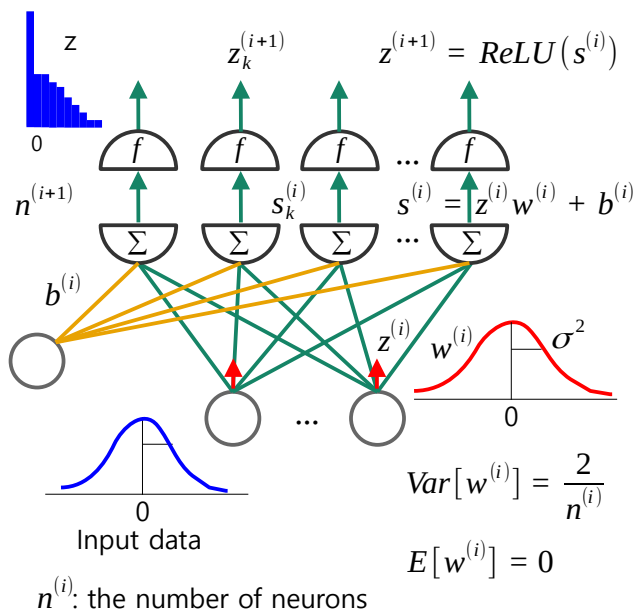




## 8. Weights Initialization

### Part 1: Observation of the distribution of hidden layer outputs



This video was produced in Korean and translated into English, and the audio was generated by AI (TTS).

[www.youtube.com/@meanxai](http://www.youtube.com/@meanxai)

[MXDL-8-01]

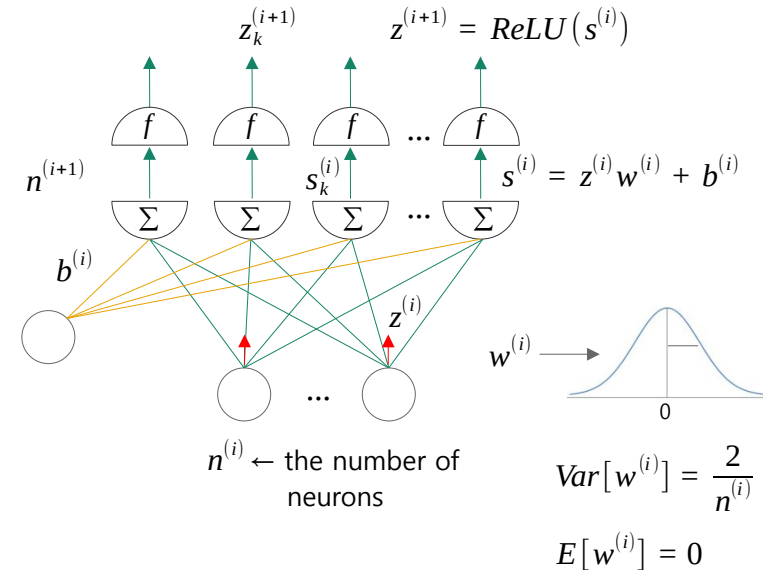
1. Parameter space and error
2. Observation of the distribution of hidden layer outputs according to initial weights
3. Formula for the variance of output

[MXDL-8-02]

4. Xavier Glorot initialization
  - Paper
  - Forward direction
  - Backward direction
  - Uniform distribution
  - Xavier Glorot initializer in Keras

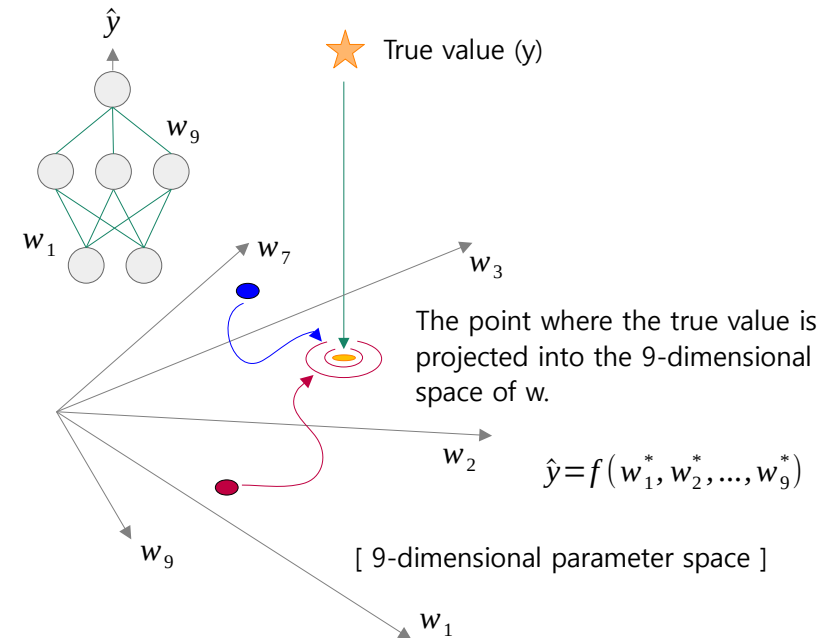
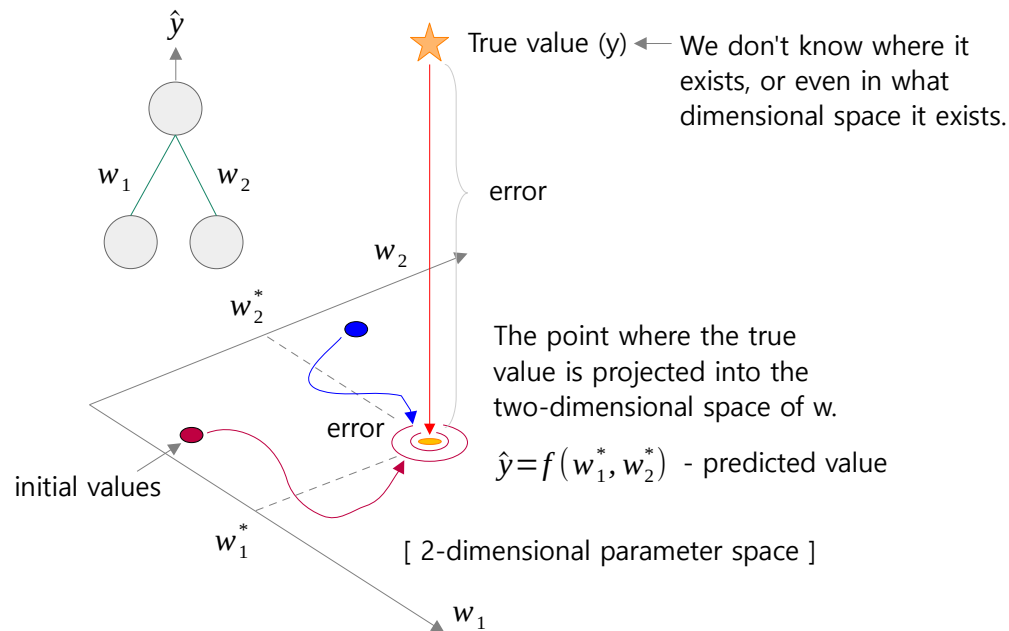
[MXDL-8-03]

5. Kaiming He initialization
  - Paper
  - Forward direction
  - Backward direction
  - Uniform distribution
  - Kaiming He initializer in Keras



