

# M2177.0043 Introduction to Deep Learning

## Lecture 8: Neural Networks<sup>1</sup>

Hyun Oh Song<sup>1</sup>

<sup>1</sup>Dept. of Computer Science and Engineering, Seoul National University

April 9, 2020

---

<sup>1</sup>Many slides and figures adapted Justin Johnson

## Last time

- ▶ Model fitting
- ▶ Computational Graphs
- ▶ Neural Networks

# Outline

Convolutional neural network

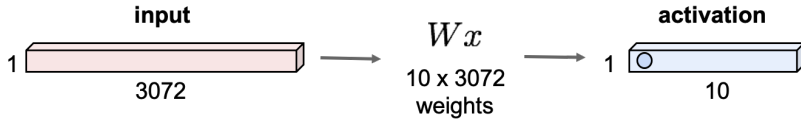
Pooling layer

Activation

Weight initialization

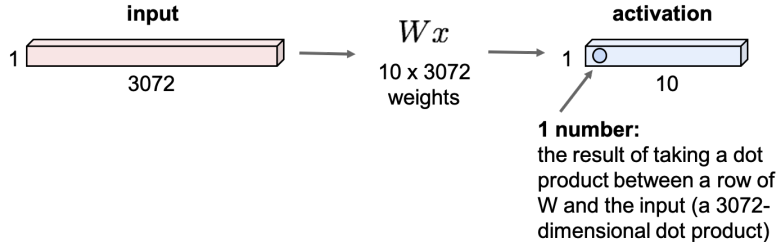
## Fully connected layer

- ▶  $32 \times 32 \times 3$  image  $\rightarrow$  stretch to  $3072 \times 1$



## Fully connected layer

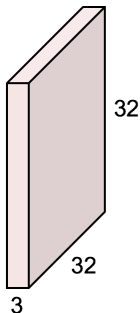
- ▶  $32 \times 32 \times 3$  image  $\rightarrow$  stretch to  $3072 \times 1$



## Convolution layer

- ▶  $32 \times 32 \times 3$  image  $\rightarrow$  preserve spatial structure of input

