

2021-1 딥러닝기술 및 응용 - Paper Review

# Weight Initialization

[Hinton2006] A fast learning algorithm for deep belief net.

[Glorot2010] Understanding the difficulty of training deep feedforward neural networks.

[He2015] Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification.



KISTI-UST  
Donghun Yang

2021.04.16.FRI.

## Contents

**01 Why do I need to initialize the weights properly?**

**02 Deep Belief Network**

**03 Initialize all weights to Zero**

**04 Initialize weights randomly**

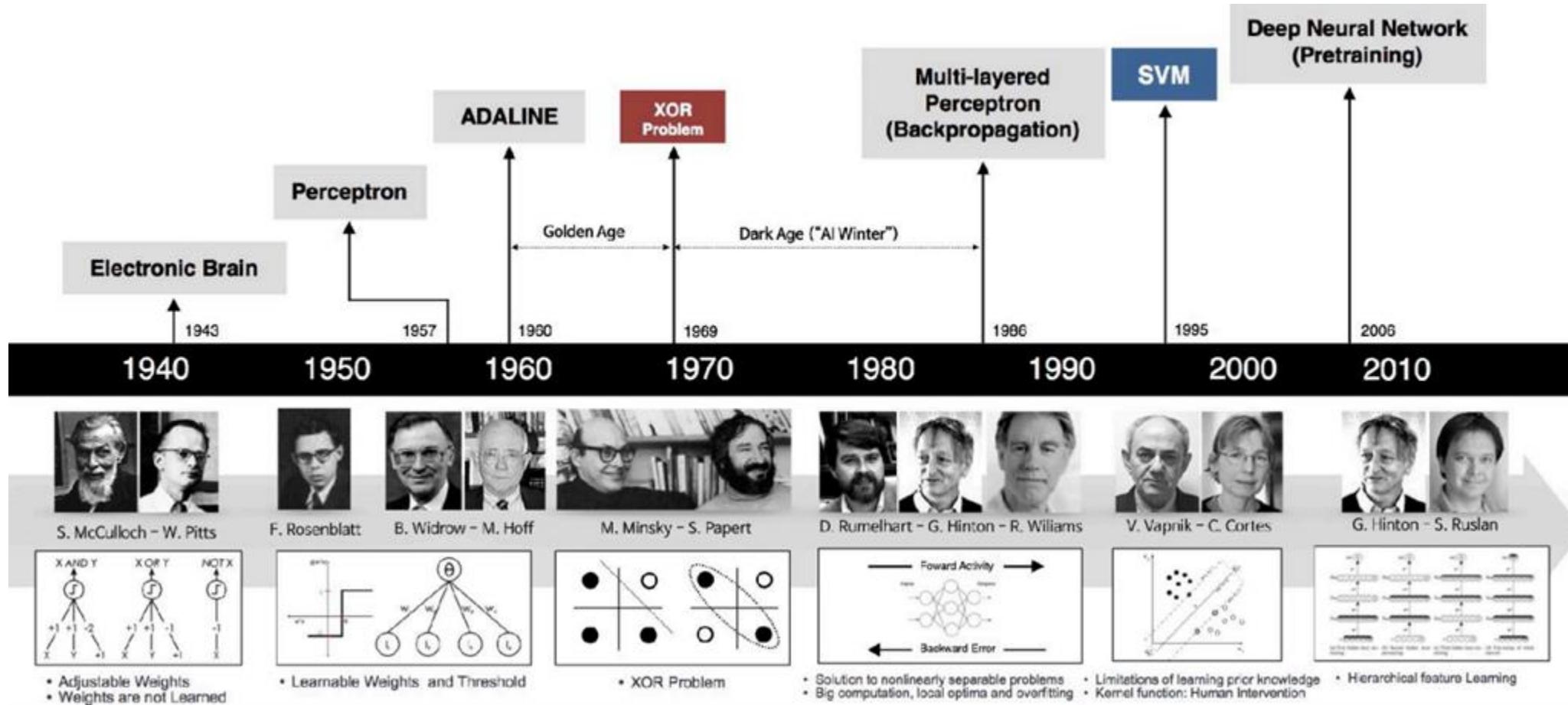
**05 Xavier initialization for linear activation function**

**06 He initialization for ReLU activation function**

**07 Conclusion**

# 01 Why do we need to initialize the weights properly?

## History of AI



## 01 Why do we need to initialize the weights properly?

Geoffrey Hinton's summary of findings up to today

### Geoffrey Hinton's summary of findings up to today

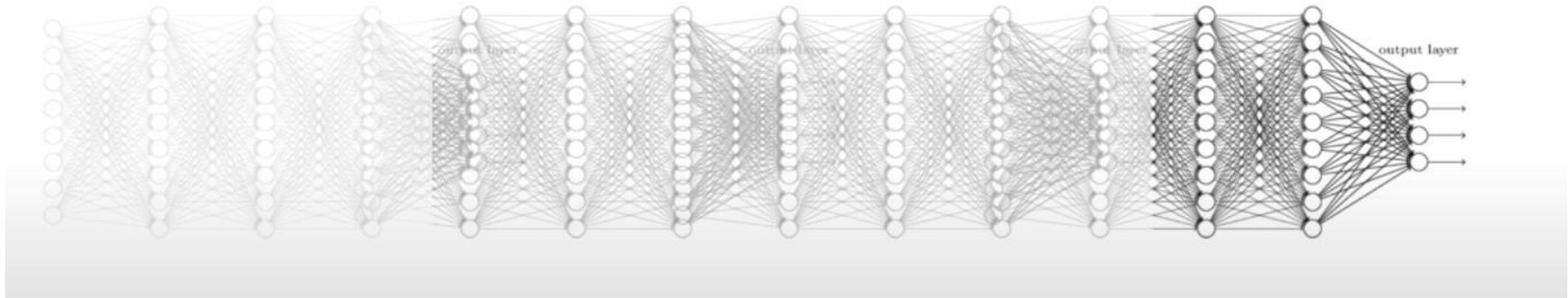
- Our labeled datasets were thousands of times too small.
- Our computers were millions of times too slow.
- We initialized the weights in a stupid way.
- We used the wrong type of non-linearity.

<https://www.skynettoday.com/overviews/neural-net-history>

# 01 Why do we need to initialize the weights properly?

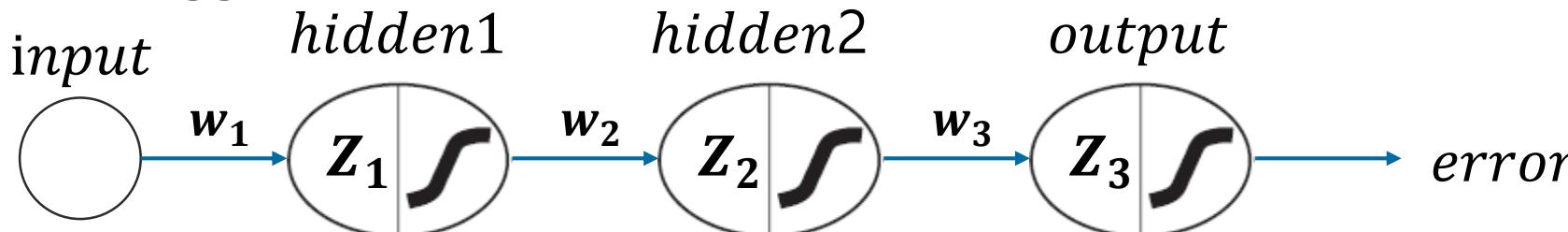
Vanishing gradient

## Vanishing gradient (NN winter2: 1986-2006)



# 01 Why do we need to initialize the weights properly?

Vanishing gradient



$$\frac{\partial \text{error}}{\partial w_3} = \frac{\partial \text{error}}{\partial \text{output}} * \frac{\partial \text{output}}{\partial w_3} = \frac{\partial \text{error}}{\partial \text{output}} * \frac{\partial \text{output}}{\partial z_3} * \frac{\partial z_3}{\partial w_3} = \frac{\partial \text{error}}{\partial \text{Sigmoid}(z_3)} * \frac{\partial \text{Sigmoid}(z_3)}{\partial z_3} * \text{hidden2}$$

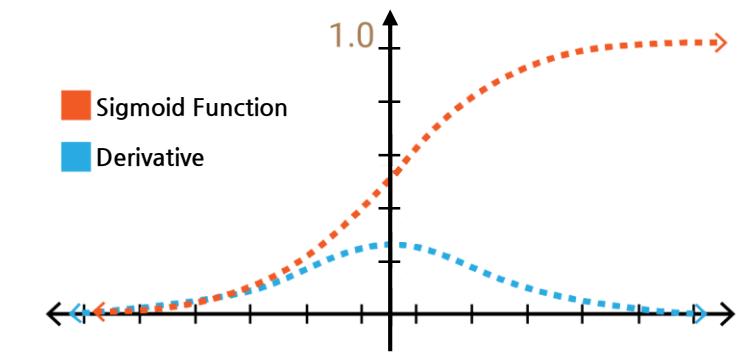
<0.4

$$\begin{aligned} \frac{\partial \text{error}}{\partial w_2} &= \frac{\partial \text{error}}{\partial \text{output}} * \frac{\partial \text{output}}{\partial \text{hidden2}} * \frac{\partial \text{hidden2}}{\partial w_2} = \frac{\partial \text{error}}{\partial \text{output}} * \frac{\partial \text{output}}{\partial z_3} * \frac{\partial z_3}{\partial \text{hidden2}} * \frac{\partial \text{hidden2}}{\partial z_2} * \frac{\partial z_2}{\partial w_2} \\ &= \frac{\partial \text{error}}{\partial \text{Sigmoid}(z_3)} * \frac{\partial \text{Sigmoid}(z_3)}{\partial z_3} * \frac{\partial \text{Sigmoid}(z_2)}{\partial z_3} * \frac{\partial \text{Sigmoid}(z_2)}{\partial z_2} * \text{hidden1} \\ &= \frac{\partial \text{error}}{\partial \text{Sigmoid}(z_3)} * \frac{\partial \text{Sigmoid}(z_3)}{\partial z_3} * w_3 * \frac{\partial \text{Sigmoid}(z_2)}{\partial z_2} * \text{hidden1} \end{aligned}$$

<0.4 \* <1 \* <0.4

$$\begin{aligned} \frac{\partial \text{error}}{\partial w_1} &= \frac{\partial \text{error}}{\partial \text{output}} * \frac{\partial \text{output}}{\partial \text{hidden2}} * \frac{\partial \text{hidden2}}{\partial \text{hidden1}} * \frac{\partial \text{hidden1}}{\partial w_1} = \frac{\partial \text{error}}{\partial \text{output}} * \frac{\partial \text{output}}{\partial z_3} * \frac{\partial z_3}{\partial \text{hidden2}} * \frac{\partial \text{hidden2}}{\partial z_2} * \frac{\partial z_2}{\partial \text{hidden1}} * \frac{\partial \text{hidden1}}{\partial z_1} * \frac{\partial z_1}{\partial w_1} \\ &= \frac{\partial \text{error}}{\partial \text{Sigmoid}(z_3)} * \frac{\partial \text{Sigmoid}(z_3)}{\partial z_3} * \frac{\partial \text{Sigmoid}(z_2)}{\partial z_3} * \frac{\partial \text{Sigmoid}(z_2)}{\partial z_2} * \frac{\partial \text{Sigmoid}(z_1)}{\partial z_2} * \frac{\partial \text{Sigmoid}(z_1)}{\partial z_1} * \text{input} \\ &= \frac{\partial \text{error}}{\partial \text{Sigmoid}(z_3)} * \frac{\partial \text{Sigmoid}(z_3)}{\partial z_3} * w_3 * \frac{\partial \text{Sigmoid}(z_2)}{\partial z_2} * w_2 * \frac{\partial \text{Sigmoid}(z_1)}{\partial z_2} * \text{input} \end{aligned}$$

<0.4 \* <1 \* <0.4 \* <1 \* <0.4



$$\begin{aligned} z_3 &= w_3 * \text{hidden2} \\ \text{output} &= \text{Sigmoid}(z_3) \end{aligned}$$

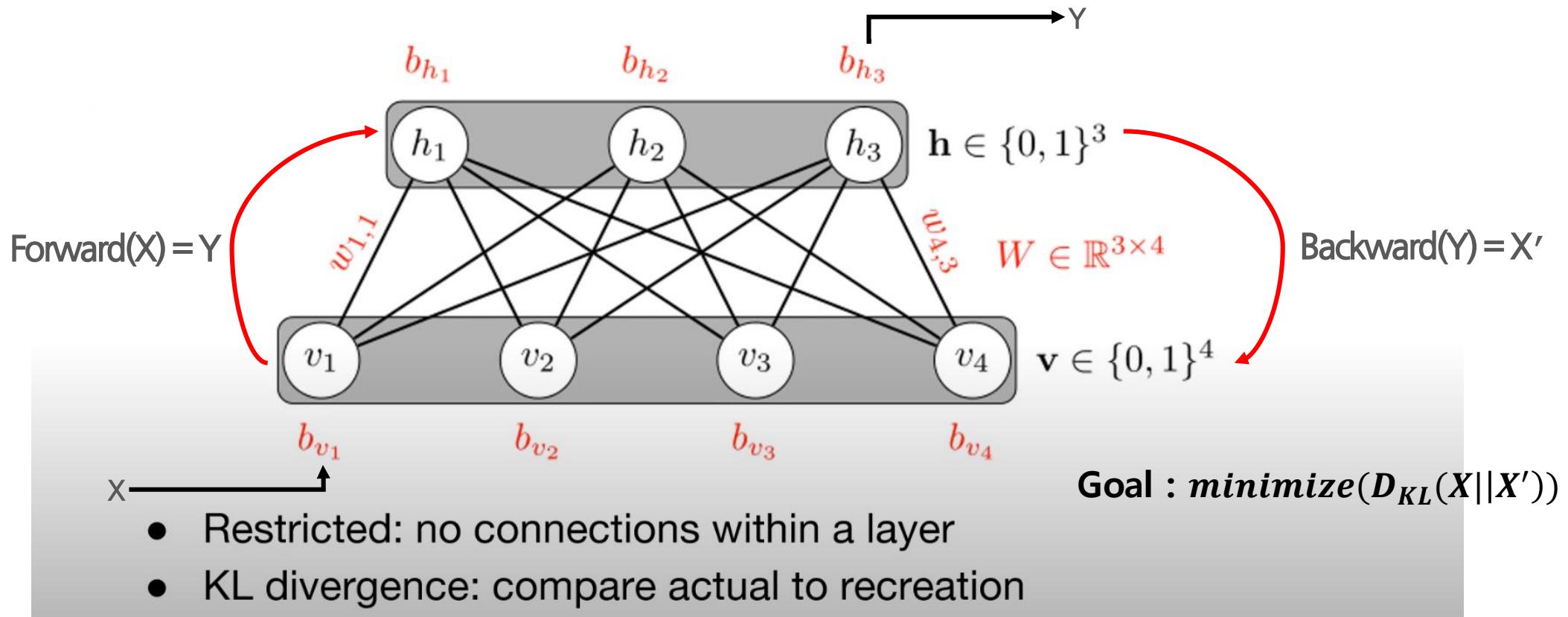
$$\begin{aligned} z_2 &= w_2 * \text{hidden1} \\ \text{hidden2} &= \text{Sigmoid}(z_2) \end{aligned}$$

$$\begin{aligned} z_1 &= w_1 * \text{input} \\ \text{hidden1} &= \text{Sigmoid}(z_1) \end{aligned}$$

## 02 Deep Belief Network

[Hinton2006] A fast learning algorithm for deep belief net - Restricted Boltzmann Machine (RBM)