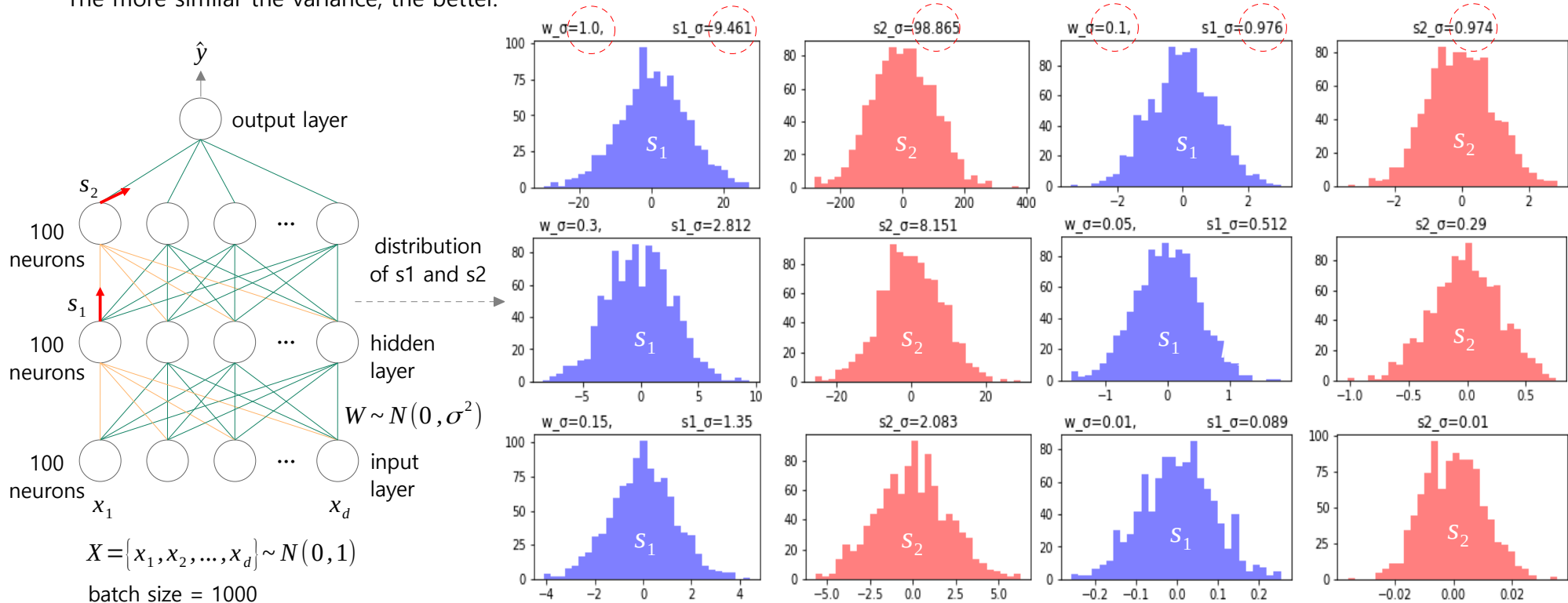
## ■ Parameter space and error

- Properly setting the number and initial values of parameters is a critical step in training a neural network.
- Since we don't know where the true values  $y$  are, we don't know how wide or deep to build a model. The number of parameters  $w$  is determined by the width and depth of a model. Models with a small number of parameters may have large errors and underfitting, while models with a large number of parameters may have small errors but overfitting.
- Additionally, the initial values of parameters also affect the performance of models. Depending on the initial values, vanishing or exploding gradients problems may occur, and local minimum problems may also increase.



## ■ Distribution of hidden layer outputs according to initial weights

- The distribution of the input data  $X$  is a normal distribution with mean 0 and variance 1, and the distribution of the initial weights in the hidden layer is a normal distribution with mean 0 and variance  $\sigma^2$ . Let's observe the distribution of hidden layer outputs as the  $\sigma$  of initial  $w$  changes.
- The distribution of a hidden layer output values changes depending on the initial values of  $w$ .
- When  $w_\sigma$  is 1.0, the variances of  $s_1$  and  $s_2$  are significantly different. When  $w_\sigma$  is 0.1, the variances of  $s_1$  and  $s_2$  are similar.
- The more similar the variance, the better.



## ■ Observation of the distribution of hidden layer outputs according to initial weights

```
# [MXDL-8-01] 1.h_distribution.py
# Observe the distribution of hidden layer outputs according
# to initial weights
```

```
import numpy as np
from tensorflow.keras.layers import Input, Dense
from tensorflow.keras import initializers
from tensorflow.keras.models import Model
import matplotlib.pyplot as plt
```

```
# Generate a simple dataset with 1000 data points
```

```
x = np.random.normal(size=(1000, 100))
```

```
# Create an ANN model with 2 hidden layers
```

```
n_input = x.shape[1]
```

```
n_hidden = 100
```

```
# Let's check the distribution of hidden layer
```

```
# outputs by changing std in  $N(0, \sigma^2)$ .
```

```
std = [1.0, 0.3, 0.15, 0.1, 0.05, 0.01]
```

```
for sigma in std:
```

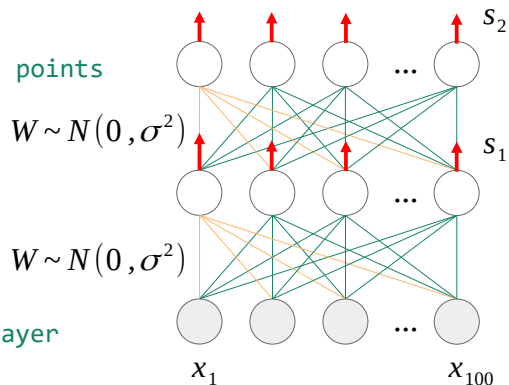
```
    # Create an ANN model
```

```
    w0 = initializers.RandomNormal(mean=0.0, stddev=sigma)
```

```
    x_input = Input(batch_shape=(None, n_input))
```

```
    s1 = Dense(n_hidden, kernel_initializer = w0)(x_input)
```

```
    s2 = Dense(n_hidden, kernel_initializer = w0)(s1)
```



$$W \sim N(0, \sigma^2)$$

$$W \sim N(0, \sigma^2)$$

$$X = \{x_1, x_2, \dots, x_d\} \sim N(0, 1)$$

```
s1_model = Model(x_input, s1)
```

```
s2_model = Model(x_input, s2)
```

```
i = 0
```

```
s1_out = s1_model.predict(x, verbose=0)[: , i]
```

```
s2_out = s2_model.predict(x, verbose=0)[: , i]
```

```
# Check the distribution of the hidden layer outputs.
```

```
plt.figure(figsize=(8,2))
```

```
plt.subplot(121)
```

```
plt.hist(s1_out, bins=30, color='blue', alpha=0.5)
```

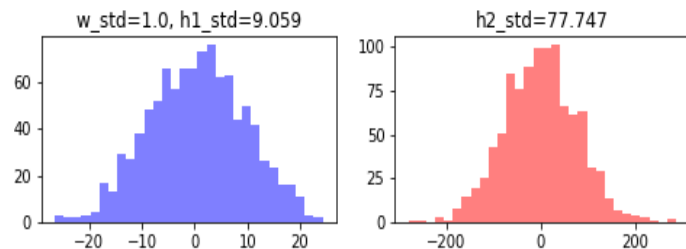
```
plt.title('w_std=' + str(sigma) + ', \
          s1_std=' + str(s1_out.std().round(3)))
```

```
plt.subplot(122)
```