32x32x3 image

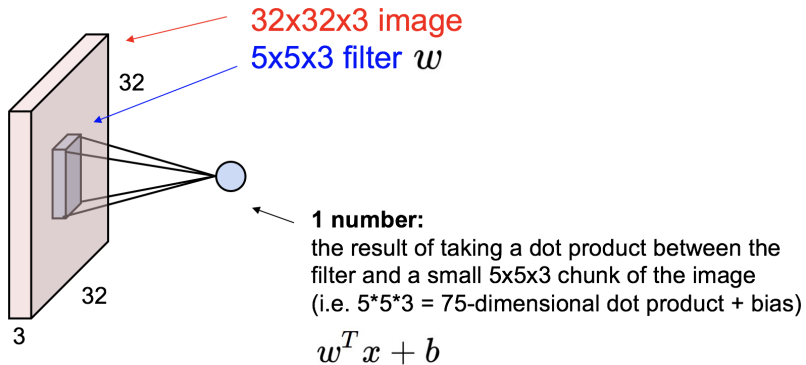


5x5x3 filter

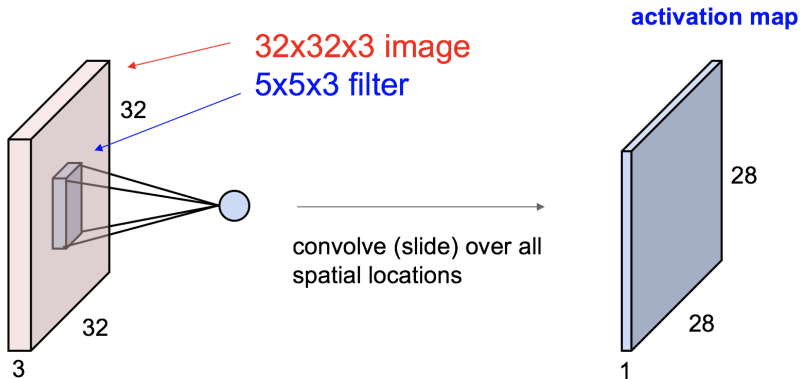


**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

## Convolution layer



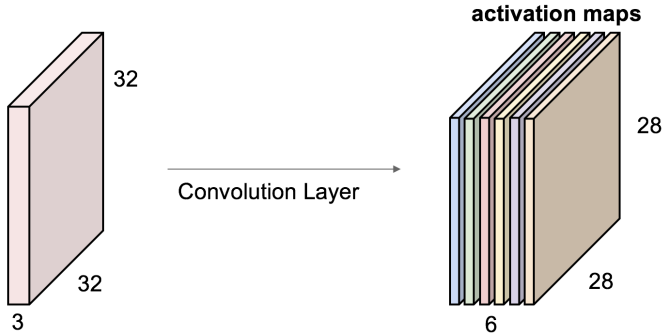
## Convolution layer





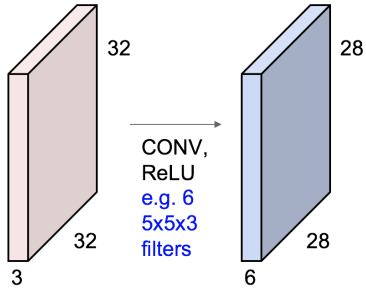
## Convolution layer

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

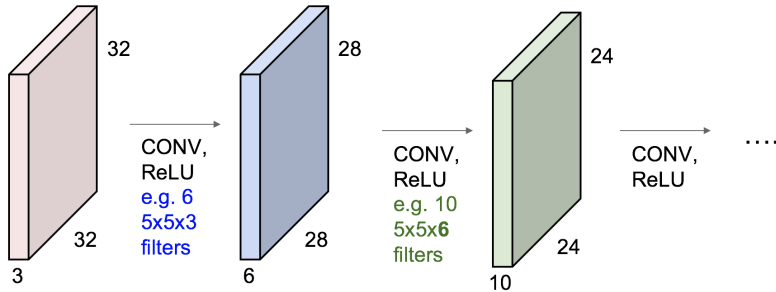


We stack these up to get a “new image” of size 28x28x6!

