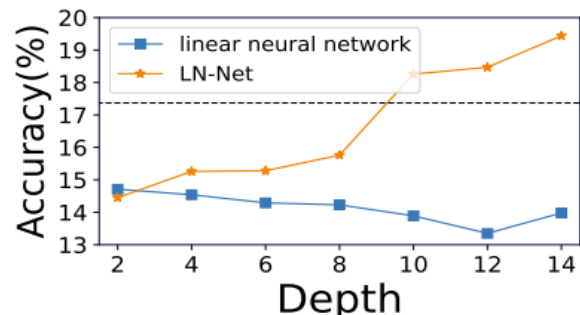
→ LN-G can amplify the nonlinearity of LN by using appropriated group number.

Comparison of Representation Capacity by Fitting Random Labels

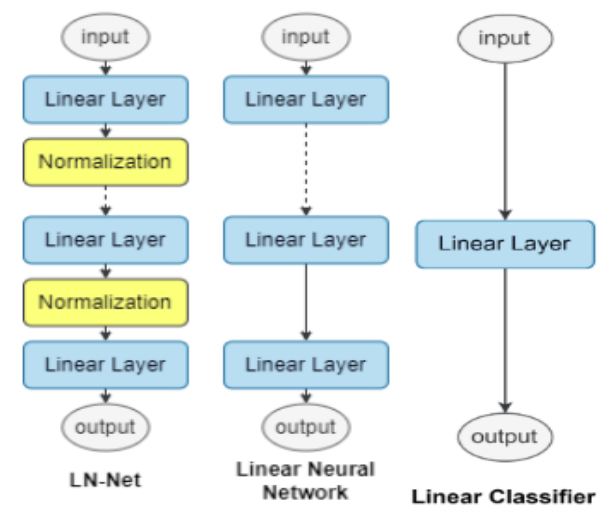
LN-Net & linear neural network & linear classifier



(a) CIFAR-10-RL.

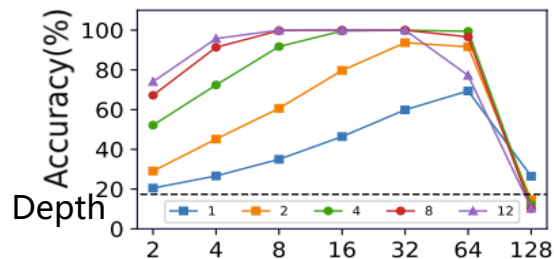


(b) MNIST-RL.

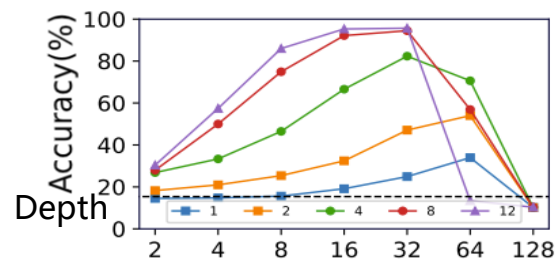


→ LN can break the bound of linearity.

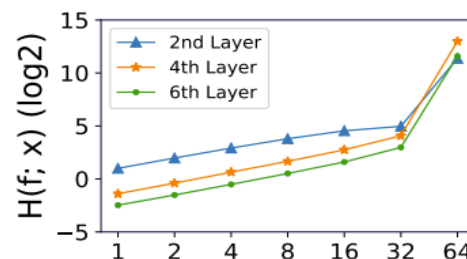
LN-Net (LN-G)



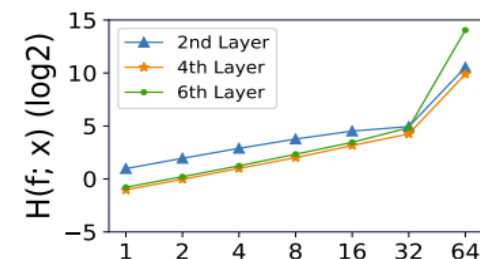
(a) Accuracy on CIFAR-10-RL.



(b) Accuracy on MNIST-RL.



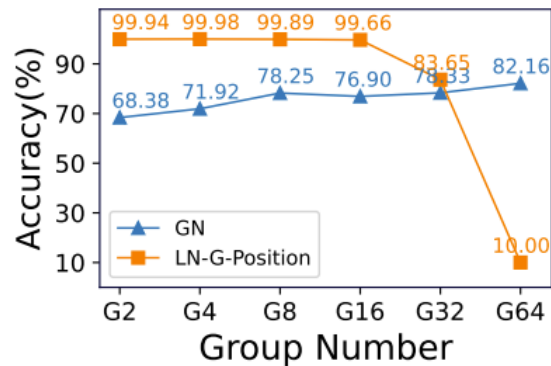
(c) $\mathcal{H}(f; \mathbf{x})$ on CIFAR-10-RL.



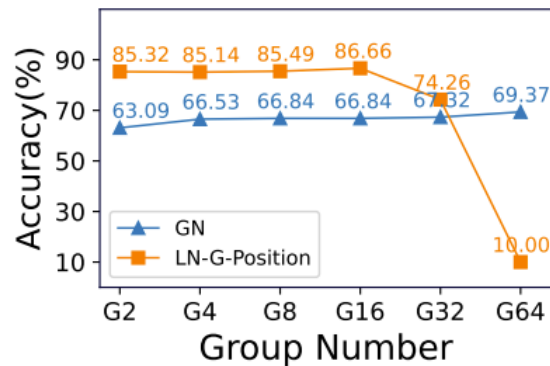
(d) $\mathcal{H}(f; \mathbf{x})$ on MNIST-RL.

Inspiration for Neural Architecture Design

ResNet without ReLU on CIFAR-10



(a) Training.



(b) Test.

Normalization methods	Train Acc(%)	Test Acc(%)
IN	10	10
BN	36.0	39.3
LN	59.5	62.85
GN	82.16	69.37
LN-G-Position	99.66	86.66

Transformer

fairseq-py on IWSLT14 De-EN:(BLEU)

LN: 35.01 ± 0.10 ; LN-G: 35.23 ± 0.07

Tiny-ViT on CIFAR-10: (test Acc)

LN: 88.81% ; LN-G: 89.26%

- Background
- The Existence of Nonlinearity in LN
- Capacity of a Network with LN
- Amplify and Exploit the Nonlinearity of LN
- **Conclusion**

Conclusion

- Mathematically demonstrated that LN is a nonlinear transformation.
- Theoretically showed the representation capacity of an LN-Net in correctly classifying samples with any label assignment.
- Call for reconsidering the analyses of the representation capacity of a network with normalization layer.

Thanks for your attention!