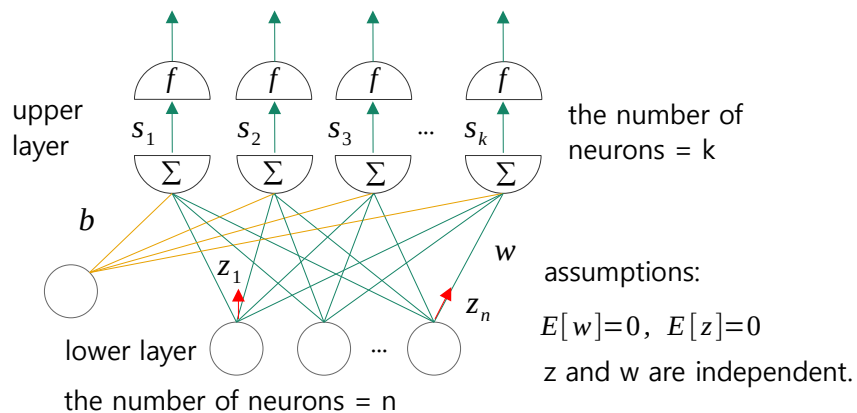
```
plt.hist(s2_out, bins=30, color='red', alpha=0.5)
```

```
plt.title('s2_std=' + str(s2_out.std().round(3)))
```

```
plt.show()
```



## ■ Formula for the variance of output



$$S = Z \cdot W + B$$

$$S = \begin{pmatrix} z_{1,1} & z_{1,2} & \dots & z_{1,n} \\ z_{2,1} & z_{2,2} & \dots & z_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ z_{m,1} & z_{m,2} & \dots & z_{m,n} \end{pmatrix} \cdot \begin{pmatrix} w_{1,1} & w_{1,2} & \dots & w_{1,k} \\ w_{2,1} & w_{2,2} & \dots & w_{2,k} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n,1} & w_{n,2} & \dots & w_{n,k} \end{pmatrix} + [b_1, b_2, b_3, \dots, b_k]$$

(m, n)                      (n, k)

(m: batch size)

$$S = \begin{pmatrix} \sum_{j=1}^n z_{1,j} w_{j,1} + b_1 & \sum_{j=1}^n z_{1,j} w_{j,2} + b_2 & \dots & \sum_{j=1}^n z_{1,j} w_{j,k} + b_k \\ \sum_{j=1}^n z_{2,j} w_{j,1} + b_1 & \sum_{j=1}^n z_{2,j} w_{j,2} + b_2 & \dots & \sum_{j=1}^n z_{2,j} w_{j,k} + b_k \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{j=1}^n z_{m,j} w_{j,1} + b_1 & \sum_{j=1}^n z_{m,j} w_{j,2} + b_2 & \dots & \sum_{j=1}^n z_{m,j} w_{j,k} + b_k \end{pmatrix} \quad (m, k)$$

$s_1$                        $s_2$                        $s_k$

$$\text{Var}[s_1] \downarrow \quad \text{Let: } z_i = z_{1:m,i} \quad (b = \text{constant})$$

$$\begin{aligned} \text{Var}\left[\sum_{j=1}^n z_j w_{j,1}\right] &= \text{Var}[z_1 w_{1,1} + z_2 w_{2,1} + \dots + z_n w_{n,1}] \\ &= \text{Var}[z_1 w_{1,1}] + \text{Var}[z_2 w_{2,1}] + \dots + \text{Var}[z_n w_{n,1}] \end{aligned}$$

Assumptions:  $s, z$  and  $w$  are i.i.d.

$$\text{Var}[z_1 w_{1,1}] \approx \text{Var}[z_2 w_{2,1}] \approx \dots \approx \text{Var}[z_n w_{n,1}]$$

$$\text{Var}[s_1] \approx \text{Var}[s_2] \approx \dots \approx \text{Var}[s_k]$$

$$\text{Var}(s_1) = \text{Var}\left[\sum_{j=1}^n z_{1:m,j} w_{j,1}\right] = n \cdot \text{Var}[z_1 w_{1,1}] = n \cdot \text{Var}[z_2 w_{2,1}] = \dots$$

$$\text{Var}[S] = n \cdot \text{Var}[Z \cdot W]$$

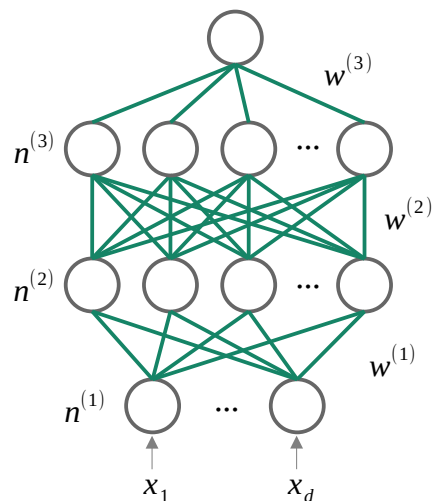
\* Reference: <https://www.pinecone.io/learn/weight-initialization/>  
<https://towardsdatascience.com/xavier-glorot-initialization-in-neural-networks-math-proof-4682bf5c6ec3>



$$w^{(i)} \sim N\left(0, \frac{2}{n^{(i)} + n^{(i+1)}}\right)$$
$$w^{(i)} \sim U\left(-\sqrt{\frac{6}{n^{(i)} + n^{(i+1)}}}, \sqrt{\frac{6}{n^{(i)} + n^{(i+1)}}}\right)$$

## 8. Weights Initialization

### Part 2: Xavier Glorot initializer



$$x = \{x_1, x_2, \dots, x_{50}\} \sim N(0, 1)$$

This video was produced in Korean and translated into English,  
and the audio was generated by AI (TTS).

[www.youtube.com/@meanxai](http://www.youtube.com/@meanxai)

- Xavier Glorot paper (2010)

## Understanding the difficulty of training deep feedforward neural networks

Xavier Glorot

Yoshua Bengio

DIRO, Universite de Montr al, Montr al, Qu bec, Canada

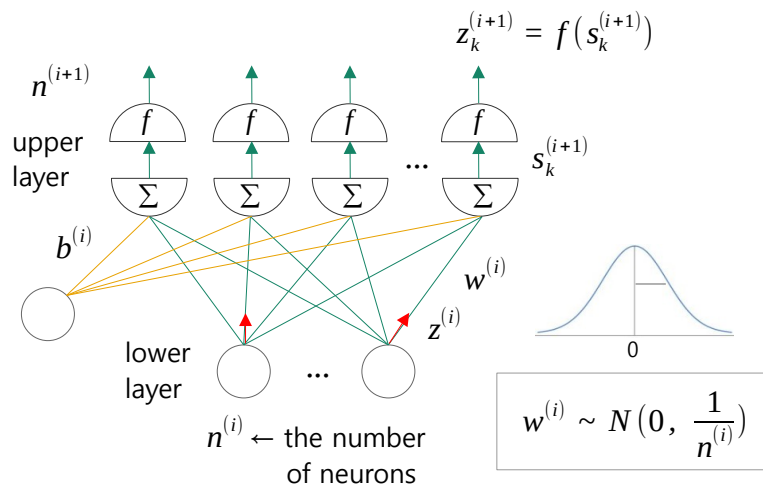
### Abstract

Whereas before 2006 it appears that deep multilayer neural networks were not successfully trained, since then several algorithms have been shown to successfully train them, with experimental results showing the superiority of deeper vs less deep architectures. All these experimental results were obtained with new initialization or training mechanisms. Our objective here is to understand better why standard gradient descent from random initialization is doing so poorly with deep neural networks, to better understand these recent relative successes and help design better algorithms in the future. We first observe the influence of the non-linear activations functions. We find that the logistic sigmoid activation is unsuited for deep networks with random initialization because of its mean value, which can drive especially the top hidden layer into saturation. Surprisingly, we find that saturated units can move out of saturation by themselves, albeit slowly, and explaining the plateaus sometimes seen when training neural networks. We find that a new non-linearity that saturates less can often be beneficial. Finally, we study how activations and gradients vary across layers and during training, with the idea that training may be more difficult when the singular values of the Jacobian associated with each layer are far from 1. Based on these considerations, we propose a **new initialization scheme** that brings substantially faster convergence.