## Convolution layer



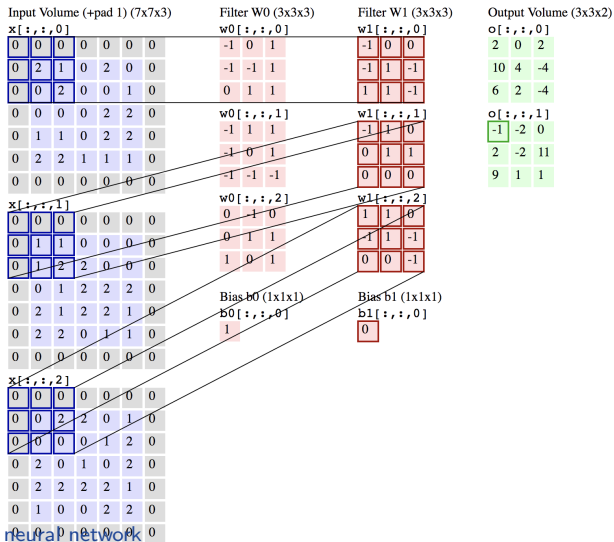
## Convolution layer



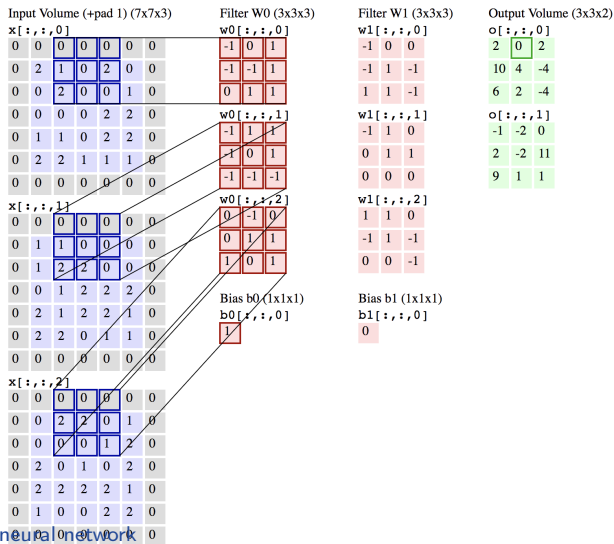
## Convolution layer hyperparameters

- ▶ Number of filters. Depth of the output volume = number of activation maps.
- ▶ Stride. When the stride is 1 then we move the filters one pixel at a time. When the stride is 2 then the filters jump 2 pixels at a time as we slide them around. This will produce smaller output volumes spatially.
- ▶ Zero-padding. The nice feature of zero padding is that it will allow us to control the spatial size of the output volumes

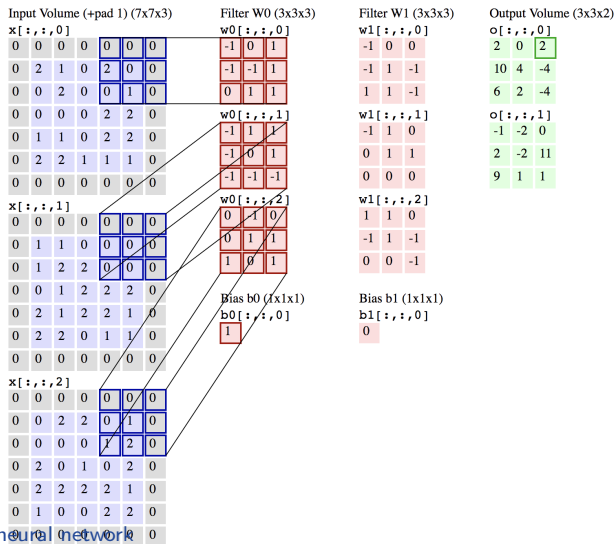
# Conv layer w/ num\_filters=2, stride=2, padding=1



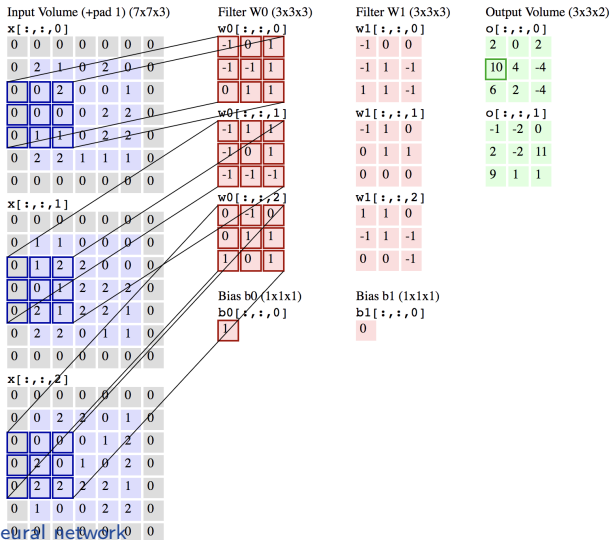
# Conv layer w/ num\_filters=2, stride=2, padding=1