### Restricted Boltzmann Machine

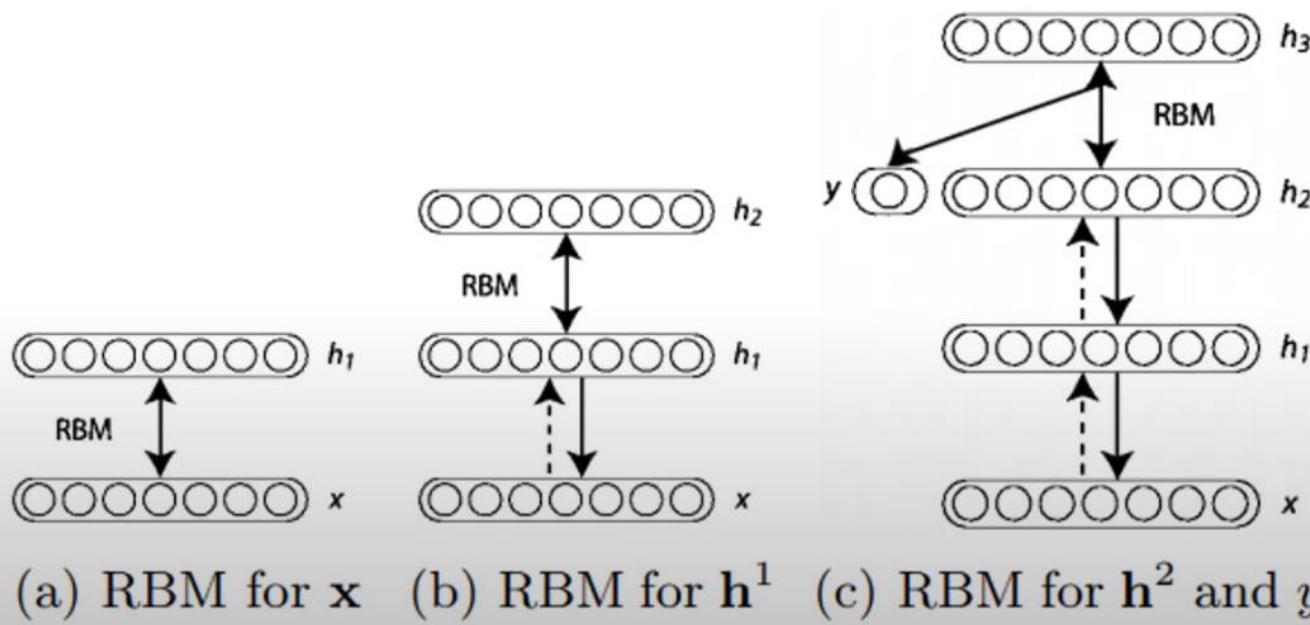


## 02 Deep Belief Network

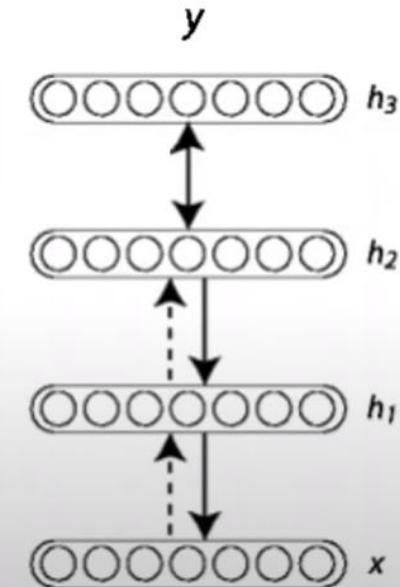
[Hinton2006] A fast learning algorithm for deep belief net. - Restricted Boltzmann Machine (RBM)

# Deep Belief Network

### Pre-training



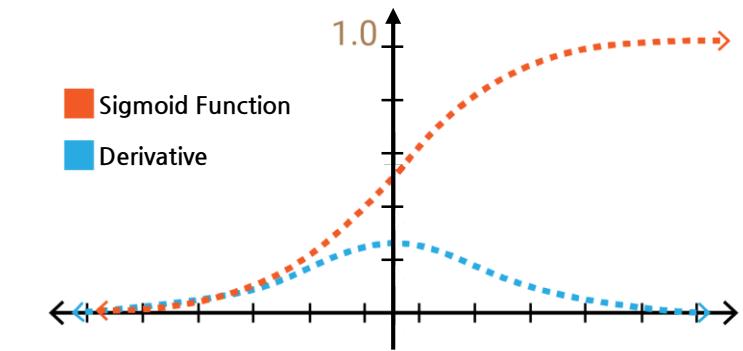
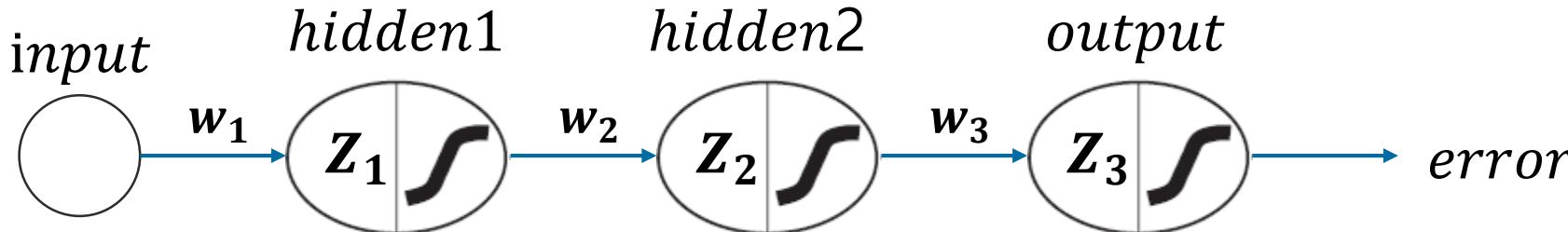
### Fine-tuning



## 03 Initialize all weights to Zero

Why not initialize all weights to ZERO?

All weights are initialized to Zero



$$\frac{\partial \text{error}}{\partial w_3} = \frac{\partial \text{error}}{\partial \text{Sigmoid}(z_3)} * \frac{\partial \text{Sigmoid}(z_3)}{\partial z_3} * \text{hidden2}$$

<0.4

$$\frac{\partial \text{error}}{\partial w_2} = \frac{\partial \text{error}}{\partial \text{Sigmoid}(z_3)} * \frac{\partial \text{Sigmoid}(z_3)}{\partial z_3} * w_3 * \frac{\partial \text{Sigmoid}(z_2)}{\partial z_2} * \text{hidden1}$$

<0.4

\*

0

\*

<0.4

$$\frac{\partial \text{error}}{\partial w_1} = \frac{\partial \text{error}}{\partial \text{Sigmoid}(z_3)} * \frac{\partial \text{Sigmoid}(z_3)}{\partial z_3} * w_3 * \frac{\partial \text{Sigmoid}(z_2)}{\partial z_2} * w_2 * \frac{\partial \text{Sigmoid}(z_1)}{\partial z_1} * \text{input}$$

<0.4

\*

0

\*

<0.4

\*

0

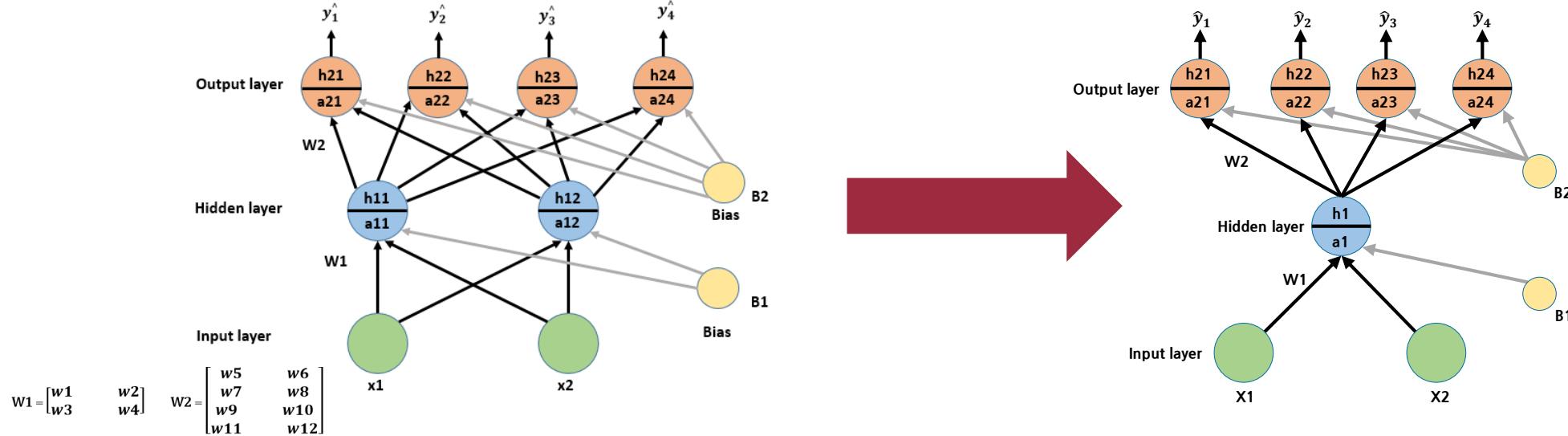
\*

<0.4

Gradient vanishing problem

## 03 Initialize all weights to Zero

Why not initialize all weights to CONSTANT VALUE?



### Forward Pass

$$a_{11} = w_1 x_1 + w_2 x_2 + b_1$$

$$a_{12} = w_3 x_1 + w_4 x_2 + b_2$$

$$a_{11} = a_{12} \text{ (assume } b_1 = b_2\text{) then } h_{11} = h_{12}$$

Similarly  $a_{21} = a_{22} = a_{23} = a_{24}$   
then  $h_{21} = h_{22} = h_{23} = h_{24}$

### Backward Pass

$$\nabla w_1 = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{21}} \frac{\partial a_{21}}{\partial h_{11}} \frac{\partial h_{11}}{\partial a_{11}} \cdot x_1$$

$$\nabla w_3 = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{21}} \frac{\partial a_{21}}{\partial h_{11}} \frac{\partial h_{11}}{\partial a_{12}} \cdot x_1$$

$$\text{here, } \frac{\partial h_{11}}{\partial a_{12}} = \frac{\partial h_{11}}{\partial a_{11}} \text{ hence } \nabla w_1 = \nabla w_3$$

*Symmetry breaking problem*

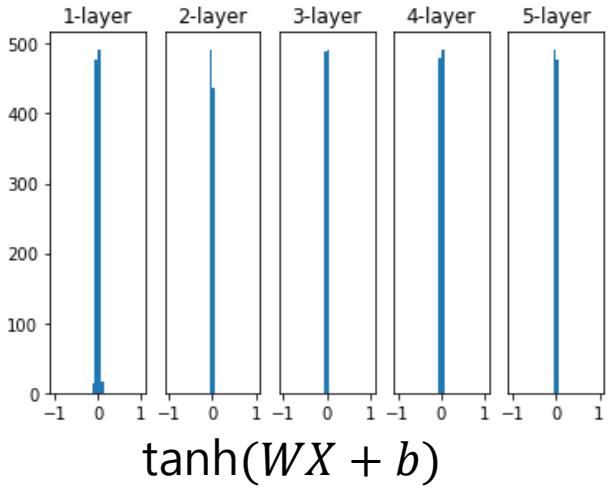
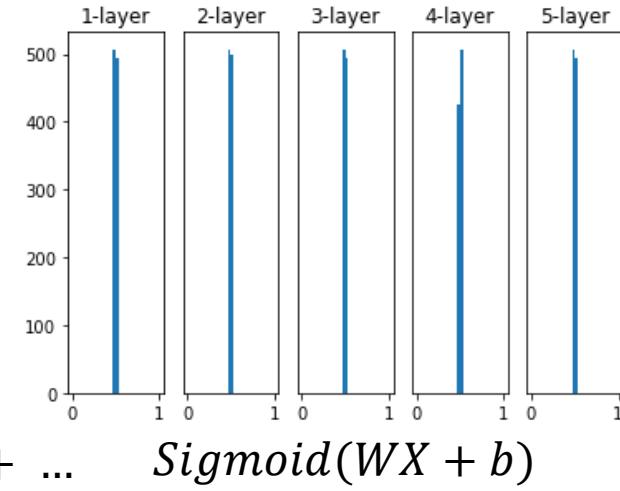
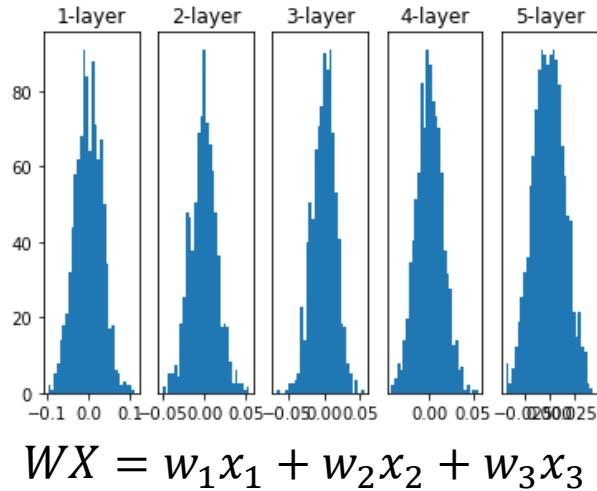
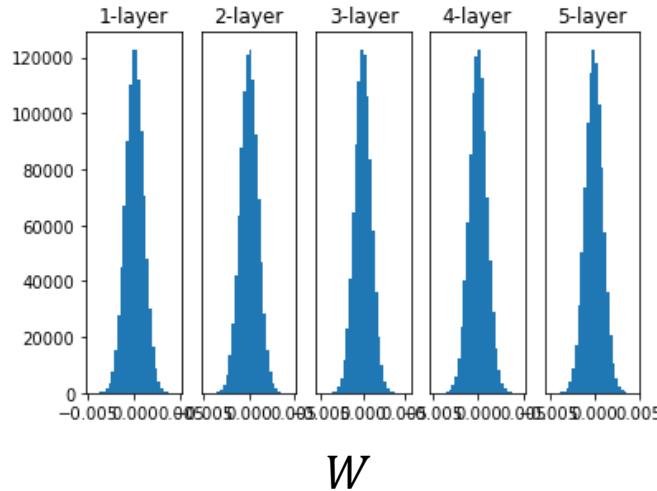
## 04 Initialize weights randomly

Can we initialize weights RANDOMLY to SMALL VALUE? -  $W \sim N(0, 0.001)$

Layer # : 5 (1000 node in each layer)

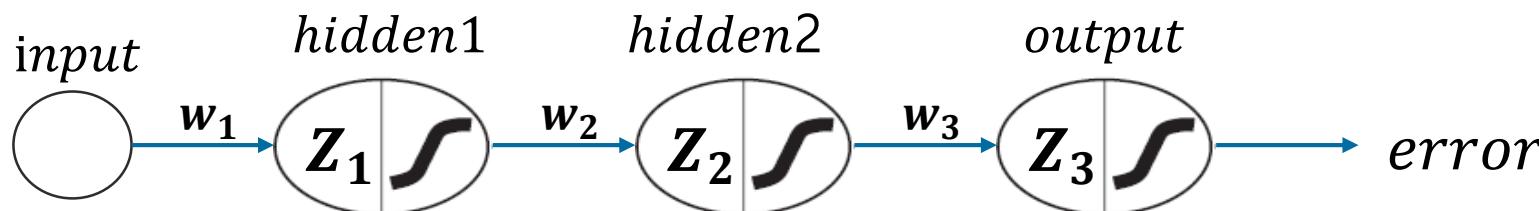
`X = np.random.randn(1, 1000)`

`W = np.random.randn(node_num, node_num)*0.001` in each layer



## 04 Initialize weights randomly

Can we initialize weights RANDOMLY to SMALL VALUE? -  $W \sim N(0, 0.001)$



$$\frac{\partial \text{error}}{\partial w_3} = \frac{\partial \text{error}}{\partial \text{activation}(z_3)} * \frac{\partial \text{activation}(z_3)}{\partial z_3} * \text{hidden2}$$