## ■ Xavier Glorot initializer: Forward direction

- During forward propagation, we would like the variance of the activations to remain constant across layers.
- Let's find the optimal variance of the initial weight distribution in the forward direction.



Assumption:  $E[z^{(i)}] = E[w^{(i)}] = 0$

Objective:  $\forall i, \text{Var}[z^{(i)}] = \text{Var}[z^{(i+1)}]$  ← We want the variance of the activations to remain constant across layers.

$\text{Var}[z^{(i+1)}] = \text{Var}[f(s^{(i+1)})] \approx \text{Var}[s^{(i+1)}]$  ← Assume  $f$  is linear near 0.  
(ex: sigmoid, tanh)

$= n^{(i)} \text{Var}[z^{(i)} w^{(i)}]$  ← by the formula for the variance of output

$= n^{(i)} [\text{Var}[z^{(i)}] \text{Var}[w^{(i)}] + \text{Var}[z^{(i)}] (E[w^{(i)}])^2 + \text{Var}[w^{(i)}] (E[z^{(i)}])^2]$

$= n^{(i)} \text{Var}[z^{(i)}] \text{Var}[w^{(i)}]$

$\text{Var}[z^{(i+1)}] = n^{(i)} \text{Var}[z^{(i)}] \cdot \text{Var}[w^{(i)}]$

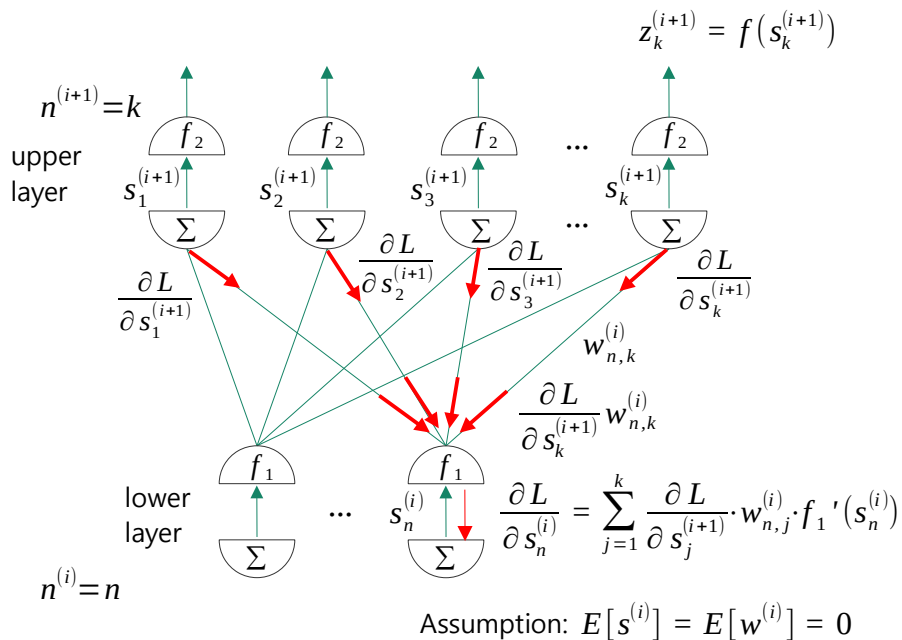
$\text{Var}[w^{(i)}] = \frac{1}{n^{(i)}}$

\* Reference: <https://www.pinecone.io/learn/weight-initialization/>  
<https://towardsdatascience.com/xavier-glorot-initialization-in-neural-networks-math-proof-4682bf5c6ec3>

For example, if the number of neurons in the lower layer is 100, the initial weights of the upper layer are set by drawing them from a normal distribution with a mean of 0 and a variance of 0.01.  $N(0, 1/100)$

## ■ Xavier Glorot initializer: Backward direction

- During backpropagation, we would like the variance of the gradients of loss with respect to the output  $s$  to remain constant across layers.



Objective:

$$\forall i, \text{Var}\left[\frac{\partial L}{\partial s^{(i+1)}}\right] = \text{Var}\left[\frac{\partial L}{\partial s^{(i)}}\right] \quad \text{We want the variance of the gradients to remain constant across layers.}$$

$$\frac{\partial L}{\partial s_n^{(i)}} = \sum_{j=1}^k \frac{\partial L}{\partial s_j^{(i+1)}} \cdot w_{n,j}^{(i)} \cdot f_1'(s_n^{(i)}) \quad \leftarrow \text{Assume } f_1 \text{ is linear near 0. (ex: sigmoid, tanh)} \\ f_1(x)=x, f_1'(x)=1$$

$$\text{Var}\left[\frac{\partial L}{\partial s_n^{(i)}}\right] = \text{Var}\left[\sum_{j=1}^k \frac{\partial L}{\partial s_j^{(i+1)}} \cdot w_{n,j}^{(i)}\right]$$

$$\text{Var}\left[\frac{\partial L}{\partial s^{(i)}}\right] = n^{(i+1)} \text{Var}\left[\frac{\partial L}{\partial s^{(i+1)}} \cdot w^{(i)}\right]$$

$$= n^{(i+1)} \left[ \text{Var}\left[\frac{\partial L}{\partial s^{(i+1)}}\right] \text{Var}[w^{(i)}] + \text{Var}\left[\frac{\partial L}{\partial s^{(i+1)}}\right] (E[w^{(i)}])^2 + \text{Var}[w^{(i)}] (E\left[\frac{\partial L}{\partial s^{(i+1)}}\right])^2 \right]$$

0                      0

(assumption)

$$\text{Var}\left[\frac{\partial L}{\partial s^{(i)}}\right] = n^{(i+1)} \text{Var}\left[\frac{\partial L}{\partial s^{(i+1)}}\right] \cdot \text{Var}[w^{(i)}]$$

$$\text{Var}[w^{(i)}] = \frac{1}{n^{(i+1)}}$$

- Find the average of the forward and backward directions.

$$\text{Backward direction} \rightarrow n^{(i+1)} \text{Var}[w^{(i)}] = 1$$

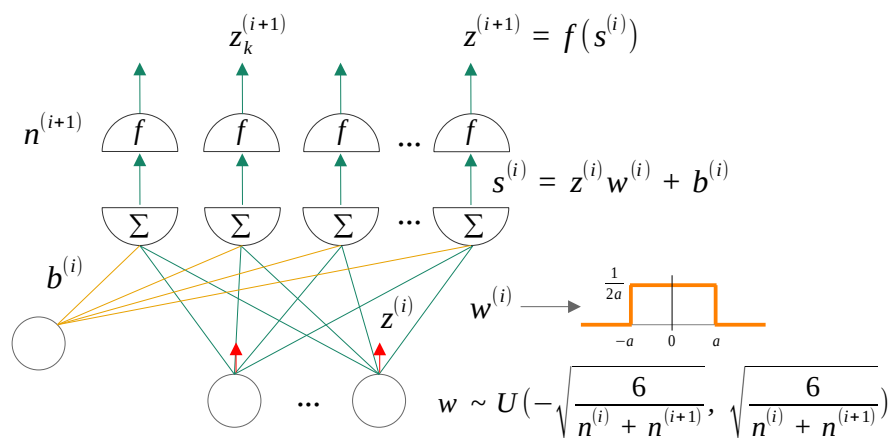
$$\text{Forward direction} \rightarrow n^{(i)} \text{Var}[w^{(i)}] = 1$$

$$\text{Var}[w^{(i)}] = \frac{2}{n^{(i)} + n^{(i+1)}}$$

For example, if the number of neurons in the lower and upper layer is 100 each, the initial weights of the upper layer are set by drawing them from a normal distribution with a mean of 0 and a variance of 0.01.  $N(0, 0.01)$

## ■ Xavier Glorot initializer: Uniform distribution

- Now that we know the optimal variance for the normal distribution, we can also find the range for the uniform distribution.