# Conv layer w/ num\_filters=2, stride=2, padding=1

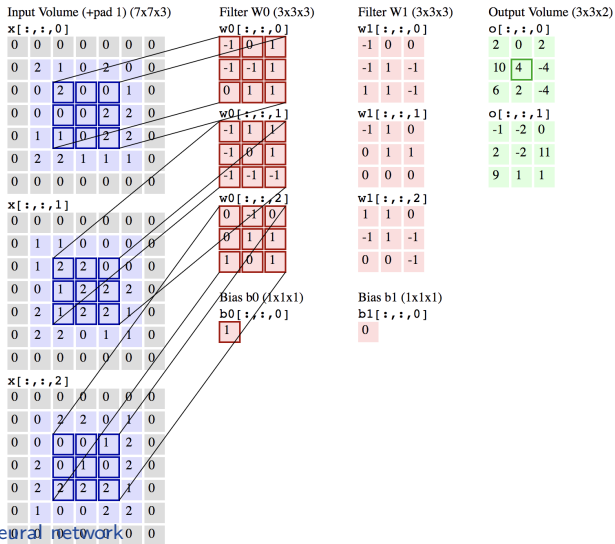


# Conv layer w/ num\_filters=2, stride=2, padding=1

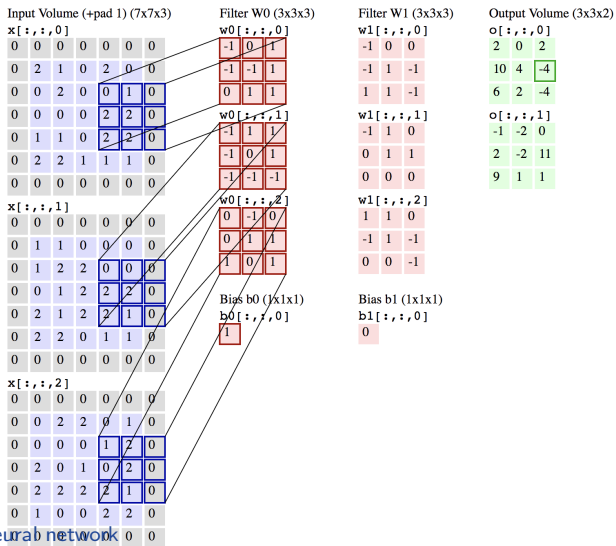




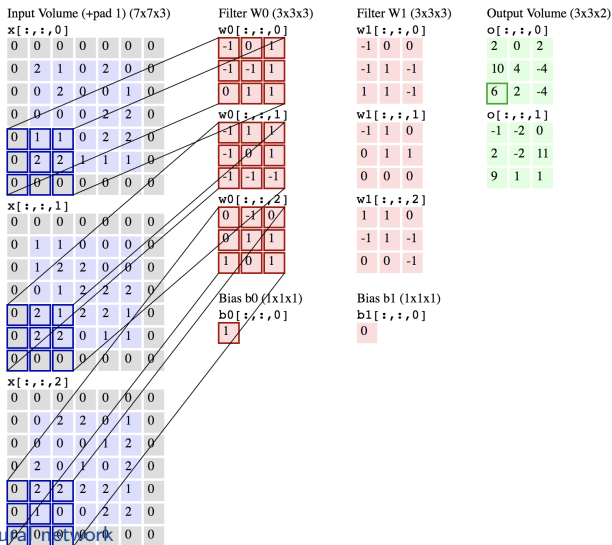
# Conv layer w/ num\_filters=2, stride=2, padding=1



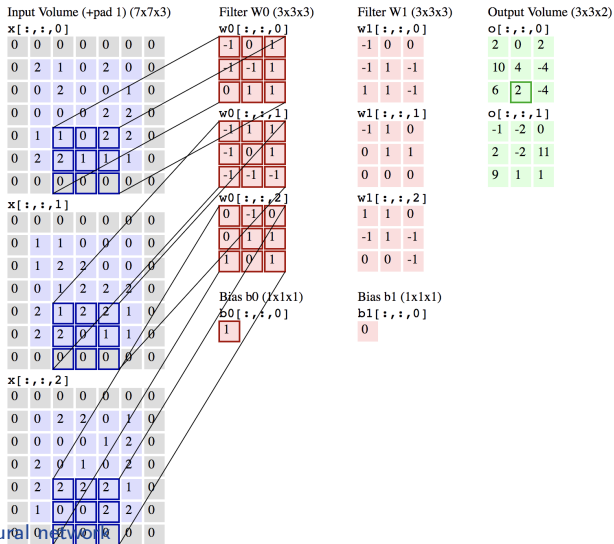
# Conv layer w/ num\_filters=2, stride=2, padding=1