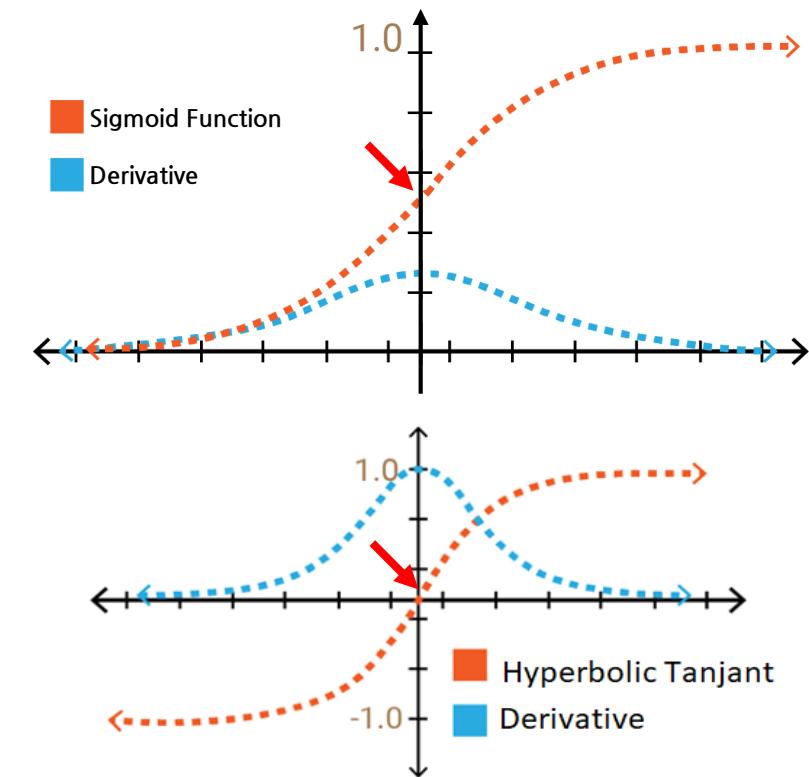
$\sim 0.4 \text{ or } 1$

$$\frac{\partial \text{error}}{\partial w_2} = \frac{\partial \text{error}}{\partial \text{activation}(z_3)} * \frac{\partial \text{activation}(z_3)}{\partial z_3} * w_3 * \frac{\partial \text{activation}(z_2)}{\partial z_2} * \text{hidden1}$$

$\sim 0.4 \text{ or } 1$        $\sim 0$        $\sim 0.4 \text{ or } 1$

$$\frac{\partial \text{error}}{\partial w_1} = \frac{\partial \text{error}}{\partial \text{activation}(z_3)} * \frac{\partial \text{activation}(z_3)}{\partial z_3} * w_3 * \frac{\partial \text{activation}(z_2)}{\partial z_2} * w_2 * \frac{\partial \text{activation}(z_1)}{\partial z_1} * \text{input}$$

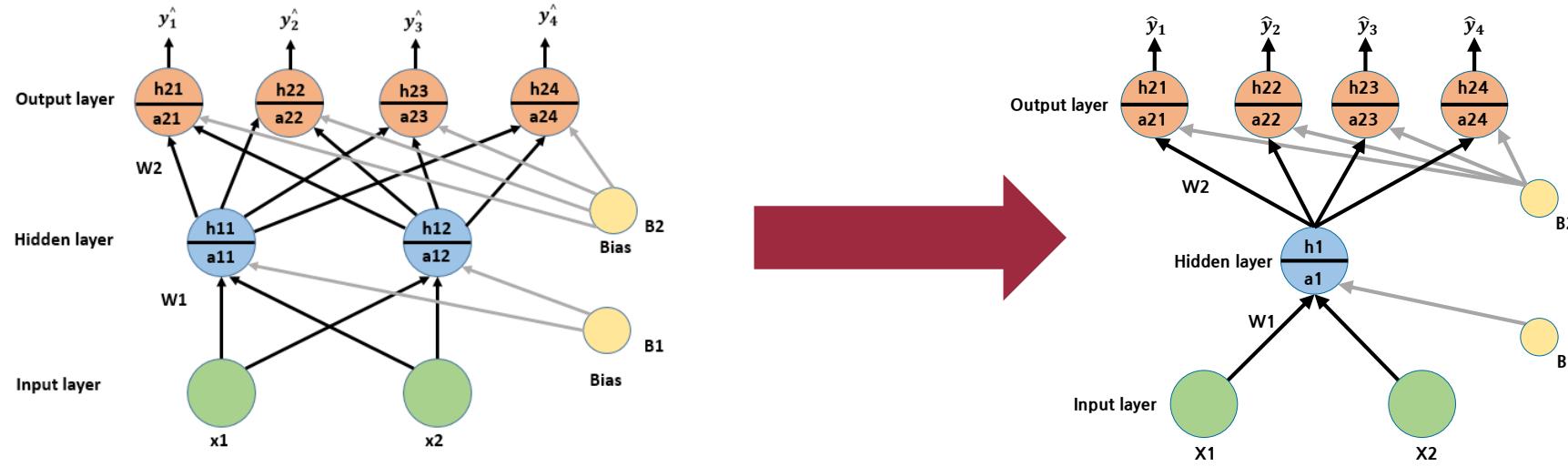
$\sim 0.4 \text{ or } 1$        $\sim 0$        $\sim 0.4 \text{ or } 1$        $\sim 0$        $\sim 0.4 \text{ or } 1$



*Alleviate gradient vanishing problem?*

## 04 Initialize weights randomly

Can we initialize weights RANDOMLY to SMALL VALUE? -  $W \sim N(0, 0.001)$



*But, still symmetry breaking problem*

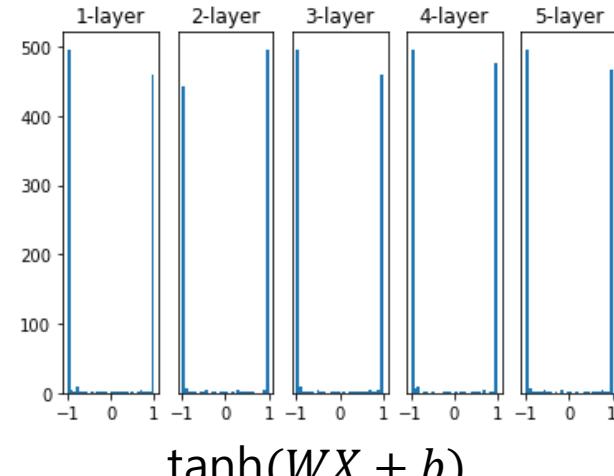
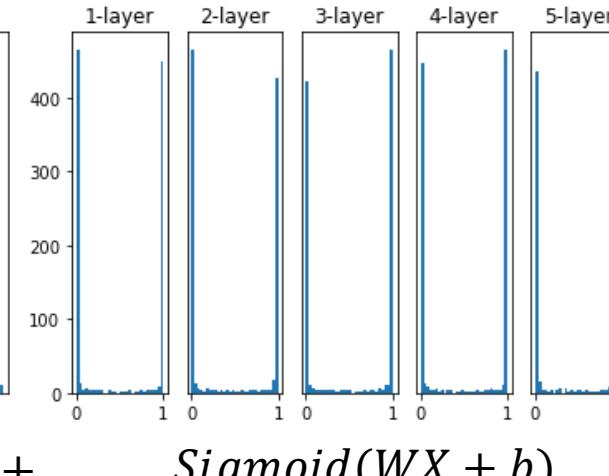
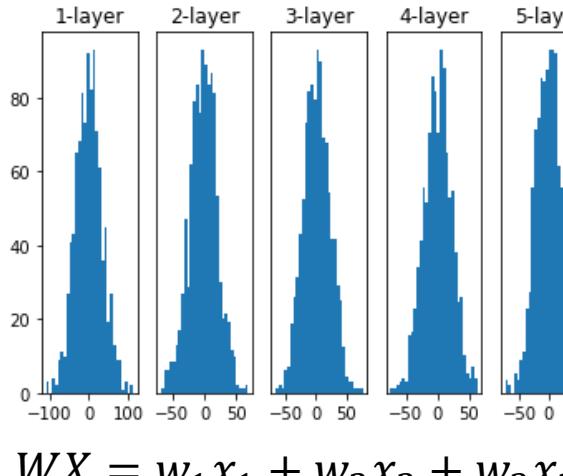
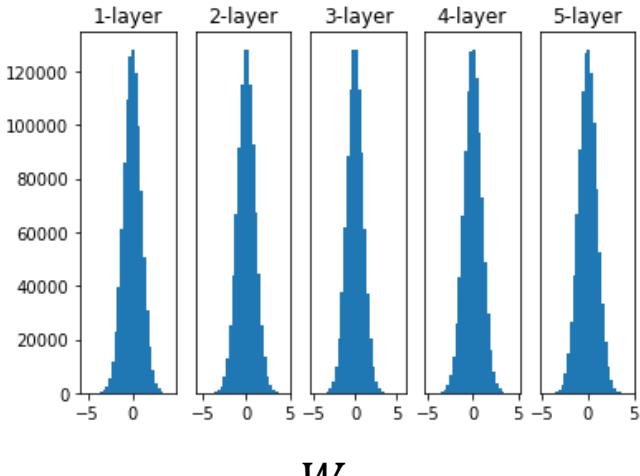
## 04 Initialize weights randomly

Can we initialize weights RANDOMLY to LARGE VALUE? -  $W \sim N(0, 1)$

Layer # : 5 (1000 node in each layer)

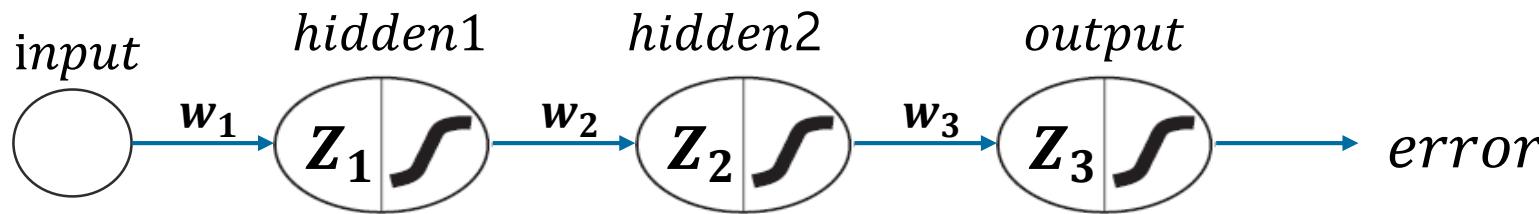
X = np.random.randn(1, 1000)

W = np.random.randn(node\_num, node\_num) in each layer



## 04 Initialize weights randomly

Can we initialize weights RANDOMLY to LARGE VALUE? -  $W \sim N(0, 1)$



$$\frac{\partial \text{error}}{\partial w_3} = \frac{\partial \text{error}}{\partial \text{activation}(z_3)} * \frac{\partial \text{activation}(z_3)}{\partial z_3} * \text{hidden2}$$

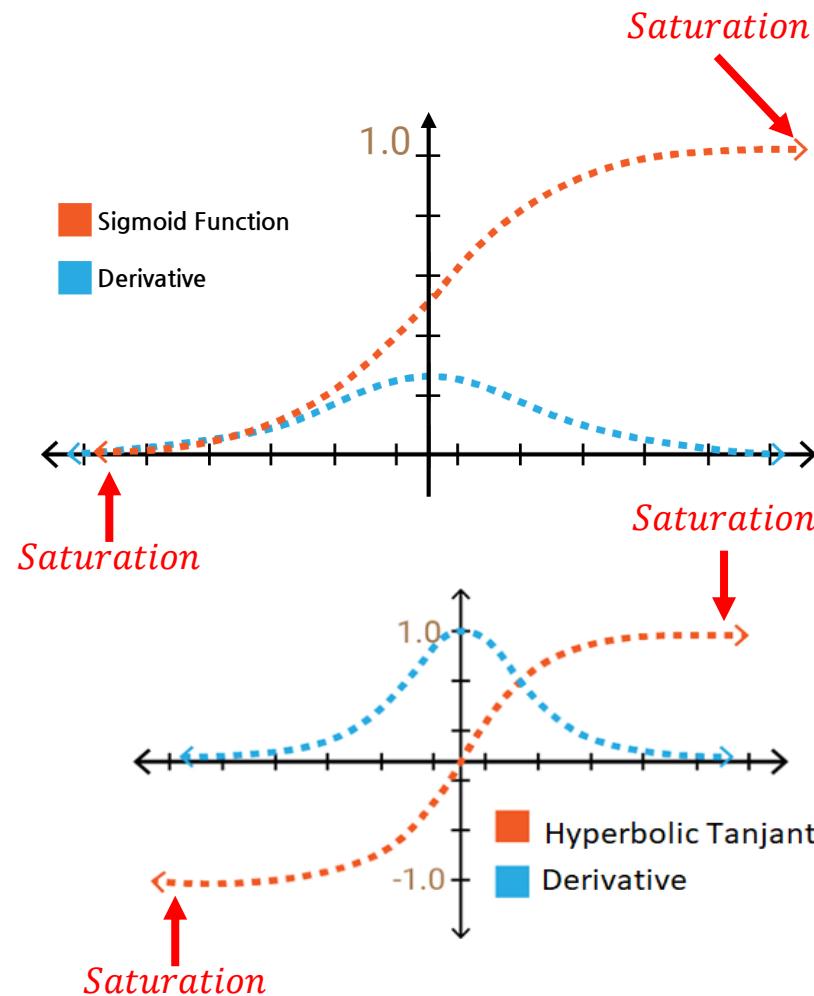
$\sim 0$

$$\frac{\partial \text{error}}{\partial w_2} = \frac{\partial \text{error}}{\partial \text{activation}(z_3)} * \frac{\partial \text{activation}(z_3)}{\partial z_3} * w_3 * \frac{\partial \text{activation}(z_2)}{\partial z_2} * \text{hidden1}$$

$\sim 0$

$$\frac{\partial \text{error}}{\partial w_1} = \frac{\partial \text{error}}{\partial \text{activation}(z_3)} * \frac{\partial \text{activation}(z_3)}{\partial z_3} * w_3 * \frac{\partial \text{activation}(z_2)}{\partial z_2} * w_2 * \frac{\partial \text{activation}(z_1)}{\partial z_1} * \text{input}$$