$n^{(i)}$  ← the number of neurons

assumption:  $E[z^{(i)}] = E[w^{(i)}] = 0$

- continuous uniform distribution

$$g(x) = \begin{cases} \frac{1}{2a}, & \text{if } -a \leq x \leq a \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Var}[x] = \int_{-a}^a x^2 \cdot \frac{1}{2a} dx = \left[ \frac{1}{3} x^3 \frac{1}{2a} \right]_{-a}^a = \frac{1}{3} a^2$$

- Forward and backward directions

$$(1) \forall i, \text{Var}[z^{(i)}] = \text{Var}[z^{(i+1)}] \quad \leftarrow \text{Forward direction}$$

$$(2) \forall i, \text{Var}\left[\frac{\partial L}{\partial s^{(i)}}\right] = \text{Var}\left[\frac{\partial L}{\partial s^{(i+1)}}\right] \quad \leftarrow \text{Backward direction}$$

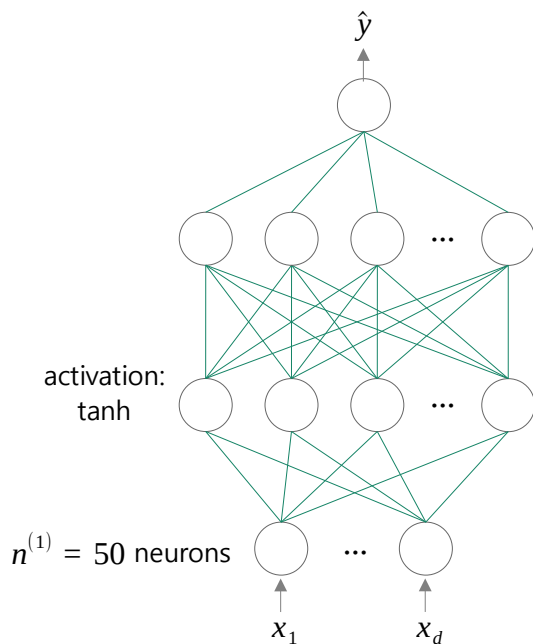
- We already know the variance that satisfies the above conditions, and the variance of the uniform distribution must also be equal to that variance.

$$\frac{1}{3} a^2 = \frac{2}{n^{(i)} + n^{(i+1)}} \longrightarrow a = \sqrt{\frac{6}{n^{(i)} + n^{(i+1)}}}$$

\* reference: <https://towardsdatascience.com/xavier-glorot-initialization-in-neural-networks-math-proof-4682bf5c6ec3>

## ■ Xavier Glorot initializer: Example

- Given an ANN network like the one below, set appropriate initial distributions of  $w^{(1)}$  and  $w^{(2)}$ .



### ■ Normal distribution

$$w \sim N(0, \sigma^2) = N\left(0, \frac{2}{n^{(i)} + n^{(i+1)}}\right)$$

$n^{(3)} = 100$  neurons

$$w^{(2)} \sim N\left(0, \frac{2}{80 + 100}\right) = N(0, 0.105^2)$$

$n^{(2)} = 80$  neurons

$$w^{(1)} \sim N\left(0, \frac{2}{50 + 80}\right) = N(0, 0.124^2)$$

### ■ Uniform distribution

$$w \sim U(-a, a) = U\left(-\sqrt{\frac{6}{n^{(i)} + n^{(i+1)}}}, \sqrt{\frac{6}{n^{(i)} + n^{(i+1)}}}\right)$$

$$w^{(2)} \sim U\left(-\sqrt{\frac{6}{80+100}}, \sqrt{\frac{6}{80+100}}\right) = U(-0.183, 0.183)$$

$$w^{(1)} \sim U\left(-\sqrt{\frac{6}{50+80}}, \sqrt{\frac{6}{50+80}}\right) = U(-0.215, 0.215)$$

## ■ Keras Xavier Glorot initializer

```
# [MXDL-8-02] 2.xavier_glorot.py
import numpy as np
from tensorflow.keras.layers import Input, Dense
from tensorflow.keras import initializers
from tensorflow.keras.models import Model
import matplotlib.pyplot as plt
```

```
method = 'Normal'
if method == 'Normal':
    init_w1 = initializers.GlorotNormal()
    init_w2 = initializers.GlorotNormal()
else:
    init_w1 = initializers.GlorotUniform()
    init_w2 = initializers.GlorotUniform()
```

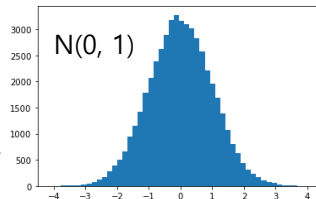
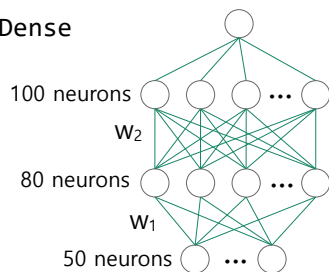
```
n_in = 50
n_s1 = 80
n_s2 = 100
x = np.random.normal(size=(1000, n_in))
```

## # Create an ANN model

```
x_input = Input(batch_shape=(None, n_in))
s1 = Dense(n_s1, kernel_initializer=init_w1, activation='tanh',
           name='W1')(x_input)
s2 = Dense(n_s2, kernel_initializer=init_w2, activation='tanh',
           name='W2')(s1)
y_output = Dense(1, activation='sigmoid')(s2)
s1_model = Model(x_input, s1)
s2_model = Model(x_input, s2)
model = Model(x_input, y_output)
```

```
w1 = model.get_layer('W1').get_weights()[0].reshape(-1,)
w2 = model.get_layer('W2').get_weights()[0].reshape(-1,)
```

```
s1_out = s1_model.predict(x, verbose=0).reshape(-1,)
s2_out = s2_model.predict(x, verbose=0).reshape(-1,)
```



```
plt.hist(x.reshape(-1,), bins=50); plt.show()
plt.hist(s1_out, bins=50); plt.show()
plt.hist(s2_out, bins=50); plt.show()
```

```
if method == 'Normal':
    print('[Keras ] σ of w1 = {:.3f}'.format(w1.std()))
    print('[Formula] σ of w1 = {:.3f}'.\
          format(np.sqrt(2 / (n_in + n_s1))))
else:
    print('[Keras ] a of w1 = {:.3f} ~ {:.3f}'.\
          format(w1.min(), w1.max()))
    print('[Formula] a of w1 = ±{:.3f}'.\
          format(np.sqrt(6 / (n_in + n_s1))))

if method == 'Normal':
    print('\n[Keras ] σ of w2 = {:.3f}'.format(w2.std()))
    print('[Formula] σ of w2 = {:.3f}'.\
          format(np.sqrt(2 / (n_s1 + n_s2))))
else:
    print('\n[Keras] a of w2 = {:.3f} ~ {:.3f}'.\
          format(w2.min(), w2.max()))
    print('[Formula] a of w2 = ±{:.3f}'.\
          format(np.sqrt(6 / (n_s1 + n_s2))))
```