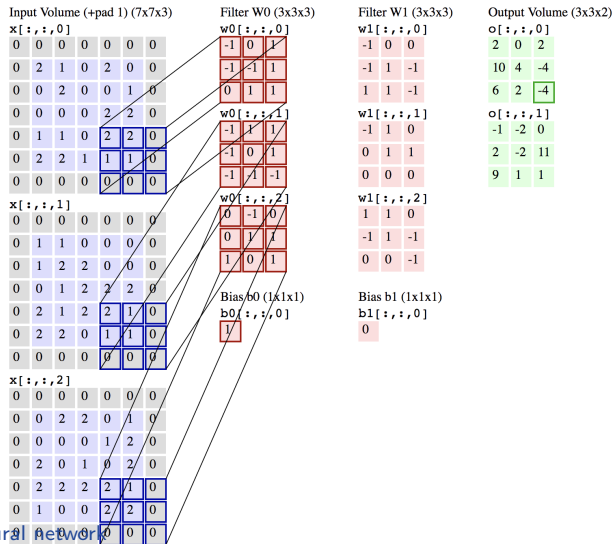
# Conv layer w/ num\_filters=2, stride=2, padding=1



# Conv layer w/ num\_filters=2, stride=2, padding=1



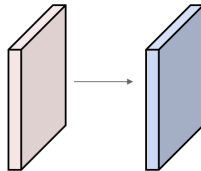
# Conv layer w/ num\_filters=2, stride=2, padding=1



## Example 1: Output activation size

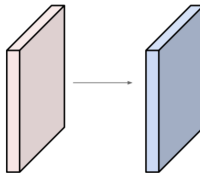
Input volume: **32x32x3**  
10 5x5 filters with stride 1, pad 2

Output volume size: ?



## Example 1: Output activation size

Input volume: **32x32x3**  
**10** **5x5** filters with stride **1**, pad **2**

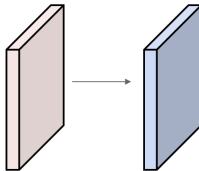


Output volume size:  
 $(32 + 2 \cdot 2 - 5) / 1 + 1 = 32$  spatially, so  
**32x32x10**

## Example 2: Number of parameters

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2



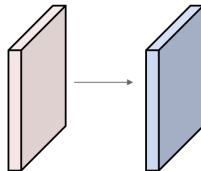
Number of parameters in this layer?



## Example 2: Number of parameters

Input volume: **32x32x3**

**10** **5x5** filters with stride 1, pad 2



Number of parameters in this layer?

each filter has  $5*5*3 + 1 = 76$  params (+1 for bias)

=>  $76*10 = 760$

# Outline

Convolutional neural network

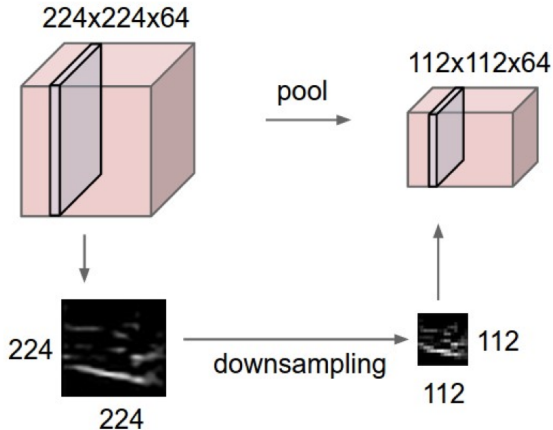
Pooling layer

Activation

Weight initialization

## Pooling layer

- ▶ makes the representations smaller and more manageable
- ▶ operates over each activation map independently