$\sim 0$



*Gradient vanishing problem*

# 05 Xavier initialization for linear activation function

[Lecun1998] Efficient BackProp

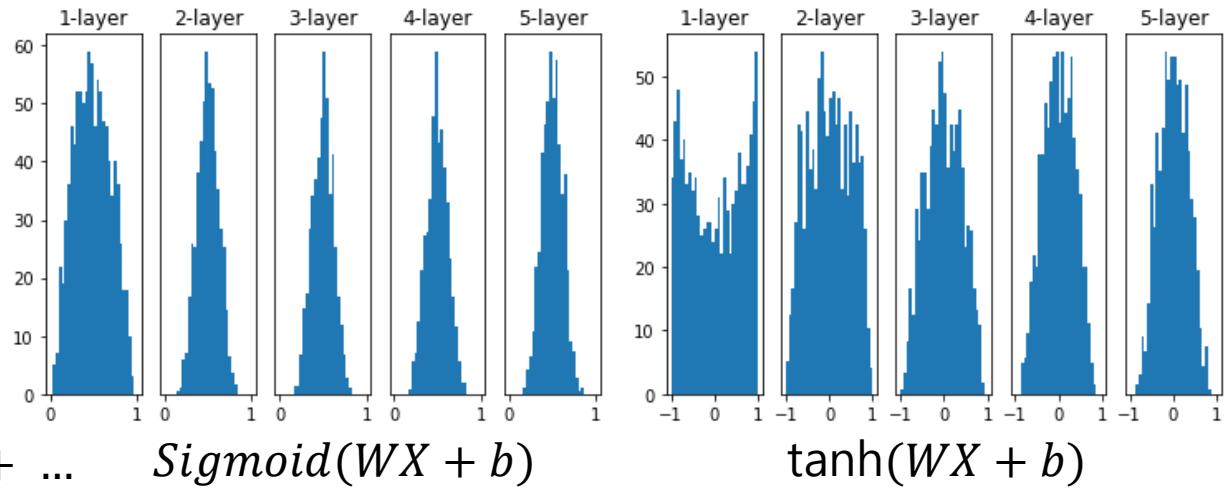
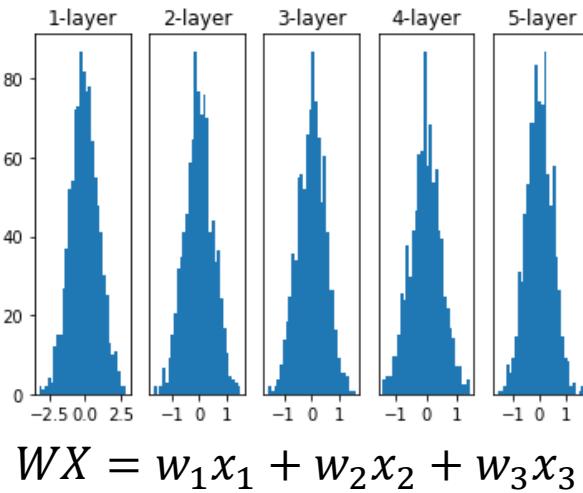
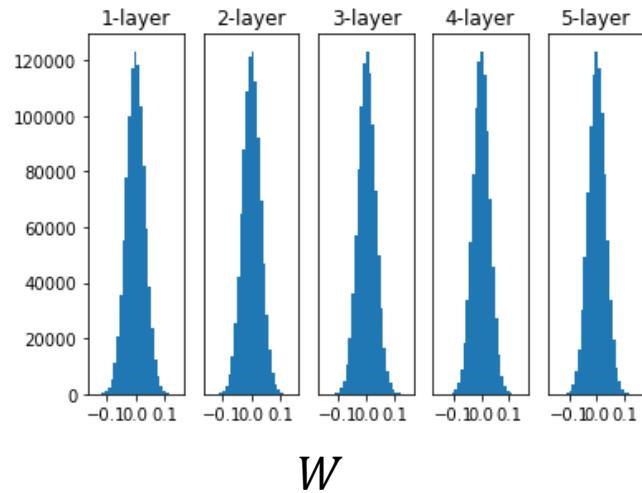
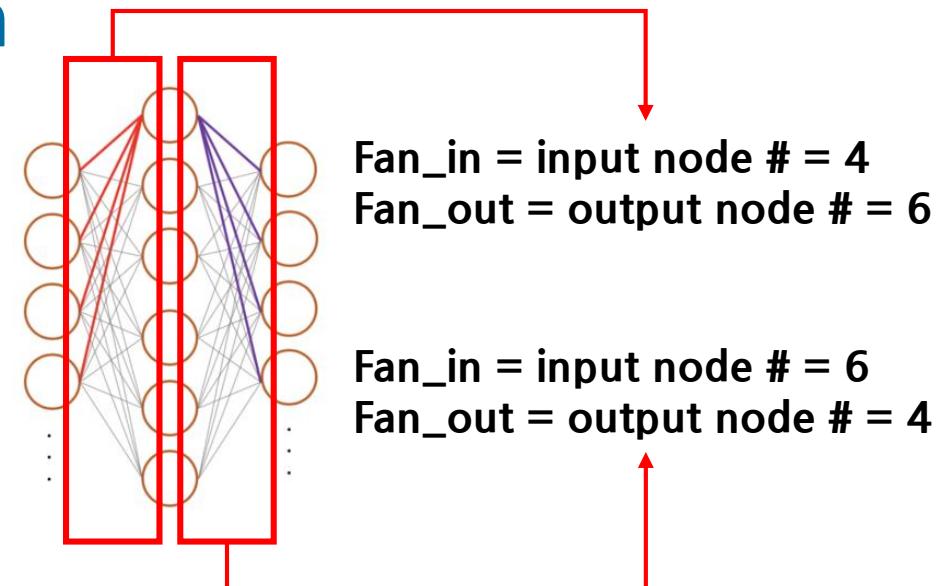
$$W \sim N(0, Var(W))$$

$$Var(W) = \frac{1}{n_{in}}$$

Layer # : 5 (1000 node in each layer)

X = np.random.randn(1, 1000)

W = np.random.randn(node\_num, node\_num) in each layer



*Solve all problem?*

## 05 Xavier initialization for linear activation function

[Glorot2010] Understanding the difficulty of training deep feedforward neural networks.

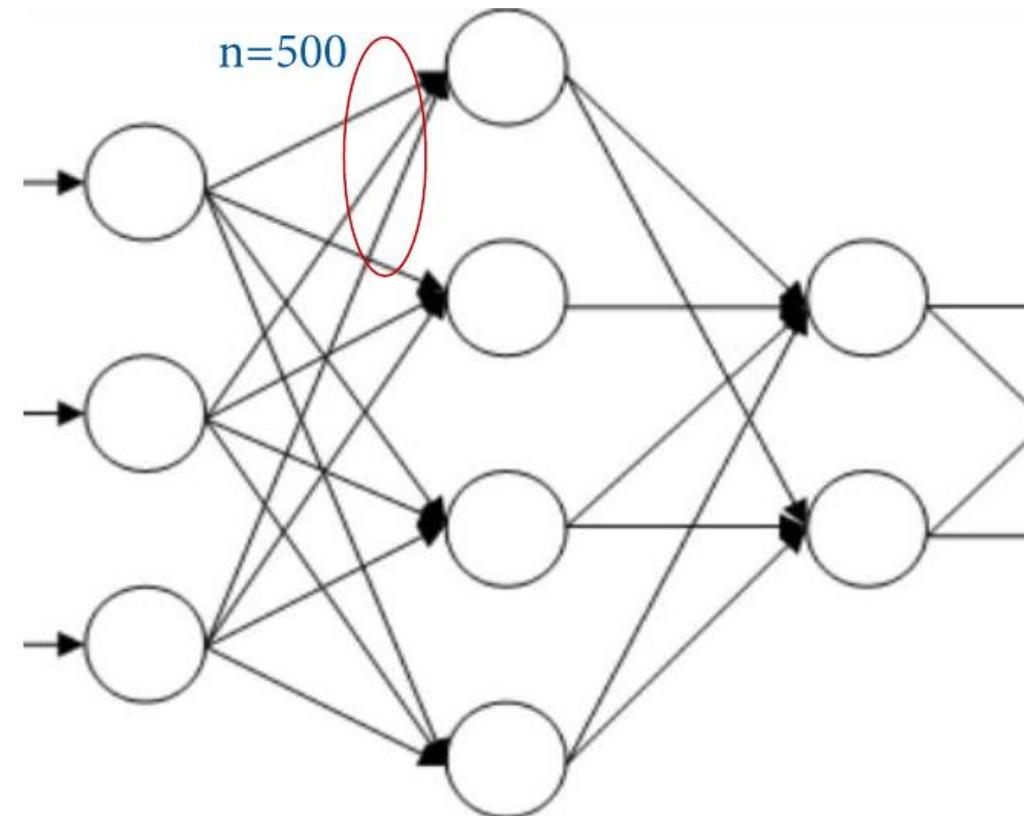
$$\cancel{WX} \sim N(0, 1)$$

$$Z = \sum WX$$

$$h = \text{Sigmoid}(Z)$$

$$E(Z) = 0$$

$$\text{Var}(Z) = 1$$



## 05 Xavier initialization for linear activation function

[Glorot2010] Understanding the difficulty of training deep feedforward neural networks.

$$\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2$$

$$\text{Var}(X + a) = \text{Var}(X)$$

$$\text{Var}(aX) = a^2 \text{Var}(X)$$

Cov(X,X)

If X,Y are uncorrelated or independent,  
we can remove these terms

$$\text{Var}(aX + bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X, Y)$$

$$\text{Var}(aX - bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) - 2ab \text{Cov}(X, Y)$$

$$\text{Var}\left(\sum_{i=1}^N X_i\right) = \sum_{i,j=1}^N \text{Cov}(X_i, X_j) = \sum_{i=1}^N \text{Var}(X_i) + \sum_{i \neq j} \text{Cov}(X_i, X_j)$$

## 05 Xavier initialization for linear activation function

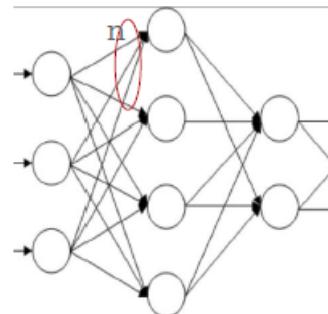
[Glorot2010] Understanding the difficulty of training deep feedforward neural networks.

$$\begin{aligned} \text{Var}(ax+by) &\equiv E[(ax+by)(ax+by)] - \underbrace{E[ax+by]}_{=0} E[ax+by] \\ &= E[a^2x^2 + 2abxy + b^2y^2] - \cancel{E[ax+by]^2} \\ &= E[a^2x^2] + 2E[abxy] + E[b^2y^2] - \cancel{E[ax+by]^2} \\ &= a^2E[x^2] + 2abE[xy] + b^2E[y^2] - \cancel{E[ax+by]^2} \\ &= a^2E[x^2] + 2abE[xy] + b^2E[y^2] - \{E(ax) + E(by)\}^2 \\ &= a^2E[x^2] + 2abE[xy] + b^2E[y^2] - \{aE(x) + bE(y)\}^2 \\ &= a^2E[x^2] + 2abE[xy] + b^2E[y^2] - (a^2E[x]^2 + 2abE[x]E[y] + b^2E[y]^2) \\ &= 2abE[xy] - 2abE[x]E[y] \\ &= a^2(E[x^2] - E[x]^2) + b^2(E[y^2] - E[y]^2) + 2ab(E[x]E[y] - E[x]E[y]) \\ &= a^2\text{Var}(x) + b^2\text{Var}(y) + 2ab\text{Cov}(x, y) \end{aligned}$$

## 05 Xavier initialization for linear activation function

[Glorot2010] Understanding the difficulty of training deep feedforward neural networks.

$$\begin{aligned} \text{Var}(s) &= \text{Var}\left(\sum_i^n w_i x_i\right) && \text{See this as one variable} \\ &\quad \downarrow && \text{Each } \{w_i x_i\} \text{ is independent} \\ z = \sum_j w_j x_j + b &= \sum_i^n \text{Var}(w_i x_i) && \text{Var}\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n \text{Var}(X_i) \\ &\quad \downarrow && \text{two variables } w \text{ and } x \text{ are independent. see next slide} \\ &= \sum_i^n [E(w_i)]^2 \text{Var}(x_i) + E[(x_i)]^2 \text{Var}(w_i) + \text{Var}(x_i) \text{Var}(w_i) && \\ &\quad \downarrow && \text{assumption: zero mean inputs and weights} \\ &= \sum_i^n \text{Var}(x_i) \text{Var}(w_i) && \\ &\quad \downarrow && \text{Assumption: } w_i, x_i \text{ are drawn from identical distributions, respectively} \\ &= (n \text{Var}(w)) \text{Var}(x) && \end{aligned}$$



## 05 Xavier initialization for linear activation function

[Glorot2010] Understanding the difficulty of training deep feedforward neural networks.

$$\text{Var}(Y) = E[XY\bar{X}\bar{Y}] - E[Y\bar{X}\bar{Y}]E[X\bar{Y}] \quad \text{by definition}$$

$$= E[X^2Y^2] - \cancel{E[X^2]\cancel{E[Y^2]}} - E[X]^2$$

$$= E[X^2]E[Y^2] - \cancel{E[X]^2} \cancel{E[Y]^2} \quad X, Y \text{ independent}$$

$$\text{Var}(Y) = E[Y^2] - E[Y]^2$$

$$(var(X) + E[X]^2)(var(Y) + E[Y]^2) - E[X]^2E[Y]^2$$

$$= \text{Var}(X)\text{Var}(Y) + E[X]^2\text{Var}(Y) + E[Y]^2\text{Var}(X) + E[X]^2E[Y]^2 - E[X]^2E[Y]^2$$

$$= E[X]^2\text{Var}(Y) + E[Y]^2\text{Var}(X) + \text{Var}(X)\text{Var}(Y)$$

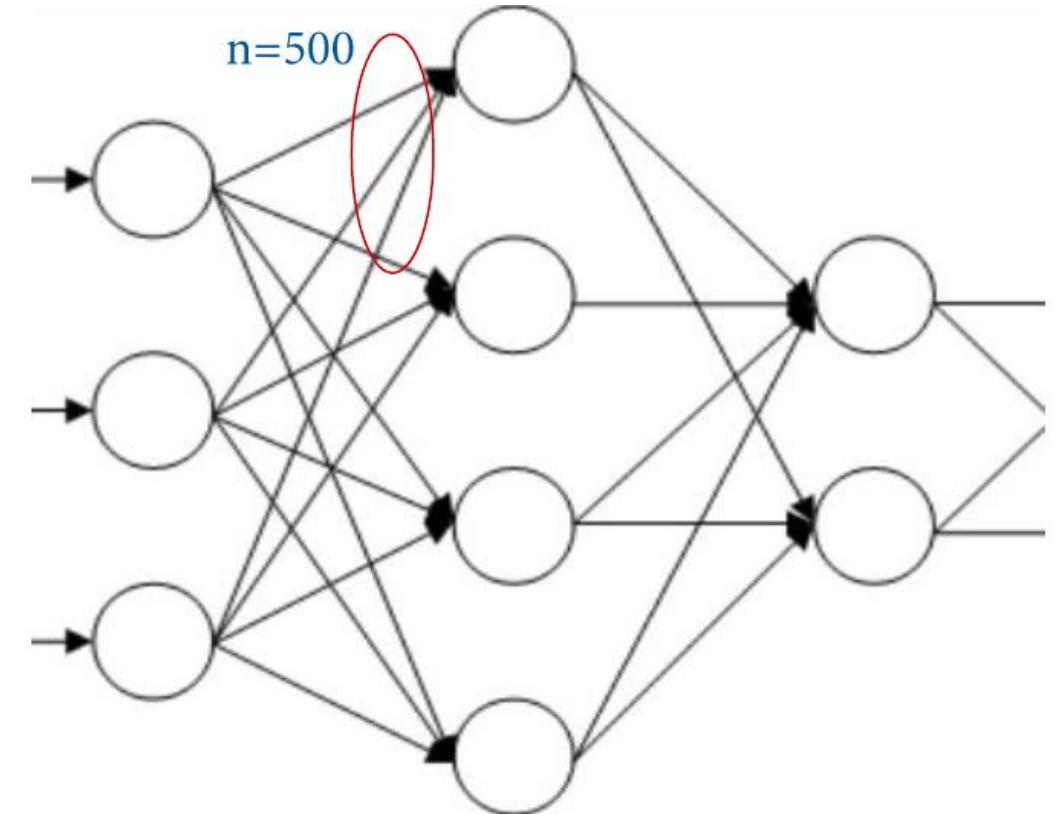
## 05 Xavier initialization for linear activation function

[Glorot2010] Understanding the difficulty of training deep feedforward neural networks.

$$Var(Z) = (nVar(W))Var(X)$$

$$Var(W) = \frac{1}{n_{in}} = \frac{1}{500}$$

$$W \sim N(0, Var(W))$$



## 05 Xavier initialization for linear activation function

[Glorot2010] Understanding the difficulty of training deep feedforward neural networks.