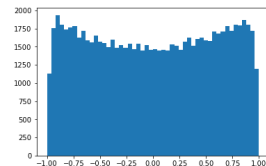
```
method = 'Normal':
[Keras ] σ of w1 = 0.125
[Formula] σ of w1 = 0.124
```

```
[Keras ] σ of w2 = 0.105
[Formula] σ of w2 = 0.105
```

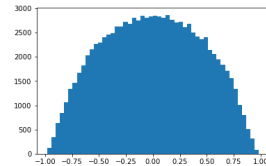
```
method = 'Uniform':
[Keras ] a of w1 = -0.215 ~ 0.214
[Formula] a of w1 = ±0.215
```

```
[Keras] a of w2 = -0.182 ~ 0.183
[Formula] a of w2 = ±0.183
```

s1\_out →



s2\_out →



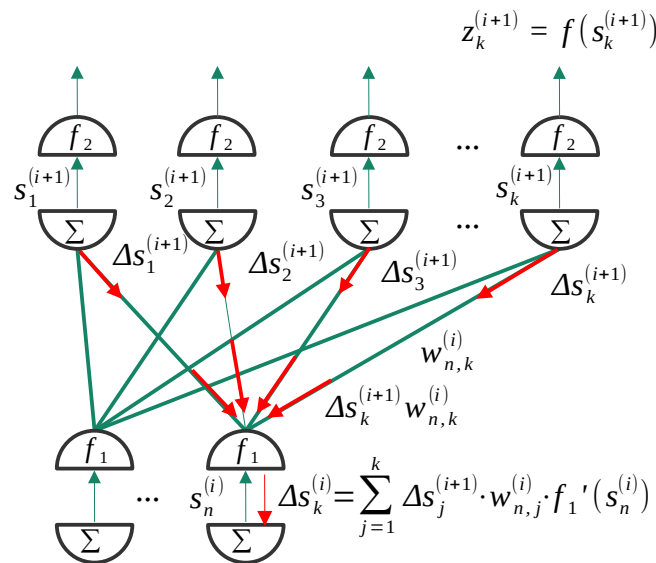


$$w^{(i)} \sim N(0, \frac{2}{n^{(i)}}), \text{ or } N(0, \frac{2}{n^{(i+1)}})$$

$$w^{(i)} \sim U(-\sqrt{\frac{6}{n^{(i)}}}, \sqrt{\frac{6}{n^{(i)}}}), \text{ or } U(-\sqrt{\frac{6}{n^{(i+1)}}}, \sqrt{\frac{6}{n^{(i+1)}}})$$

## 8. Weights Initialization

### Part 3: Kaiming He initializer



This video was produced in Korean and translated into English, and the audio was generated by AI (TTS).

[www.youtube.com/@meanxai](http://www.youtube.com/@meanxai)

- Kaiming He paper (2015)

## Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification

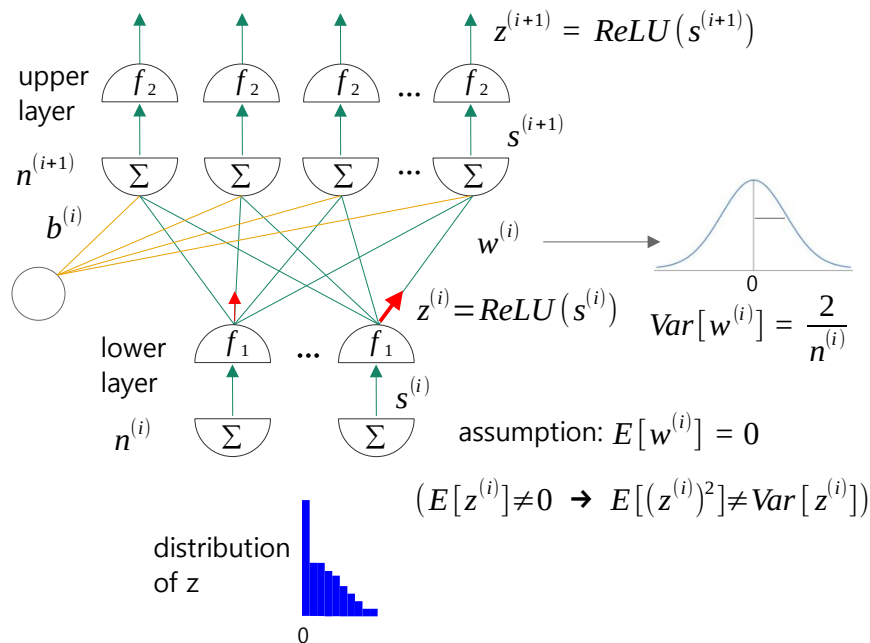
Kaiming He Xiangyu Zhang Shaoqing Ren Jian Sun  
Microsoft Research  
{kahe, v-xiangz, v-shren, jiansun}@microsoft.com

### Abstract

Rectified activation units (rectifiers) are essential for state-of-the-art neural networks. In this work, we study rectifier neural networks for image classification from two aspects. **First**, we propose a Parametric Rectified Linear Unit (PReLU) that generalizes the traditional rectified unit. PReLU improves model fitting with nearly zero extra computational cost and little overfitting risk. **Second**, we derive a robust initialization method that particularly considers the rectifier nonlinearities. This method enables us to train extremely deep rectified models directly from scratch and to investigate deeper or wider network architectures. Based on our PReLU networks (PReLU-nets), we achieve 4.94% top-5 test error on the ImageNet 2012 classification dataset. This is a 26% relative improvement over the ILSVRC 2014 winner (GoogLeNet, 6.66% [29]). To our knowledge, our result is the first to surpass human-level performance (5.1%, [22]) on this visual recognition challenge.

## ■ Kaiming He initializer: Forward direction

- During forward propagation, we would like the variance of the activations to remain constant across layers.
- Let's find the optimal variance of the initial weight distribution in the forward direction.



Objective:  $\forall i, \text{Var}[s^{(i+1)}] = \text{Var}[s^{(i)}]$  ← We want the variance of the activations to remain constant across layers.

$$\text{Var}[s^{(i+1)}] = n^{(i)} \text{Var}[z^{(i)} w^{(i)}]$$

$$= n^{(i)} [\text{Var}[z^{(i)}] \text{Var}[w^{(i)}] + \text{Var}[z^{(i)}] (E[w^{(i)}])^2 + \text{Var}[w^{(i)}] (E[z^{(i)}])^2]$$

$$= n^{(i)} [\text{Var}[z^{(i)}] \text{Var}[w^{(i)}] + \text{Var}[w^{(i)}] (E[z^{(i)}])^2]$$

$$= n^{(i)} [\text{Var}[z^{(i)}] \text{Var}[w^{(i)}] + \text{Var}[w^{(i)}] (E[(z^{(i)})^2] - \text{Var}[z^{(i)}])]$$