## Max pooling, average pooling, etc.

Single depth slice

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pool with 2x2 filters  
and stride 2



6	8
3	4

# Outline

Convolutional neural network

Pooling layer

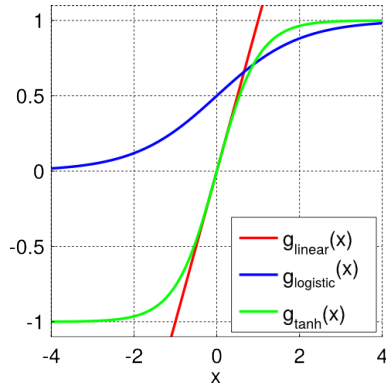
Activation

Weight initialization

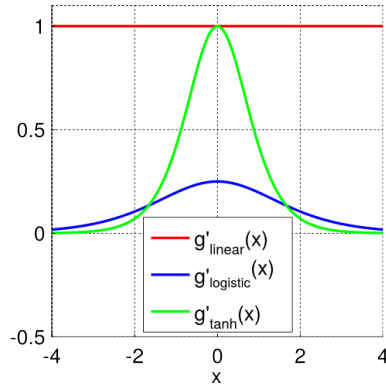
## Activation functions<sup>2</sup>

- ▶  $g_{\text{logistic}} = 1/(1 + e^{-z})$
- ▶  $g_{\text{tanh}} = (e^z - e^{-z})/(e^z + e^{-z})$

Some Common Activation Functions



Activation Function Derivatives



# Outline

Convolutional neural network

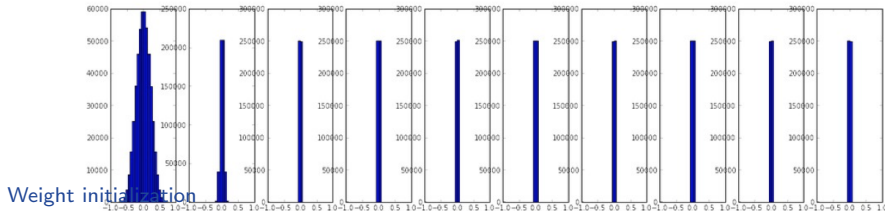
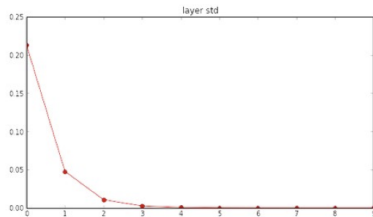
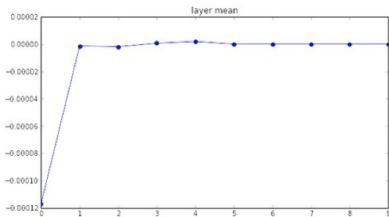
Pooling layer

Activation

Weight initialization

## 10 Layer MLP w/ tanh: small weights

- ▶ `W = 0.01 * np.random.randn(n_in, n_out)`
- ▶ Activations get collapsed to zero. Why?



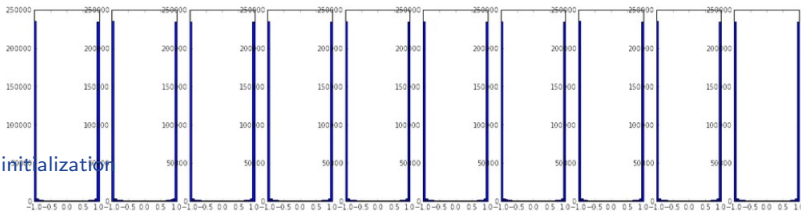
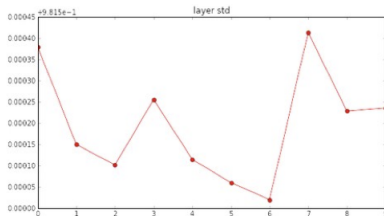
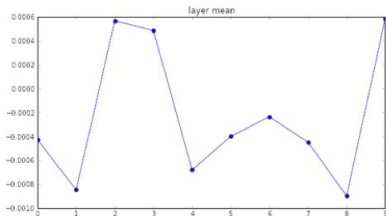


