

# **Deep Learning Bootcamp Day 2**

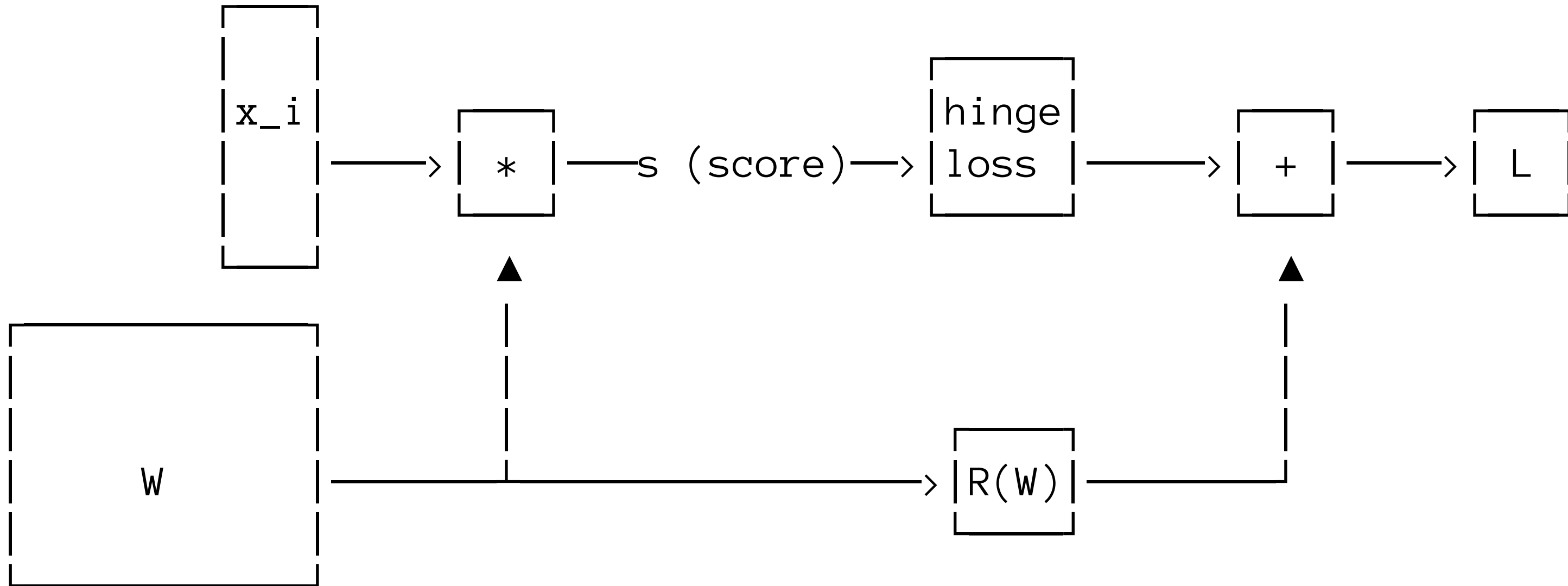
**Dr. Kashif Rasul<sup>1</sup>**

**25.09.2018 at Center for Systems Biology Dresden**

---

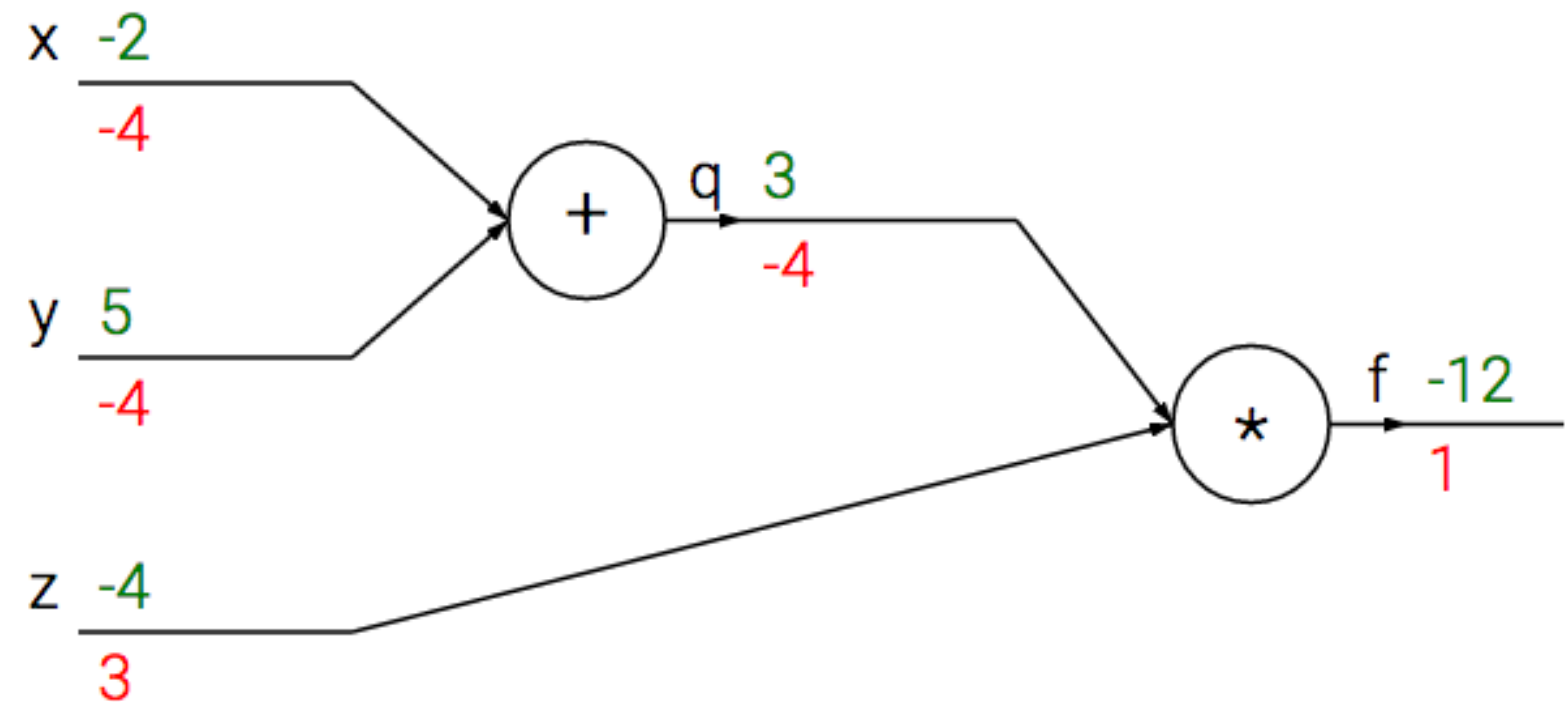
<sup>1</sup>[Zalando Research](#) Twitter: [@krasul](#) Github: [kashif](#)

# Computational Graph



# Compound expression

- $f(x, y, z) = (x + y)z$
- E.g.  $x = -2, y = 5, z = -4$
- $q = x + y, \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$
- $f = qz, \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$