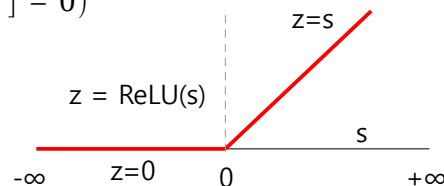
$$= n^{(i)} \text{Var}[w^{(i)}] E[(z^{(i)})^2]$$

$$E[(z^{(i)})^2] = \int_{-\infty}^{\infty} \text{ReLU}(s^{(i)})^2 p(s) ds = \int_0^{\infty} (s^{(i)})^2 p(s) ds = \frac{1}{2} \int_{-\infty}^{\infty} (s^{(i)})^2 p(s) ds$$

$$= \frac{1}{2} E[(s^{(i)})^2] = \frac{1}{2} \text{Var}[s^{(i)}] \quad (E[s^{(i)}] = 0)$$

$$\text{Var}[s^{(i+1)}] = \frac{1}{2} n^{(i)} \text{Var}[w^{(i)}] \text{Var}[s^{(i)}]$$

$$\text{Var}[w^{(i)}] = \frac{2}{n^{(i)}}$$

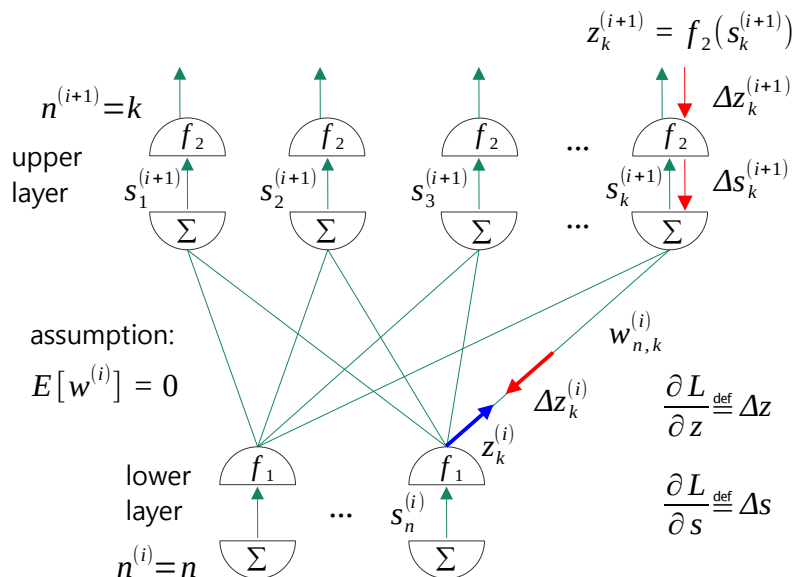


For example, if the number of neurons in the lower layer is 100, the initial weights of the upper layer are set by drawing them from a normal distribution with a mean of 0 and a variance of 0.02.  $N(0, 2/100)$



## ■ Kaiming He initializer: Backward direction

- During backpropagation, we would like the variance of the gradients of loss with respect to the output  $z$  to remain constant across layers.



For example, if the number of neurons in the lower layer is 80 and the number of neurons in the upper layer is 100, the initial weights between the lower and upper layer are set by drawing them from a normal distribution with a mean of 0, and a variance of either 2/100 or 2/80.

$$\text{Objective: } \forall i, \text{Var}\left[\frac{\partial L}{\partial z^{(i+1)}}\right] = \text{Var}\left[\frac{\partial L}{\partial z^{(i)}}\right] \rightarrow \text{Var}[\Delta z^{(i+1)}] = \text{Var}[\Delta z^{(i)}]$$

$$\Delta z^{(i)} = \frac{\partial L}{\partial s^{(i+1)}} \frac{\partial s^{(i+1)}}{\partial z^{(i)}} = \Delta s^{(i+1)} w^{(i)} \quad (s^{(i+1)} = z^{(i)} \cdot w^{(i)} + b^{(i)})$$

$$E[\Delta z^{(i)}] = E[\Delta s^{(i+1)} w^{(i)}] = E[\Delta s^{(i+1)}] E[w^{(i)}] = 0 \quad (\Delta s \text{ and } w \text{ are independent, } E[w] = 0)$$

$$\Delta s^{(i+1)} = \frac{\partial L}{\partial s^{(i+1)}} = \frac{\partial L}{\partial z^{(i+1)}} \frac{\partial z^{(i+1)}}{\partial s^{(i+1)}} = \Delta z^{(i+1)} f_2'(s^{(i+1)}) \quad \left[ \begin{array}{l} f'(s) = 1, \text{ if } s > 0 \\ f'(s) = 0, \text{ if } s < 0 \end{array} \right] \text{ with equal probabilities}$$

$$E[\Delta s^{(i+1)}] = E[\Delta z^{(i+1)} f_2'(s^{(i+1)})] = \frac{1}{2} E[\Delta z^{(i+1)} \times 1] + \frac{1}{2} E[\Delta z^{(i+1)} \times 0] = 0$$

$$\begin{aligned} \text{Var}[\Delta s^{(i+1)}] &= \text{Var}[\Delta z^{(i+1)} f_2'(s^{(i+1)})] \\ &= E[(\Delta z^{(i+1)})^2] E[f_2'(s^{(i+1)})^2] - E[\Delta z^{(i+1)}]^2 E[f_2'(s^{(i+1)})]^2 \\ &= \frac{1}{2} \text{Var}[\Delta z^{(i+1)}] \quad (E[f_2'(s^{(i+1)})^2] = \frac{0^2}{2} + \frac{1^2}{2} = \frac{1}{2}) \end{aligned}$$

$$\begin{aligned} \text{Var}[\Delta z^{(i)}] &= n^{(i+1)} \text{Var}[\Delta s^{(i+1)} w^{(i)}] \quad \leftarrow \Delta z^{(i)} = \sum_{i=1}^k \Delta s_i^{(i+1)} w^{(i)} \\ &= n^{(i+1)} \text{Var}[\Delta s^{(i+1)}] \text{Var}[w^{(i)}] = n^{(i+1)} \frac{1}{2} \text{Var}[\Delta z^{(i+1)}] \text{Var}[w^{(i)}] \end{aligned}$$

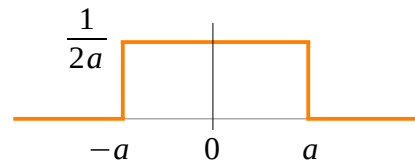
$$\text{Var}[w^{(i)}] = \frac{2}{n^{(i+1)}}$$

$$\begin{aligned} \text{Forward direction: } \text{Var}[w^{(i)}] &= \frac{2}{n^{(i)}} \\ \text{Backward direction: } \text{Var}[w^{(i)}] &= \frac{2}{n^{(i+1)}} \end{aligned} \quad \left. \vphantom{\begin{aligned} \text{Forward direction: } \text{Var}[w^{(i)}] &= \frac{2}{n^{(i)}} \\ \text{Backward direction: } \text{Var}[w^{(i)}] &= \frac{2}{n^{(i+1)}} \end{aligned}} \right\} \text{ Either one can be used.}$$

## ■ Kaiming He initializer: Uniform distribution

- Now that we know the optimal variance for the normal distribution, we can also find the range for the uniform distribution.