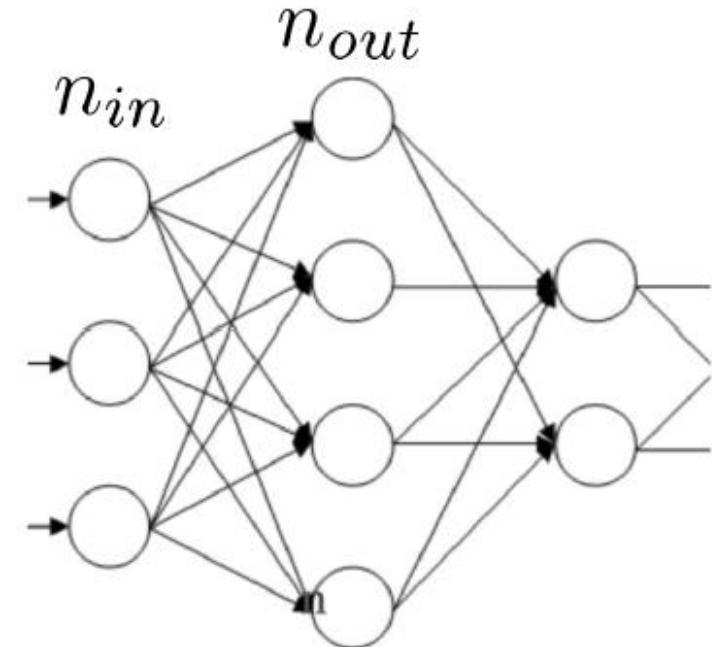
- What this initialization method is considering

- Not only the forward (as before)
  - for forward, there are  $n_{in}$  inputs for each neuron
- But also backpropagation
  - for backward, there are  $n_{out}$  input for each neuron

$$Var(W) = \frac{2}{n_{in} + n_{out}}$$

*Harmonic Mean*

$$W \sim N(0, Var(W))$$

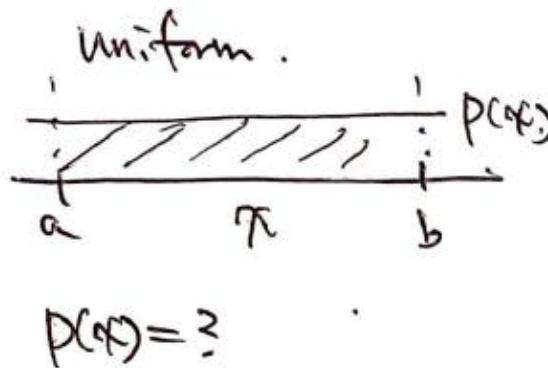


<https://www.programmersought.com/article/73876303829/>

## 05 Xavier initialization for linear activation function

[Glorot2010] Understanding the difficulty of training deep feedforward neural networks.

$$\begin{aligned} \text{Var}(x) &= E[x^2] - E[x]E[x] = E[x^2] - \left(\frac{a+b}{2}\right)\left(\frac{a+b}{2}\right) \\ &= \int_a^b p(x)x^2 - \left(\frac{a+b}{2}\right)^2 = \int_a^b \frac{1}{b-a}x^2 - \left(\frac{a+b}{2}\right)^2 \\ &= \frac{1}{(b-a)} \frac{1}{3}x^3 \Big|_a^b - \left(\frac{a+b}{2}\right)^2 = \frac{1}{b-a} \frac{1}{3} (b^3 - a^3) - \left(\frac{a+b}{2}\right)^2 \\ &= \frac{1}{3} \frac{1}{b-a} (b-a)(b^2 + ab + a^2) - \left(\frac{a+b}{2}\right)^2 = \frac{1}{3} (a^2 + ab + b^2) - \frac{a^2 + 2ab + b^2}{4} \\ &= \frac{4a^2 + 4ab + 4b^2 - 3a^2 - 6ab - 3b^2}{12} = \frac{a^2 - 2ab + b^2}{12} = \boxed{\frac{(a-b)^2}{12}}. \end{aligned}$$



$$P(x) = ?$$

$$\downarrow$$

$$P(x)(b-a) = 1.$$

$$P(x) = \frac{1}{b-a}.$$

## 05 Xavier initialization for linear activation function

[Glorot2010] Understanding the difficulty of training deep feedforward neural networks.

$$n \text{Var}(W) = 1$$

$$n \frac{(a+b)^2}{12} = 1.$$

$$b=-a \Rightarrow n \frac{(a+b)^2}{12} = n \frac{(2a)^2}{12} = n \cdot \frac{4a^2}{12} = \frac{n a^2}{3}$$

$$\frac{n a^2}{3} = 1 \quad a^2 = \frac{3}{n} \quad a = \pm \sqrt{\frac{3}{n}}$$

∴ Draw sample from  $[-\sqrt{\frac{3}{n}}, \sqrt{\frac{3}{n}}]$  uniformly

## 05 Xavier initialization for linear activation function

[Glorot2010] Understanding the difficulty of training deep feedforward neural networks.

- 1) draw sample from Uniform

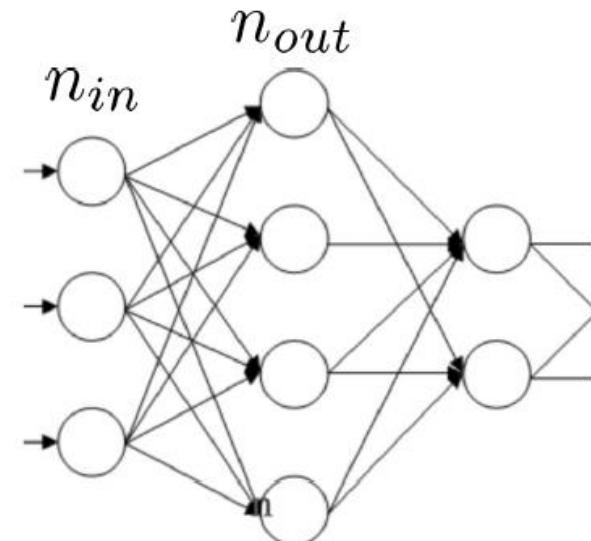
$$\left[ -\sqrt{\frac{6}{n_i + n_{i+1}}}, \sqrt{\frac{6}{n_i + n_{i+1}}} \right]$$

$$\text{Var}(W) = \frac{1}{n} \quad (\text{Gaussian})$$

$$\text{Var}(W) = \frac{3}{n} \quad (\text{Uniform})$$

$$(\text{Gaussian})\text{Var}(W) = \frac{2}{n_i + n_{i+1}} \Rightarrow (\text{Uniform})\text{Var}(W) = \frac{2}{n_i + n_{i+1}} \times 3 = \frac{6}{n_i + n_{i+1}}$$

thus when Gaussian  $\rightarrow$  Uniform,  
we can multiply by "3"



# 05 Xavier initialization for linear activation function

[Glorot2010] Understanding the difficulty of training deep feedforward neural networks.

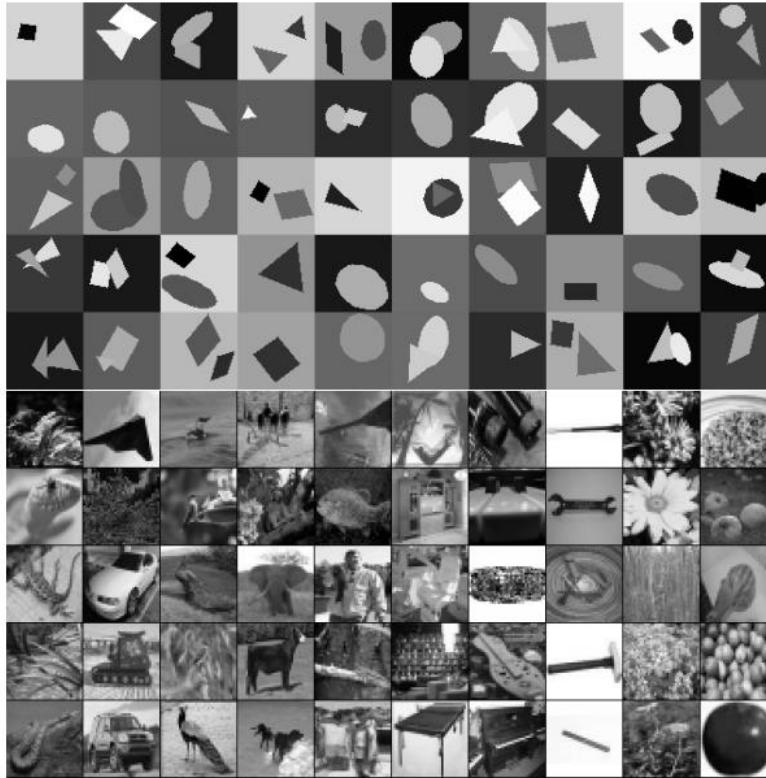


Figure 1: Top: Shapeset-3  $\times$  2 images at  $64 \times 64$  resolution. The examples we used are at  $32 \times 32$  resolution. The learner tries to predict which objects (parallelogram, triangle, or ellipse) are present, and 1 or 2 objects can be present, yielding 9 possible classifications. Bottom: Small-ImageNet images at full resolution.

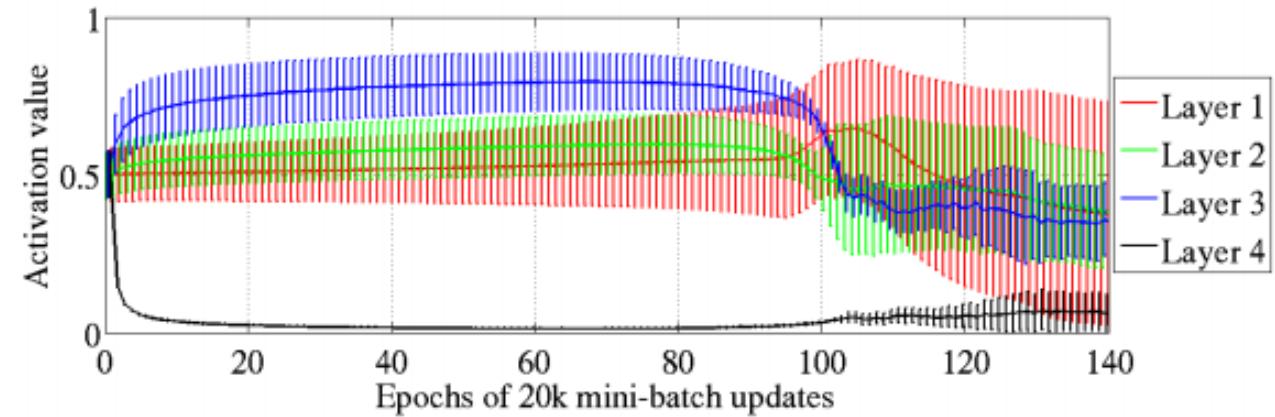


Figure 2: Mean and standard deviation (vertical bars) of the activation values (output of the sigmoid) during supervised learning, for the different hidden layers of a deep architecture. The top hidden layer quickly saturates at 0 (slowing down all learning), but then slowly desaturates around epoch 100.

# 05 Xavier initialization for linear activation function

[Glorot2010] Understanding the difficulty of training deep feedforward neural networks.

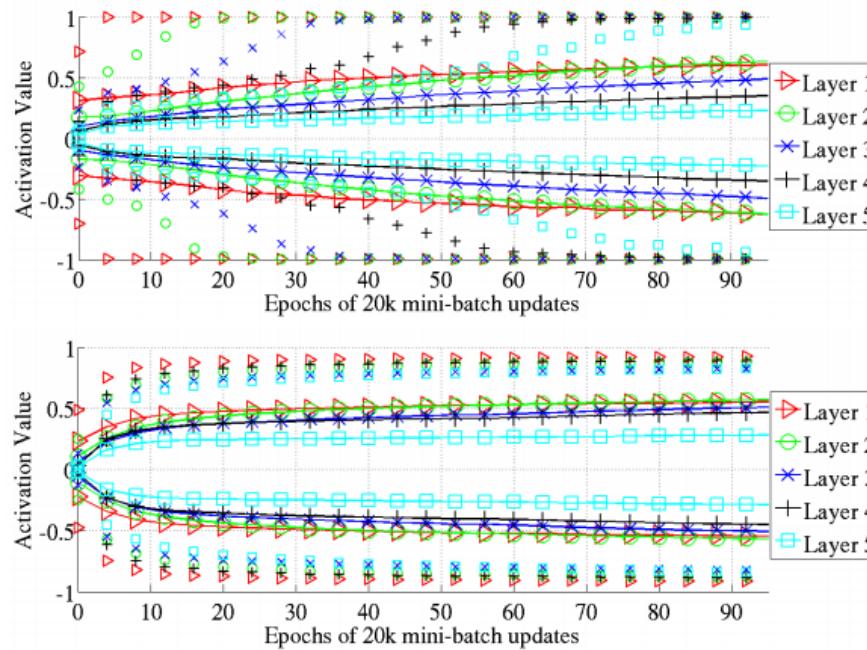


Figure 3: Top: 98 percentiles (markers alone) and standard deviation (solid lines with markers) of the distribution of the activation values for the hyperbolic tangent networks in the course of learning. We see the first hidden layer saturating first, then the second, etc. Bottom: 98 percentiles (markers alone) and standard deviation (solid lines with markers) of the distribution of activation values for the soft-sign during learning. Here the different layers saturate less and do so together.

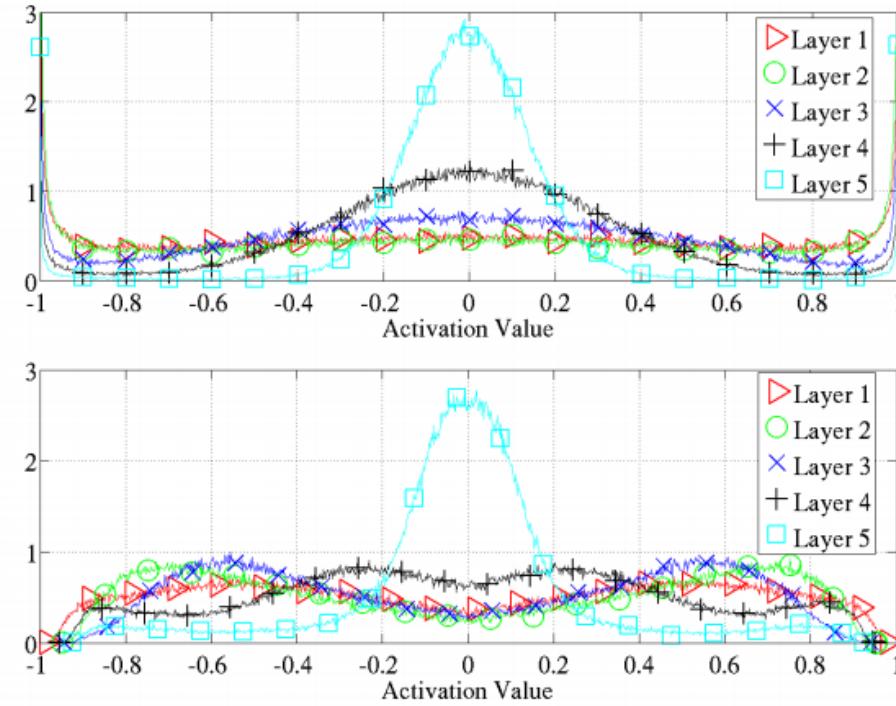


Figure 4: Activation values normalized histogram at the end of learning, averaged across units of the same layer and across 300 test examples. Top: activation function is hyperbolic tangent, we see important saturation of the lower layers. Bottom: activation function is softsign, we see many activation values around (-0.6,-0.8) and (0.6,0.8) where the units do not saturate but are non-linear.