- **Want:**  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$
- **Chain rule:**  $\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$



activations

$$x$$

“local gradient”

$$\frac{\partial z}{\partial x}$$

f

$$\frac{\partial z}{\partial y}$$

$$z$$

$$\frac{\partial L}{\partial z}$$

gradients

$$y$$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial y}$$

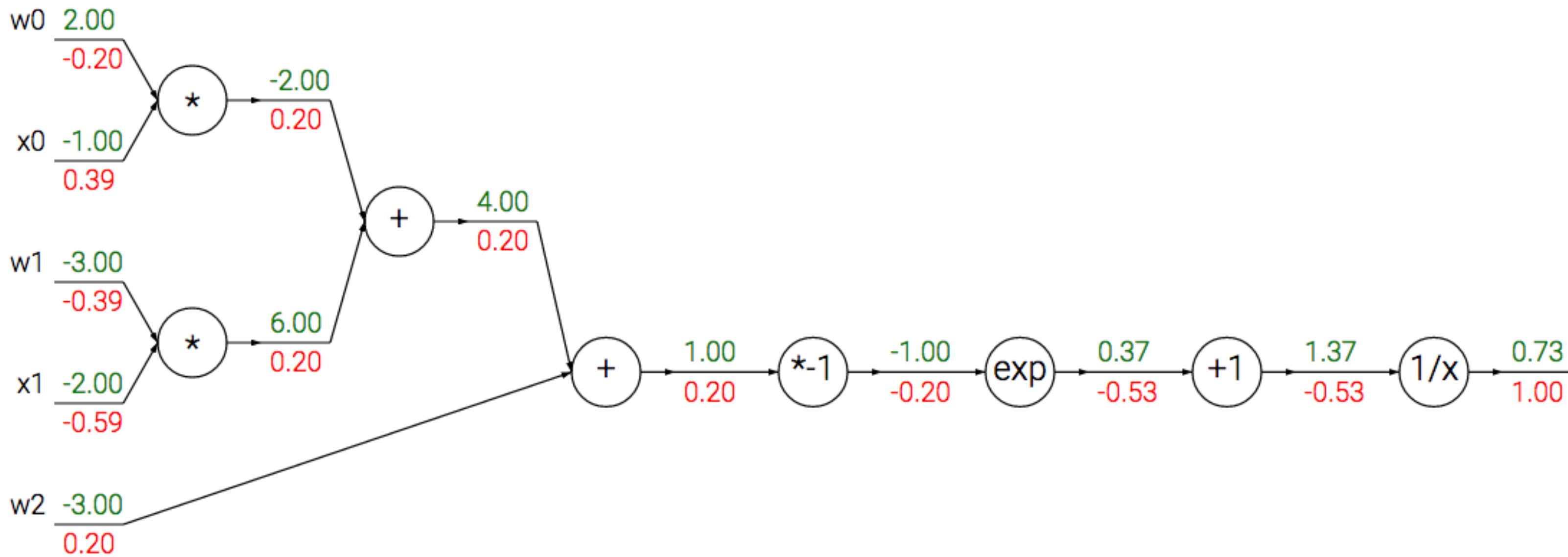
$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x}$$