- continuous uniform distribution



$$g(x) = \begin{cases} \frac{1}{2a}, & \text{if } -a \leq x \leq a \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Var}[x] = \int_{-a}^a x^2 \cdot \frac{1}{2a} dx = \left[ \frac{1}{3} x^3 \frac{1}{2a} \right]_{-a}^a = \frac{1}{3} a^2$$

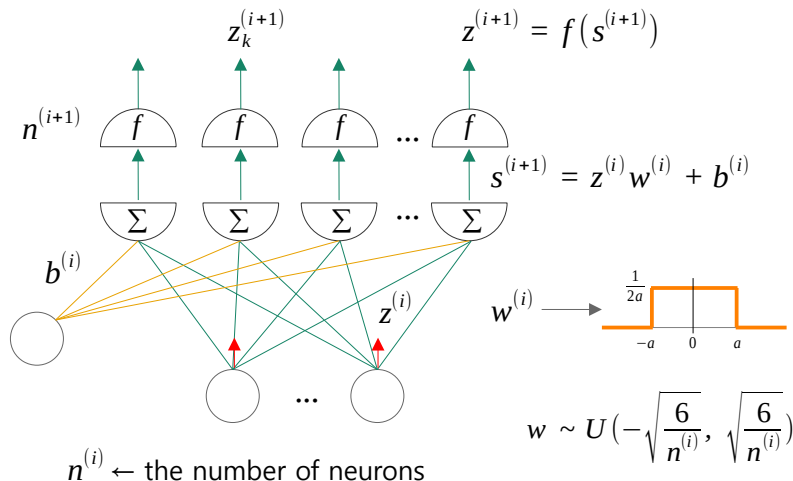
- Forward and backward directions

$$(1) \forall i, \text{Var}[z^{(i)}] = \text{Var}[z^{(i+1)}] \quad \leftarrow \text{Forward direction}$$

$$(2) \forall i, \text{Var}\left[\frac{\partial L}{\partial z^{(i)}}\right] = \text{Var}\left[\frac{\partial L}{\partial z^{(i+1)}}\right] \quad \leftarrow \text{Backward direction}$$

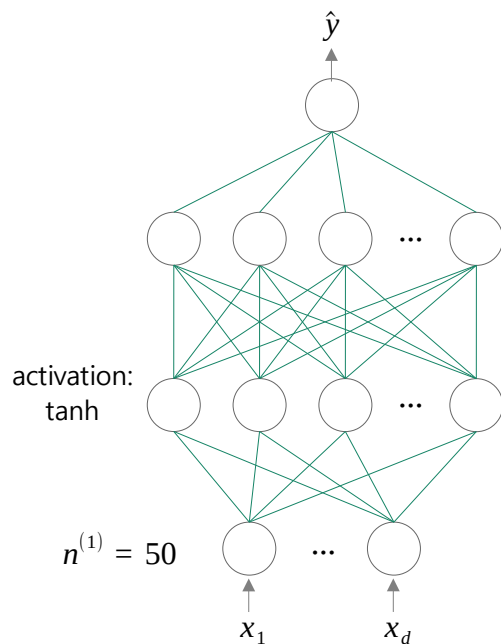
- We already know the variance that satisfies the above conditions, and the variance of the uniform distribution must also be equal to that variance.

$$\frac{1}{3} a^2 = \frac{2}{n^{(i)}} \longrightarrow a = \sqrt{\frac{6}{n^{(i)}}}$$



## ■ Kaiming He initializer: Example

- Given an ANN network like the one below, set appropriate initial distributions of  $w^{(1)}$  and  $w^{(2)}$ .



### ■ Normal distribution

$$w \sim N(0, \sigma^2) = N(0, \frac{2}{n^{(i)}})$$

$$n^{(3)} = 100$$

$$w^{(2)} \sim N(0, \frac{2}{80}) = N(0, 0.158^2)$$

$$n^{(2)} = 80$$

$$w^{(1)} \sim N(0, \frac{2}{50}) = N(0, 0.2^2)$$

### ■ Uniform distribution

$$w \sim U(-a, a) = U(-\sqrt{\frac{6}{n^{(i)}}}, \sqrt{\frac{6}{n^{(i)}}})$$

$$w^{(2)} \sim U(-\sqrt{\frac{6}{80}}, \sqrt{\frac{6}{80}}) = U(-0.274, 0.274)$$

$$w^{(1)} \sim U(-\sqrt{\frac{6}{50}}, \sqrt{\frac{6}{50}}) = U(-0.346, 0.346)$$

## ■ Keras Kaiming He initializer

# [MXDL-8-03] 3.kaiming\_he.py

```
import numpy as np
from tensorflow.keras.layers import Input, Dense
from tensorflow.keras import initializers
from tensorflow.keras.models import Model
import matplotlib.pyplot as plt
```

method = 'Normal'

if method == 'Normal':

init\_w1 = initializers.HeNormal()

init\_w2 = initializers.HeNormal()

else:

init\_w1 = initializers.HeUniform()

init\_w2 = initializers.HeUniform()

n\_in = 50

n\_s1 = 80

n\_s2 = 100

x = np.random.normal(size=(1000, n\_in))

# Create an ANN model

x\_input = Input(batch\_shape=(None, n\_in))

s1 = Dense(n\_s1, kernel\_initializer=init\_w1, activation='relu',  
name='W1')(x\_input)s2 = Dense(n\_s2, kernel\_initializer=init\_w2, activation='relu',  
name='W2')(s1)

y\_output = Dense(1, activation='sigmoid')(s2)

s1\_model = Model(x\_input, s1)

s2\_model = Model(x\_input, s2)

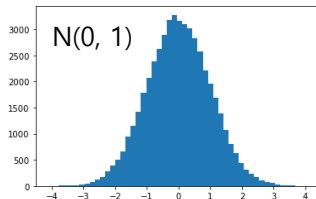
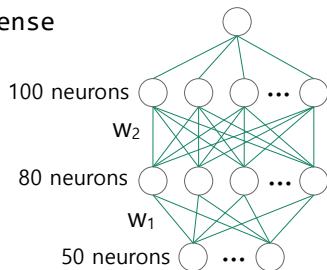
model = Model(x\_input, y\_output)

w1 = model.get\_layer('W1').get\_weights()[0].reshape(-1,)

w2 = model.get\_layer('W2').get\_weights()[0].reshape(-1,)

s1\_out = s1\_model.predict(x, verbose=0).reshape(-1,)

s2\_out = s2\_model.predict(x, verbose=0).reshape(-1,)



plt.hist(x.reshape(-1,), bins=50); plt.show()

plt.hist(s1\_out, bins=50); plt.show()

plt.hist(s2\_out, bins=50); plt.show()

if method == 'Normal':

print('[Keras ]  $\sigma$  of w1 = {:.3f}'.format(w1.std()))print('[Formula]  $\sigma$  of w1 = {:.3f}'.\n

format(np.sqrt(2 / n\_in)))

else:

print('[Keras ] a of w1 = {:.3f} ~ {:.3f}'.\n

format(w1.min(), w1.max()))

print('[Formula] a of w1 =  $\pm$ {:.3f}'.format(np.sqrt(6 / n\_in)))

if method == 'Normal':

print('\n[Keras ]  $\sigma$  of w2 = {:.3f}'.format(w2.std()))print('[Formula]  $\sigma$  of w2 = {:.3f}'.format(np.sqrt(2 / n\_s1)))

else:

print('\n[Keras] a of w2 = {:.3f} ~ {:.3f}'.\n

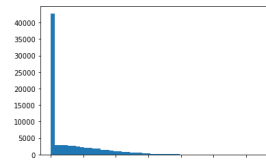
format(w2.min(), w2.max()))

print('[Formula] a of w2 =  $\pm$ {:.3f}'.format(np.sqrt(6 / n\_s1)))

method = 'Normal':

[Keras ]  $\sigma$  of w1 = 0.198[Formula]  $\sigma$  of w1 = 0.200[Keras ]  $\sigma$  of w2 = 0.158[Formula]  $\sigma$  of w2 = 0.158

s1\_out →



method = 'Uniform':

[Keras ] a of w1 = -0.346 ~ 0.346

[Formula] a of w1 =  $\pm$ 0.346

[Keras] a of w2 = -0.274 ~ 0.274

[Formula] a of w2 =  $\pm$ 0.274

s2\_out →