## 10 Layer MLP w/ tanh: small weights

- ▶ Activations get saturated at 0 and gradient doesn't flow. Why?
- ▶ Forward: Keep multiplying near zero matrices.
- ▶ Backward: Keep multiplying near zero matrices as well. Upstream gradient collapsing to zero as well.
- ▶ Activations collapse to all zeros.
- ▶ So the if the problem is because  $W$  is initialized with **small weights**, what if we initialize with **big weights instead**?

## 10 Layer MLP w/ tanh: large weights

- ▶ `W = 1.0 * np.random.randn(n_in, n_out)`
- ▶ Activations get saturated at  $-1$  or  $1$  and gradient doesn't flow. Why?

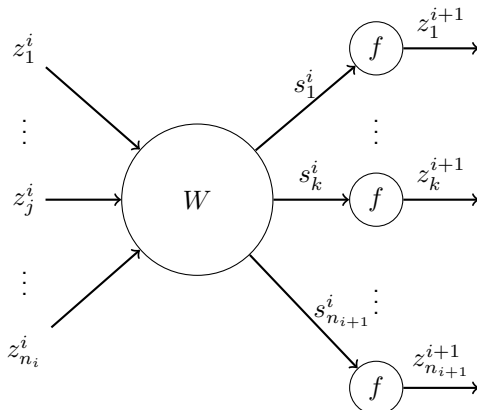


## 10 Layer MLP w/ tanh: large weights

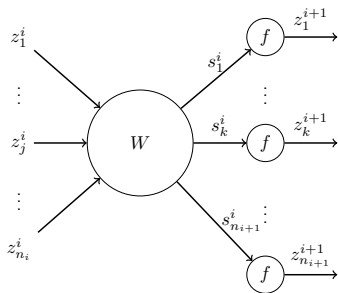
- ▶ Activations get saturated at  $-1$  or  $1$  and gradient doesn't flow.
- ▶ Forward: Now, **tanh** saturates large activations to  $-1$  or  $+1$
- ▶ Backward: **tanh** has near zero gradient at  $-1$  or  $+1$  and the weights do not update

## Xavier initialization

- ▶ Assume MLP network, weights and activations have zero mean, weights drawn i.i.d.
- ▶  $f(\cdot)$  is the non-linearity, centered & linear & has unit gradient around zero, *i.e.* tanh layer.



## Xavier initialization - forward variance



$$z_k^{i+1} = f(s_k^{i+1}) = f\left(\sum_{j=1}^{n_i} W_{kj}^i z_j^i\right) \approx \sum_{j=1}^{n_i} W_{kj}^i z_j^i$$

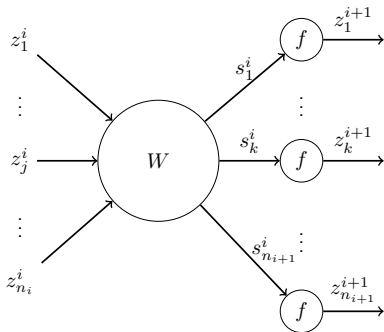
Note the property of variance and by independence ( $X \perp\!\!\!\perp Y$ ),  
 $Var(XY) = E[Y]^2 Var(X) + E[X]^2 Var(Y) + Var(X)Var(Y)$ .  
Thus,