## 05 Xavier initialization for linear activation function

[Glorot2010] Understanding the difficulty of training deep feedforward neural networks.

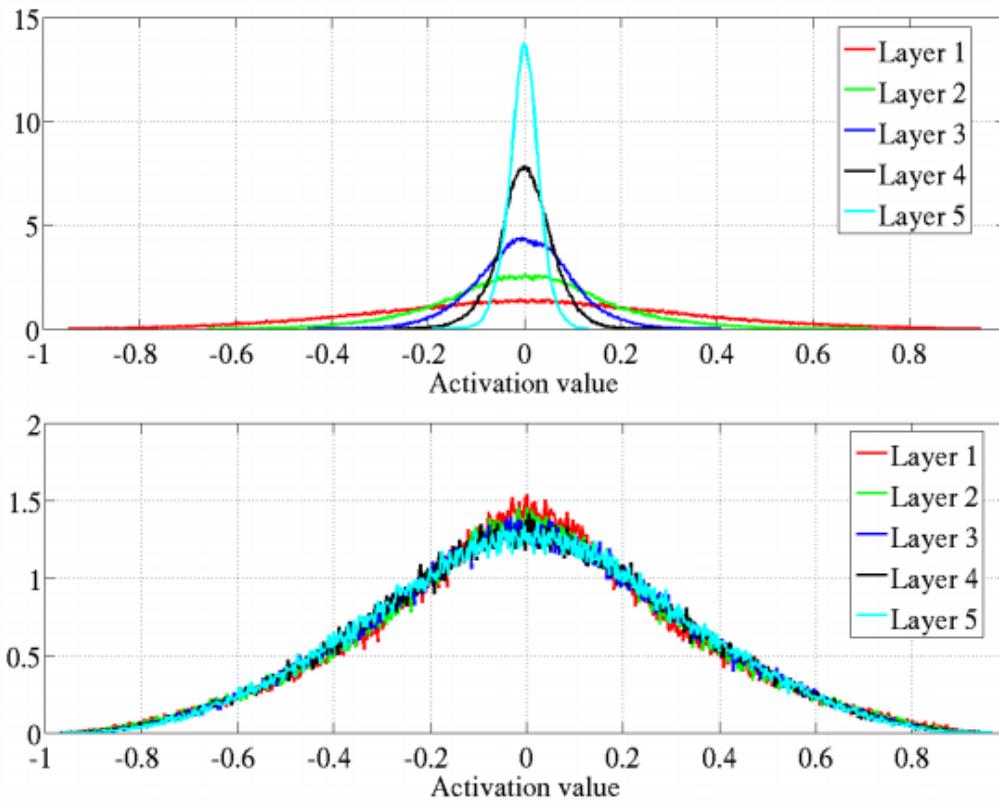


Figure 6: Activation values normalized histograms with hyperbolic tangent activation, with standard (top) vs normalized initialization (bottom). Top: 0-peak increases for higher layers.

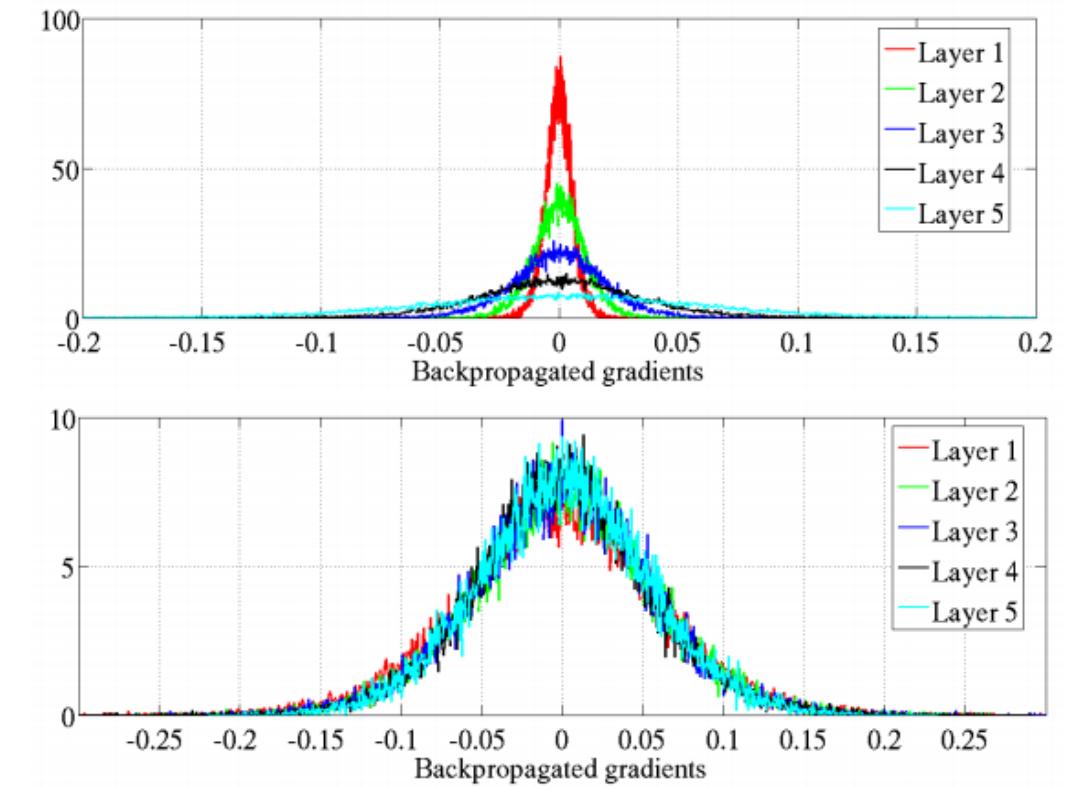


Figure 7: Back-propagated gradients normalized histograms with hyperbolic tangent activation, with standard (top) vs normalized (bottom) initialization. Top: 0-peak decreases for higher layers.

## 05 Xavier initialization for linear activation function

[Glorot2010] Understanding the difficulty of training deep feedforward neural networks.

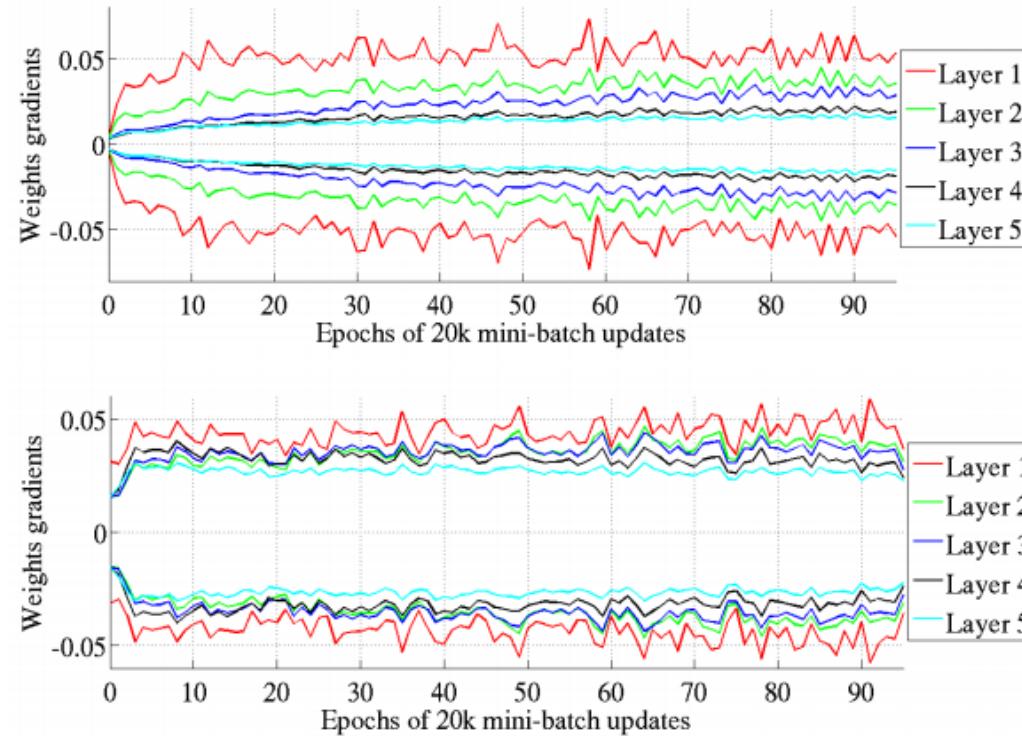


Figure 9: Standard deviation intervals of the weights gradients with hyperbolic tangents with standard initialization (top) and normalized (bottom) during training. We see that the normalization allows to keep the same variance of the weights gradient across layers, during training (top: smaller variance for higher layers).

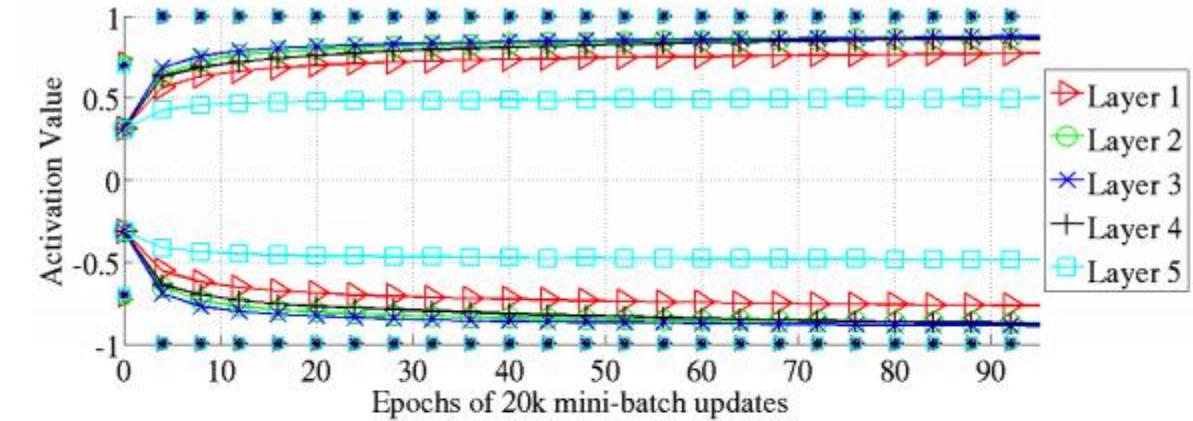


Figure 10: 98 percentile (markers alone) and standard deviation (solid lines with markers) of the distribution of activation values for hyperbolic tangent with normalized initialization during learning.

activations (flowing upward) and gradients (flowing backward).

## 05 Xavier initialization for linear activation function

[Glorot2010] Understanding the difficulty of training deep feedforward neural networks.

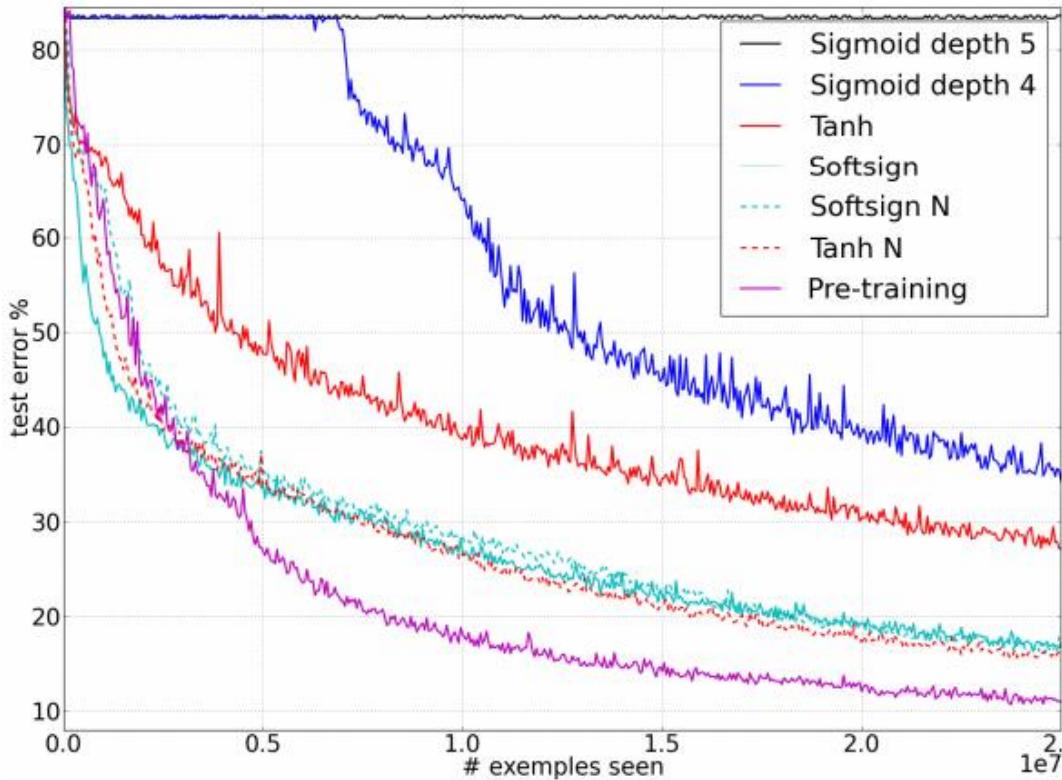


Figure 11: Test error during online training on the **Shapenet-3x2** dataset, for various activation functions and initialization schemes (ordered from top to bottom in decreasing final error). N after the activation function name indicates the use of normalized initialization.

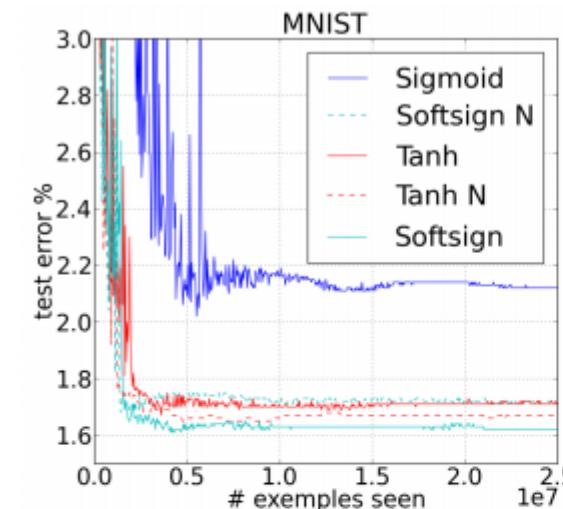
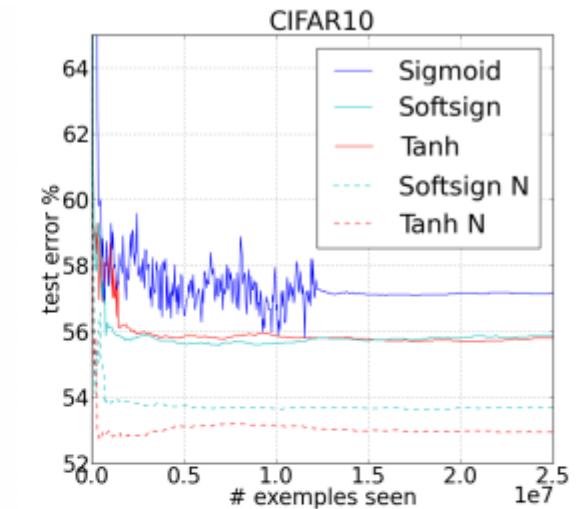


Figure 12: Test error curves during training on **MNIST** and **CIFAR10**, for various activation functions and initialization schemes (ordered from top to bottom in decreasing final error). N after the activation function name indicates the use of normalized initialization.



## 05 Xavier initialization for linear activation function

[Glorot2010] Understanding the difficulty of training deep feedforward neural networks.