$$Var(z^{i+1}) \approx Var\left(\sum_{j=1}^{n_i} W_{kj}^i z_j^i\right)$$

$$\implies Var(z^{i+1}) = n_i Var(W^i) Var(z^i)$$

To achieve  $Var(z^i) = Var(z^{i+1})$ ,  
set  $Var(W^i) = \frac{1}{n_i}$

## Xavier initialization - backward variance



$$\begin{aligned}
 \frac{\partial \ell}{\partial z_j^i} &= \sum_{k=1}^{n_{i+1}} \overbrace{\frac{\partial \ell}{\partial z_k^{i+1}}}^{\text{upstream grad}} \overbrace{\frac{\partial z_k^{i+1}}{\partial z_j^i}}^{\text{local grad}} \\
 &= \sum_{k=1}^{n_{i+1}} \frac{\partial \ell}{\partial z_k^{i+1}} \underbrace{\frac{\partial z_k^{i+1}}{\partial s_k^{i+1}} \frac{\partial s_k^{i+1}}{\partial z_j^i}}_{\approx 1} \\
 &= \sum_{k=1}^{n_{i+1}} \frac{\partial \ell}{\partial z_k^{i+1}} W_{kj}^i
 \end{aligned}$$

$$\Rightarrow \text{Var}\left(\frac{\partial \ell}{\partial z^i}\right) = n_{i+1} \text{Var}\left(\frac{\partial \ell}{\partial z^{i+1}}\right) \text{Var}(W^i)$$

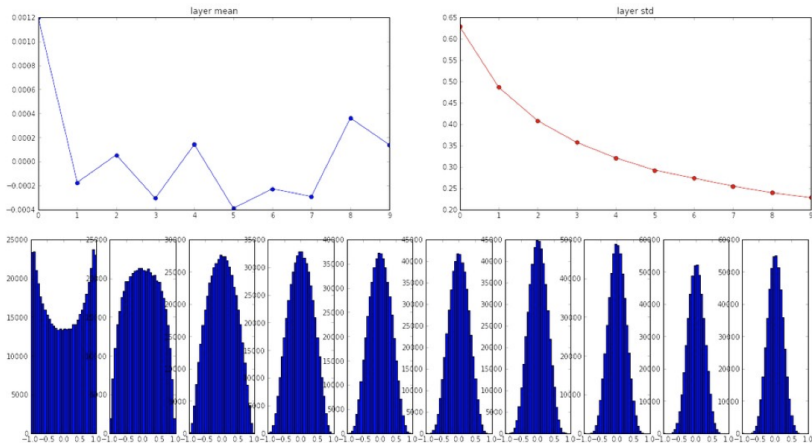
To achieve  $\text{Var}\left(\frac{\partial \ell}{\partial z^i}\right) = \text{Var}\left(\frac{\partial \ell}{\partial z^{i+1}}\right)$ ,  
 set  $\text{Var}(W^i) = \frac{1}{n_{i+1}}$

## Xavier initialization

- ▶ unless  $n_i = n_{i+1}$ , we have to compromise between these two conditions, and a reasonable choice is the harmonic mean,  
$$Var(W^i) = \frac{2}{n_i + n_{i+1}}$$
- ▶ Note, some implementations (Caffe) just use  $Var(W^i) = \frac{1}{n_i}$
- ▶ if we sample from  $N(0, 1)$ , multiply the samples by  $\sqrt{\frac{2}{n_i + n_{i+1}}}$ .
- ▶ if we sample from  $Uni(-a, a)$ , take  $a = \sqrt{\frac{6}{n_i + n_{i+1}}}$ . Why?
- ▶ Read Glorot & Bengio 2010 for more details.

## 10 Layer MLP w/ tanh

► `W = np.random.randn(n_in, n_out) / np.sqrt(n_in)`





## He initialization

- ▶ The assumptions do not hold for  $f(\cdot)$  is ReLU activation.
- ▶ Roughly, ReLU zeroes out half the inputs.
- ▶ Set  $Var(W^i) = \frac{2}{n_i}$ .
- ▶ Read He *et al.* 2015 for more details.
- ▶ How to best initialize deep networks is an active research area.