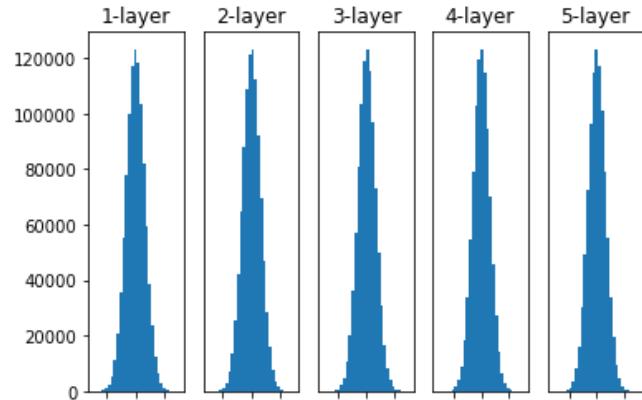
$$W \sim N(0, Var(W))$$

$$Var(W) = \frac{2}{n_{in} + n_{out}}$$

Layer # : 5 (1000 node in each layer)

X = np.random.randn(1, 1000)

W = np.random.randn(node\_num, node\_num) in each layer

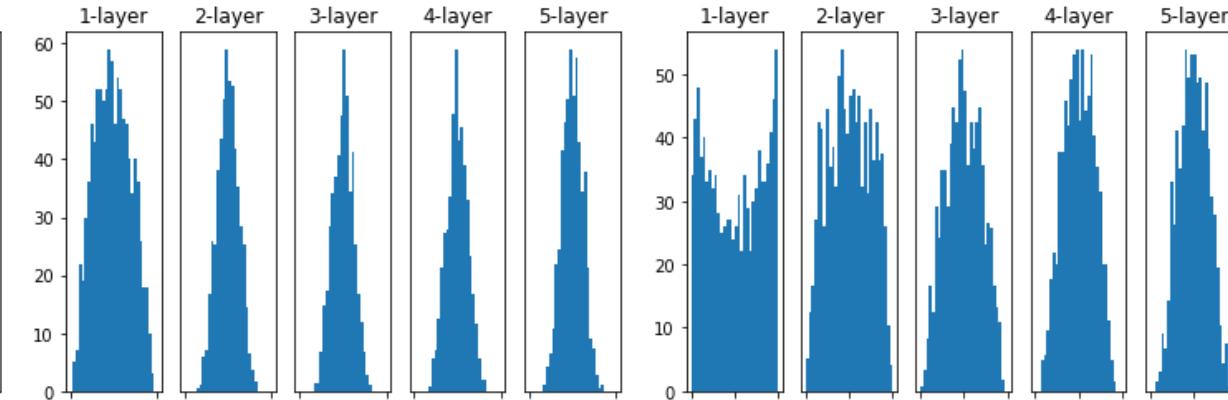
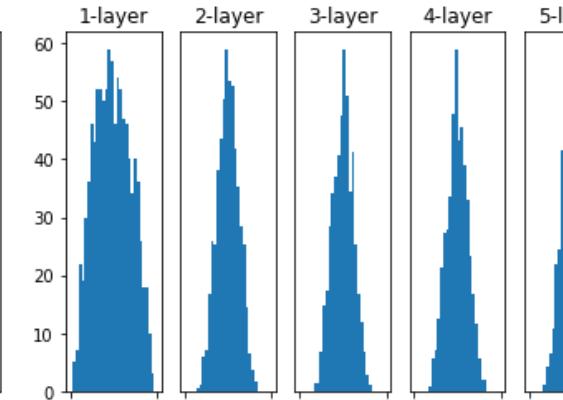
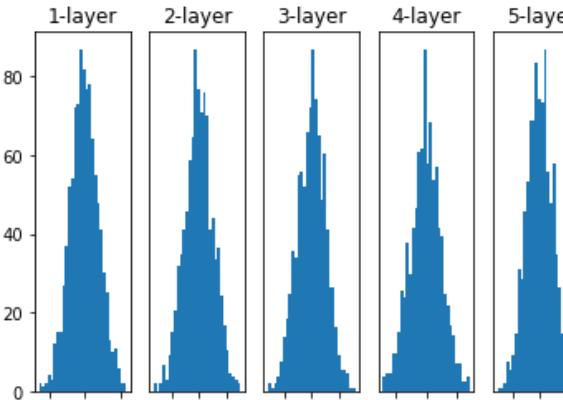


$W$

$$WX = w_1x_1 + w_2x_2 + w_3x_3 + \dots$$

$$\text{Sigmoid}(WX + b)$$

$$\tanh(WX + b)$$



# 06 He initialization for ReLU activation function

[He2015] Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification.

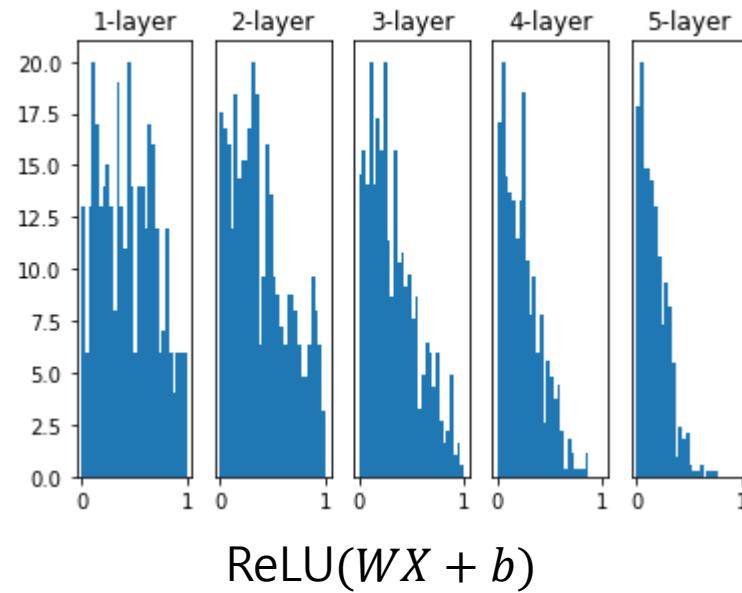
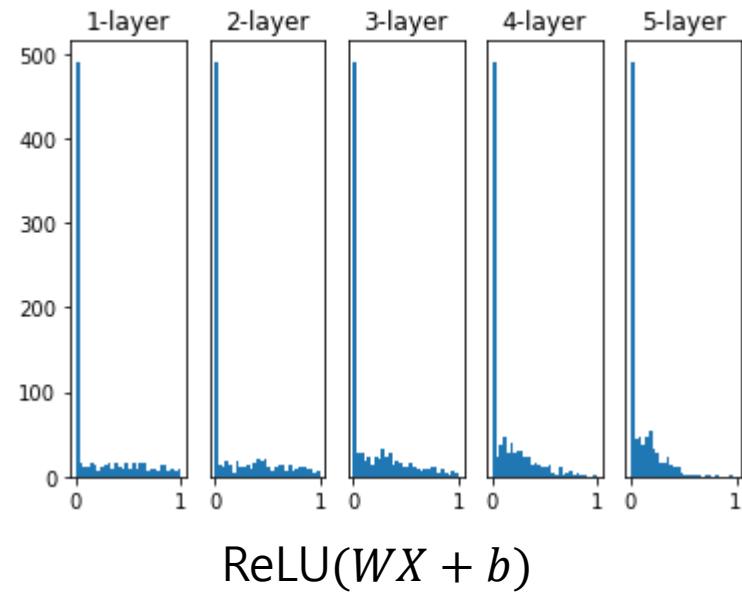
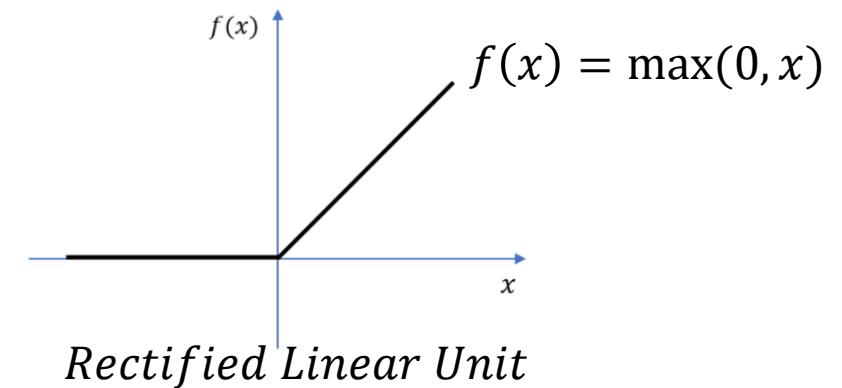
$$W \sim N(0, Var(W))$$

$$Var(W) = \frac{2}{n_{in} + n_{out}}$$

Layer # : 5 (1000 node in each layer)

X = np.random.randn(1, 1000)

W = np.random.randn(node\_num, node\_num) in each layer



# 06 He initialization for ReLU activation function

[He2015] Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification.

**Forward Pass :**

$$\mathbf{y}_l = \mathbf{W}_l \mathbf{x}_l + \mathbf{b}_l$$

$$\mathbf{x}_l = f(\mathbf{y}_{l-1})$$

$$Y = W_1 X_1 + W_2 X_2 + \dots + W_n X_n$$

$$\text{Var}[y_l] = n_l \text{Var}[w_l x_l]$$

$$\text{Var}[y_l] = n_l \text{Var}[w_l] E[x_l^2]$$

$$\text{Var}[y_l] = \frac{1}{2} n_l \text{Var}[w_l] \text{Var}[y_{l-1}]$$

$$\frac{1}{2} n_l \text{Var}[w_l] = 1$$

$$\text{Var}(W) = \frac{2}{n_l}$$

$$W \sim N(0, \text{Var}(W))$$

*Assume all  $x, w$  are independent*

$$\text{Var}[w_l x_l] = E[w_l^2] E[x_l^2] - [E[w_l]]^2 [E[x_l]]^2$$

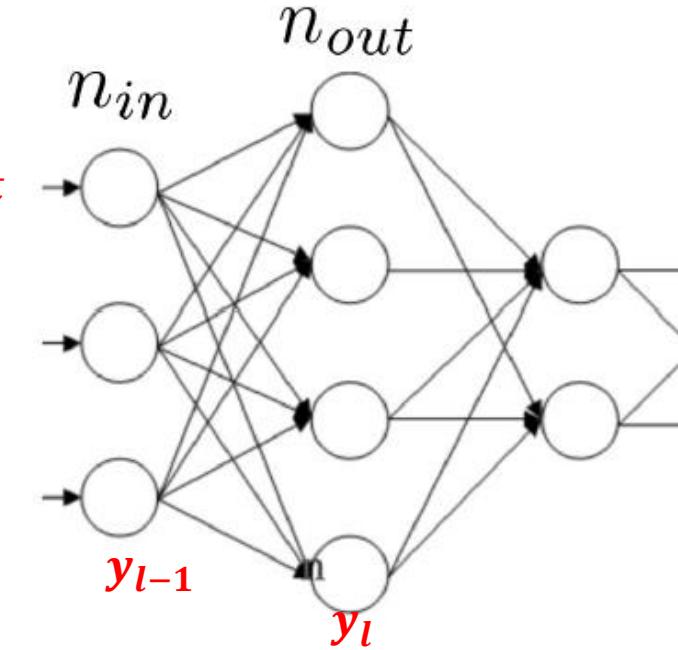
$$\text{Var}[w_l x_l] = \boxed{E[w_l^2]} E[x_l^2]$$

$$\text{Var}[w_l] = E[w_l^2] - [E[w_l]]^2$$

$$\text{Var}[w_l] = \boxed{E[w_l^2]}$$

$$\text{Var}[w_l x_l] = \boxed{E[w_l^2]} \boxed{E[x_l^2]}$$

$$\text{Var}[w_l x_l] = \underline{\underline{\text{Var}[w_l] E[x_l^2]}}$$



$$\begin{aligned} E[x^2] &= \int_{-\infty}^{\infty} x^2 P(x) dx \\ &= \int_{-\infty}^{\infty} \max(0, y)^2 P(y) dy \\ &= \int_0^{\infty} y^2 P(y) dy \\ &= 0.5 * \int_{-\infty}^{\infty} y^2 P(y) dy \\ &= 0.5 * \text{Var}(y) \end{aligned}$$

$$\begin{aligned} E(x_l^2) &= \frac{1}{\sigma_{y_{l-1}} \sqrt{2\pi}} \int_0^{\infty} x^2 e^{-\frac{x^2}{2\sigma_{y_{l-1}}^2}} dx \\ &= \left(\frac{1}{2}\right) \frac{1}{\sigma_{y_{l-1}} \sqrt{2\pi}} \int_{-\infty}^{\infty} x^2 e^{-\frac{x^2}{2\sigma_{y_{l-1}}^2}} dx \\ &= \frac{\sigma_{y_{l-1}}^2}{2} = \frac{1}{2} \text{Var}(y_{l-1}) \end{aligned}$$

# 06 He initialization for ReLU activation function

[He2015] Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification.

## Backward Pass :

$$\Delta \mathbf{x}_l = \hat{\mathbf{W}}_l \Delta \mathbf{y}_l$$

$$\Delta y_l = f'(y_l) \Delta x_{l+1}$$

$$Var[\Delta x_l] = \hat{n}_l Var[w_l] Var[\Delta y_l]$$

$$= \frac{1}{2} \hat{n}_l Var[w_l] Var[\Delta x_{l+1}]$$

ReLU: 1/2

$$E[\Delta y_l] = E[f'(y_l) \Delta x_{l+1}] = E[f'(y_l)] E[\Delta x_{l+1}] \\ \doteq \frac{1}{2} E[\Delta x_{l+1}] = 0$$

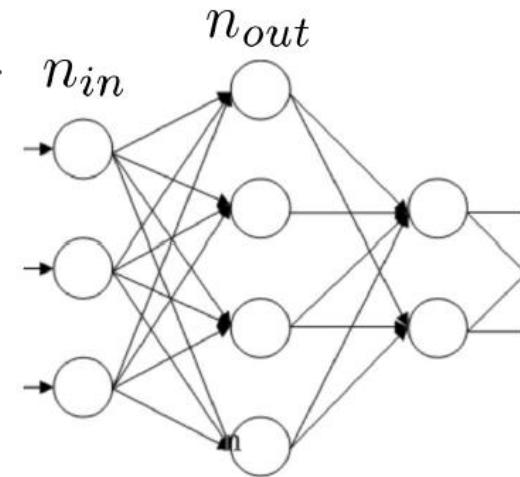
$$Var(w_l \cdot \Delta y_l) = [E(w_l)]^2 Var(\Delta y_l) + [E(\Delta y_l)]^2 Var(w_l) \\ + Var(w_l) Var(\Delta y_l)$$

$$\frac{1}{2} n_l Var[w_l] = 1$$

$$Var(W) = \frac{2}{n_l}$$

$$W \sim N(0, Var(W))$$

$$\underline{Var[\Delta y_l]} = \underline{Var[f'(y_l) \Delta x_{l+1}]} \\ = E[(f'(y_l))^2] E[(\Delta x_{l+1})^2] - [E[f'(y_l)]]^2 [E[\Delta x_{l+1}]]^2 \\ = E[(f'(y_l))^2] E[(\Delta x_{l+1})^2] - 0 \\ = E[(f'(y_l))^2] E[(\Delta x_{l+1})^2] \\ = \frac{1}{2} E[(\Delta x_{l+1})^2]$$



## 06 He initialization for ReLU activation function

[He2015] Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification.

$$Var(W) = \frac{2}{n_{in} + n_{out}}$$

$$W \sim N(0, Var(W))$$

*Xavier Normal*

$$Var(W) = \frac{6}{n_{in} + n_{out}}$$

$$W \sim U(0, Var(W))$$

*Xavier Uniform*

***Xavier initialization***

$$Var(W) = \frac{2}{n_{in}}$$

$$W \sim N(0, Var(W))$$

*He Normal*

$$Var(W) = \frac{6}{n_{in}}$$

$$W \sim U(0, Var(W))$$

*He Uniform*

***He initialization***

*Only Fan in is enough*



# 06 He initialization for ReLU activation function

[He2015] Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification.

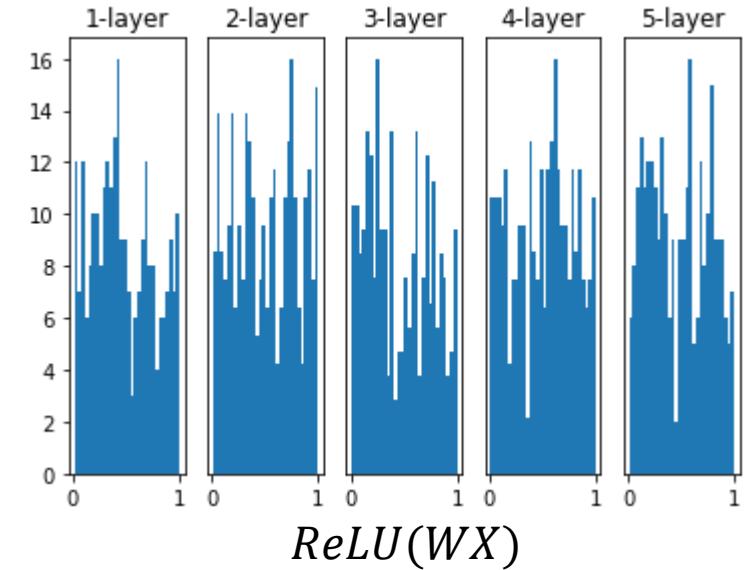
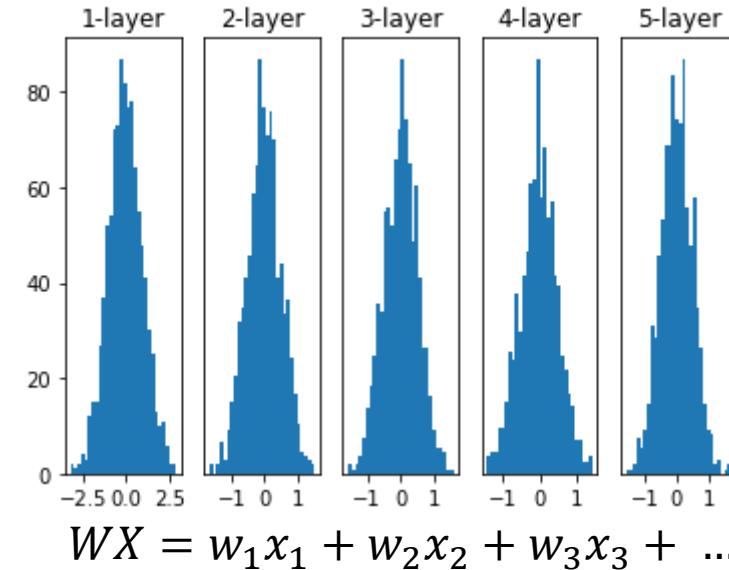
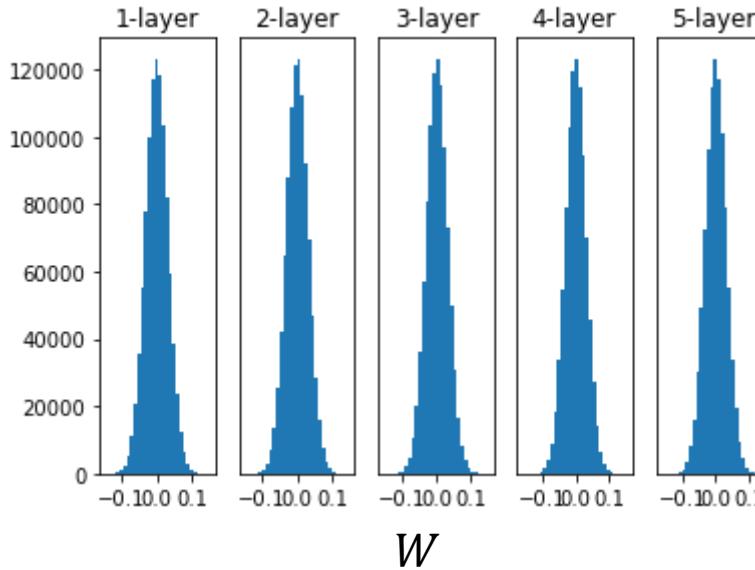
$$W \sim N(0, Var(W))$$

$$Var(W) = \frac{2}{n_{in}}$$

**Layer # : 5 (1000 node in each layer)**

X = np.random.randn(1, 1000)

W = np.random.randn(node\_num, node\_num) in each layer



## 06 He initialization for ReLU activation function

[He2015] Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification.

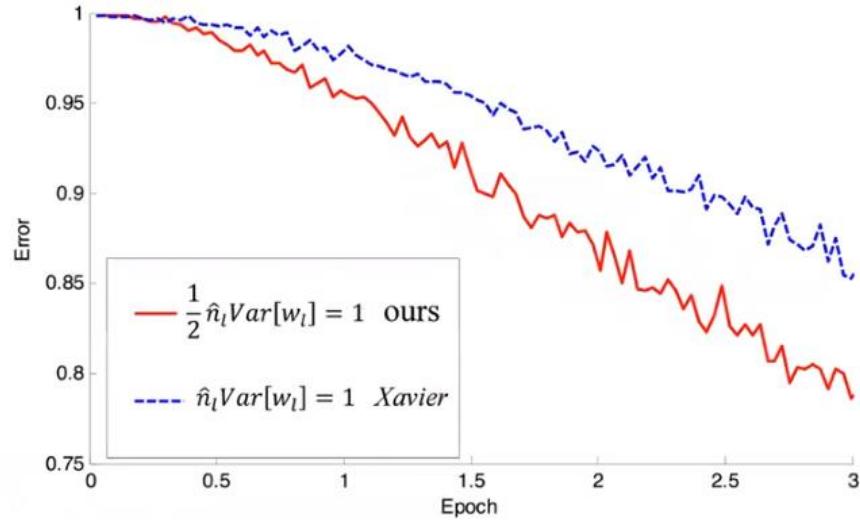


Figure 2. The convergence of a **22-layer** large model (B in Table 3). The x-axis is the number of training epochs. The y-axis is the top-1 error of 3,000 random val samples, evaluated on the center crop. We use ReLU as the activation for both cases. Both our initialization (red) and “Xavier” (blue) [7] lead to convergence, but ours starts reducing error earlier.

19 conv + 3 fc

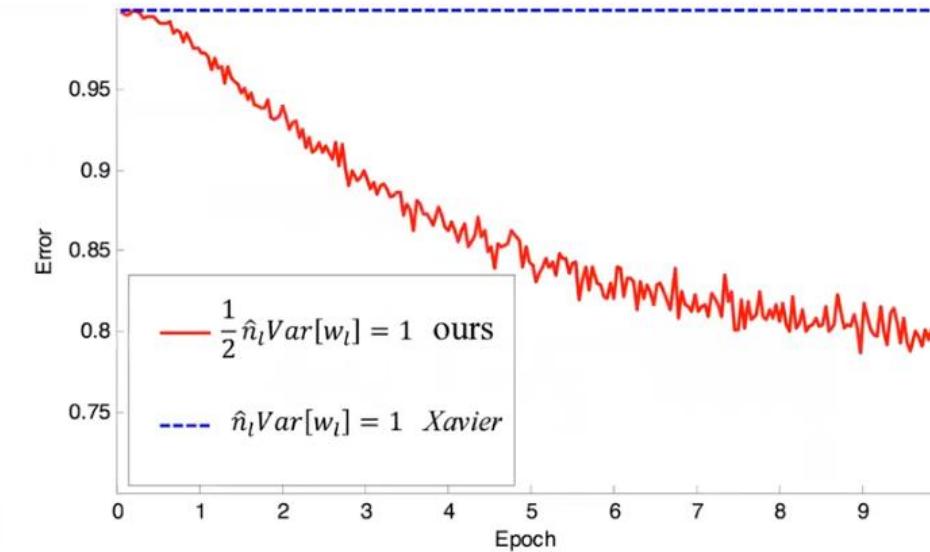


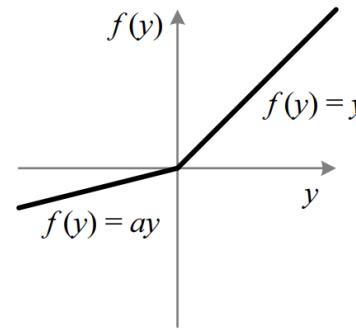
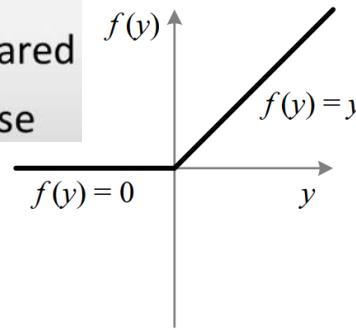
Figure 3. The convergence of a **30-layer** small model (see the main text). We use ReLU as the activation for both cases. Our initialization (red) is able to make it converge. But “Xavier” (blue) [7] completely stalls - we also verify that its gradients are all diminishing. It does not converge even given more epochs.

27 conv + 3 fc

# 06 He initialization for ReLU activation function

[He2015] Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification.

- per layer
- channel-shared
- channel-wise



- ReLU: when  $a_i = 0$

$$f(x_i) = \begin{cases} x_i, & \text{if } x_i > 0 \\ 0, & \text{if } x_i \leq 0 \end{cases}$$

- PReLU: when  $a_i$  is a *learnable parameter*

$$f(x_i) = \max(0, x_i) + a_i \min(0, x_i)$$

- LReLU: Leaky ReLU, when  $a_i = 0.01$

$$f(x_i) = \max(0, x_i) + 0.01 \min(0, x_i)$$

Figure 1. ReLU vs. PReLU. For PReLU, the coefficient of the negative part is not constant and is adaptively learned.

## Initialization in PReLU Case

- Forward Propagation Case:

$$\frac{1}{2}(1 + a^2)n_l \underline{\text{Var}[w_l]} = 1$$

- Backward Propagation Case:

$$\frac{1}{2}(1 + a^2)\hat{n}_l \underline{\text{Var}[w_l]} = 1$$

# 06 He initialization for ReLU activation function

[He2015] Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification.

## Architectures

### Models

- A: VGG-19 modified for speed
- B: deeper (+3 conv) than A
- C: wider (more filters) than B

input size	VGG-19 [25]	model A	model B	model C
224	$3 \times 3, 64$	$7 \times 7, 96, /2$	$7 \times 7, 96, /2$	$7 \times 7, 96, /2$
	$3 \times 3, 64$			
	$2 \times 2 \text{ maxpool}, /2$			
112	$3 \times 3, 128$			
	$3 \times 3, 128$	$2 \times 2 \text{ maxpool}, /2$	$2 \times 2 \text{ maxpool}, /2$	$2 \times 2 \text{ maxpool}, /2$
	$2 \times 2 \text{ maxpool}, /2$			
56	$3 \times 3, 256$	$3 \times 3, 256$	$3 \times 3, 256$	$3 \times 3, 384$
	$3 \times 3, 256$	$3 \times 3, 256$	$3 \times 3, 256$	$3 \times 3, 384$
	$3 \times 3, 256$	$3 \times 3, 256$	$3 \times 3, 256$	$3 \times 3, 384$
	$3 \times 3, 256$	$3 \times 3, 256$	$3 \times 3, 256$	$3 \times 3, 384$
	$3 \times 3, 256$	$3 \times 3, 256$	$3 \times 3, 256$	$3 \times 3, 384$
	$2 \times 2 \text{ maxpool}, /2$			
28	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 768$
	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 768$
	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 768$
	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 768$
	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 768$
	$2 \times 2 \text{ maxpool}, /2$			
14	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 896$
	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 896$
	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 896$
	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 896$
	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 512$	$3 \times 3, 896$
	$2 \times 2 \text{ maxpool}, /2$	$\text{spp}, \{7, 3, 2, 1\}$	$\text{spp}, \{7, 3, 2, 1\}$	$\text{spp}, \{7, 3, 2, 1\}$
		4096	4096	1000
fc <sub>1</sub>				
fc <sub>2</sub>				
fc <sub>3</sub>				
depth (conv+fc)	19	19	22	22
complexity (ops., $\times 10^{10}$ )	1.96	1.90	2.32	5.30

Table 3. Architectures of large models. Here “/2” denotes a stride of 2.

## 06 He initialization for ReLU activation function

[He2015] Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification.

# Experiments on ImageNet

Performed experiments on the 1000-class ImageNet 2012 dataset.

- 1.2M training images, 40K validation images, 100K test images
- Results measured by top-1/top-5 error rates
- **Top-5 error rate** is the metric officially used to rank the methods in the classification challenge.

## 06 He initialization for ReLU activation function

[He2015] Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification.

### Comparisons between ReLU and PReLU

PReLU reduces top-1 error by 1.05% and top-5 error by 0.23%.

- Large model A
- Channel-wise PReLU
- PReLU improves both small and large models consistently

model A	ReLU		PReLU	
	scale $s$	top-1	top-5	top-1
256	26.25	8.25	<b>25.81</b>	<b>8.08</b>
384	24.77	7.26	<b>24.20</b>	<b>7.03</b>
480	25.46	7.63	<b>24.83</b>	<b>7.39</b>
multi-scale	24.02	6.51	<b>22.97</b>	<b>6.28</b>

Table 4. Comparisons between ReLU/PReLU on model A in ImageNet 2012 using dense testing.

## 06 He initialization for ReLU activation function

[He2015] Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification.

	team	top-1	top-5
in competition ILSVRC 14	MSRA [11]	27.86	9.08 <sup>†</sup>
	VGG [25]	-	8.43 <sup>†</sup>
	GoogLeNet [29]	-	7.89
post-competition	VGG [25] (arXiv v2)	24.8	7.5
	VGG [25] (arXiv v5)	24.4	7.1
	Baidu [32]	24.88	7.42
	MSRA (A, ReLU)	24.02	6.51
	MSRA (A, PReLU)	22.97	6.28
	MSRA (B, PReLU)	22.85	6.27
	MSRA (C, PReLU)	<b>21.59</b>	<b>5.71</b>

Table 6. The **single-model** results for ImageNet 2012 val set. <sup>†</sup>: Evaluated from the test set.

	team	top-5 (test)
in competition ILSVRC 14	MSRA, SPP-nets [11]	8.06
	VGG [25]	7.32
	GoogLeNet [29]	6.66
post-competition	VGG [25] (arXiv v5)	6.8
	Baidu [32]	5.98
	<b>MSRA, PReLU-nets</b>	<b>4.94</b>

Table 7. The **multi-model** results for the ImageNet 2012 test set.

O. Russakovsky et al.,  
“Imagenet large scale visual  
recognition challenge”, 2014.

Human performance yields  
**5.1% top-5 error.**

“Our result is the first published  
instance of surpassing humans on  
this visual recognition challenge.”

# 07 Conclusion

## Conclusion

- ✓ It is very important to initialize weights properly in deep network.
- ✓ Don't initialize all weights to zero or constant value.
- ✓ Don't initialize all weights randomly to small value or large value.
- ✓ Xavier initialization is good for linear activation function (Sigmoid, TanH, etc.).
- ✓ He initialization is good for ReLU.
- ✓ [Kumar2017] On weight initialization in deep neural networks.



# Thank you !

E-Mail : [yangdonghun3@gmail.com](mailto:yangdonghun3@gmail.com)  
Phone : +82-10-2484-5894