# Modularity

- Any differentiable function can act as a gate
- Can group multiple gates into a single gate
- Decompose a function into multiple gates

$$f(\vec{w}, \vec{x}) = \frac{1}{1 + \exp(-(w_0 x_0 + w_1 x_1 + w_2))}$$

- We have these new gates:
  - $f(x) = 1/x$  and  $df/dx = -1/x^2$
  - $f(x) = c + x$  and  $df/dx = 1$
  - $f(x) = \exp(x)$  and  $df/dx = \exp(x)$
  - $f(x) = ax$  and  $df/dx = a$



# Sigmoid function

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

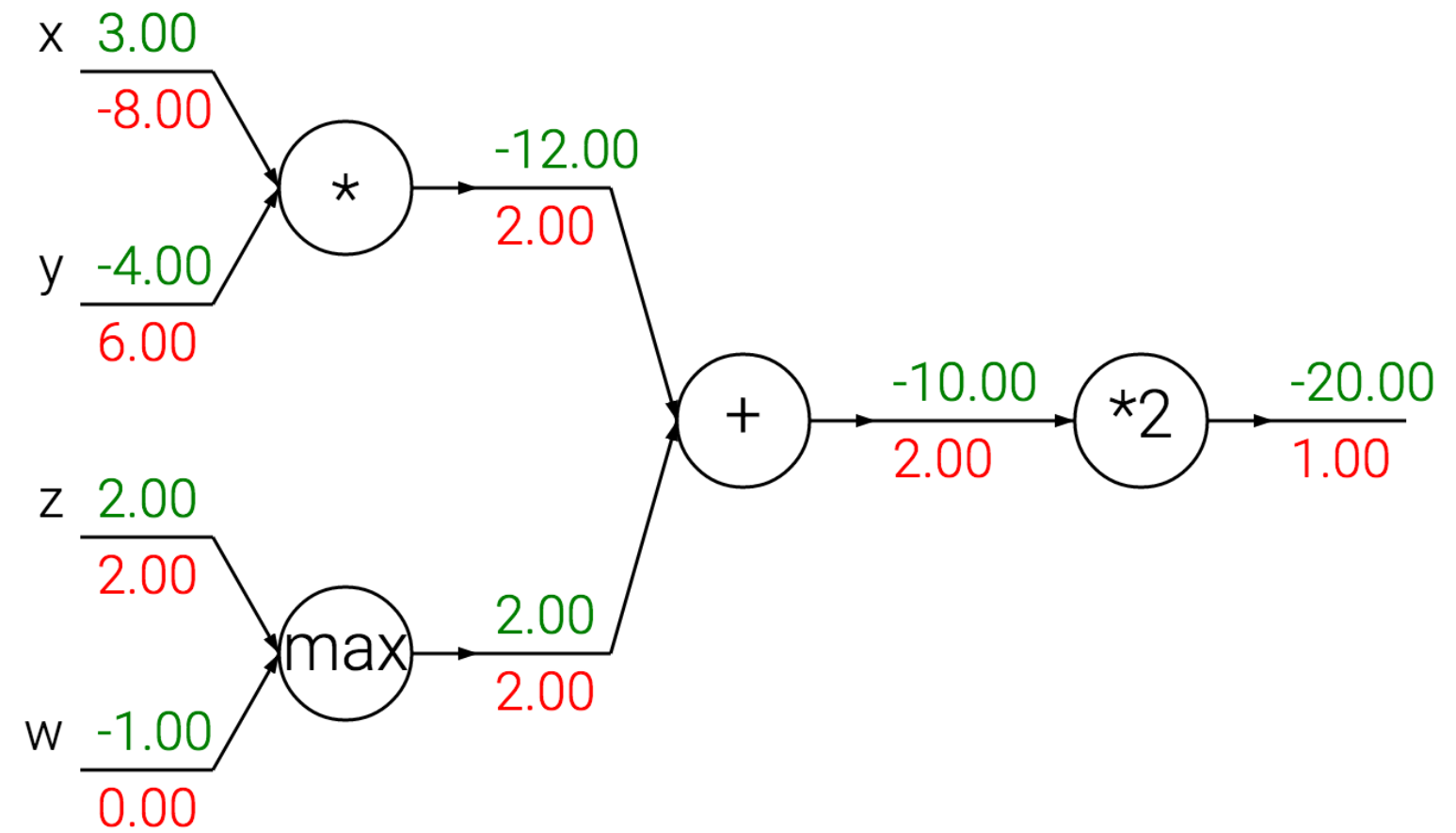
for which

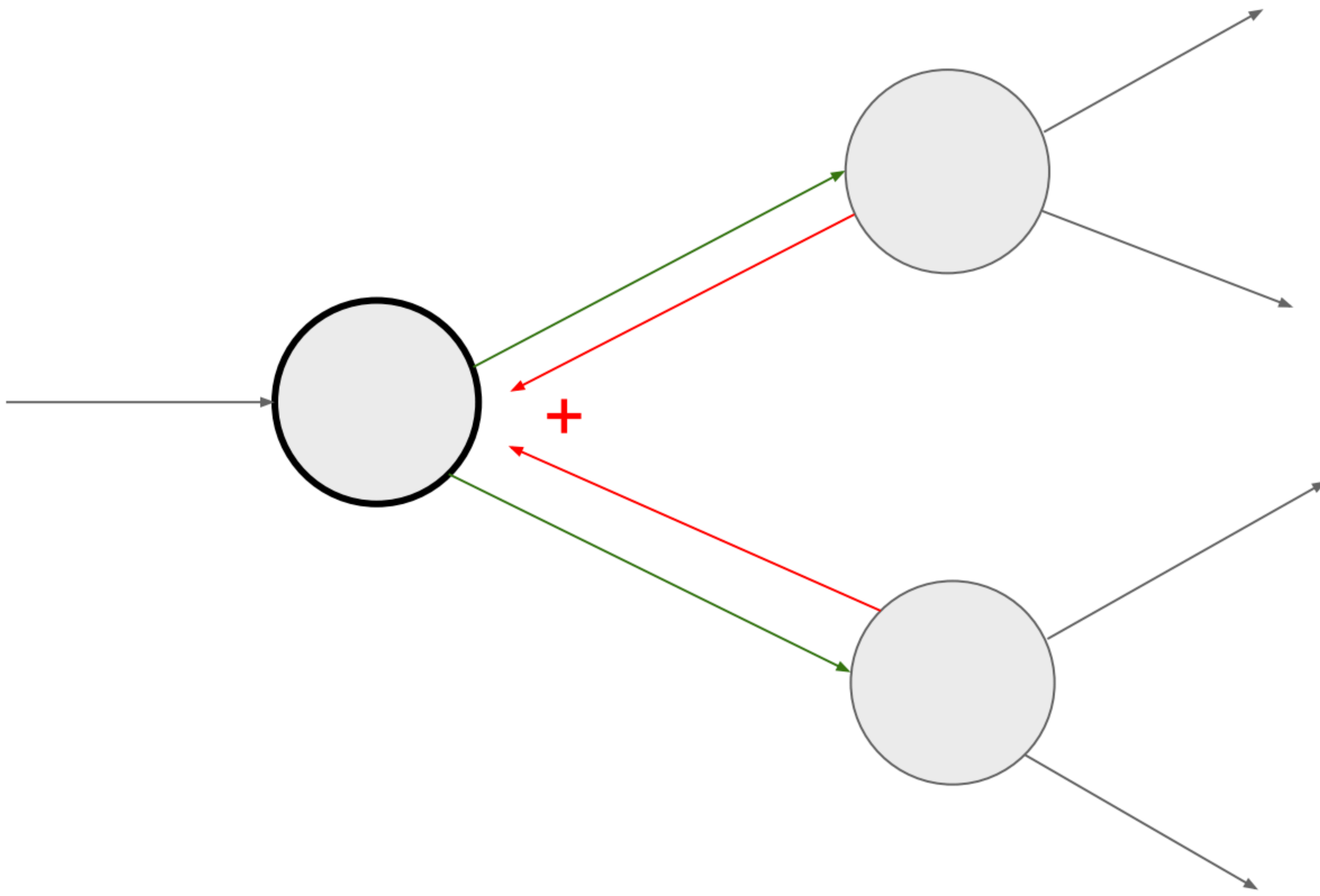
$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \sigma(x)(1 - \sigma(x))$$



# Patterns

- Add gate: gradient distributor
- Max gate: gradient router
- Multiply gate: gradient "switcher"





# Unintuitive effects

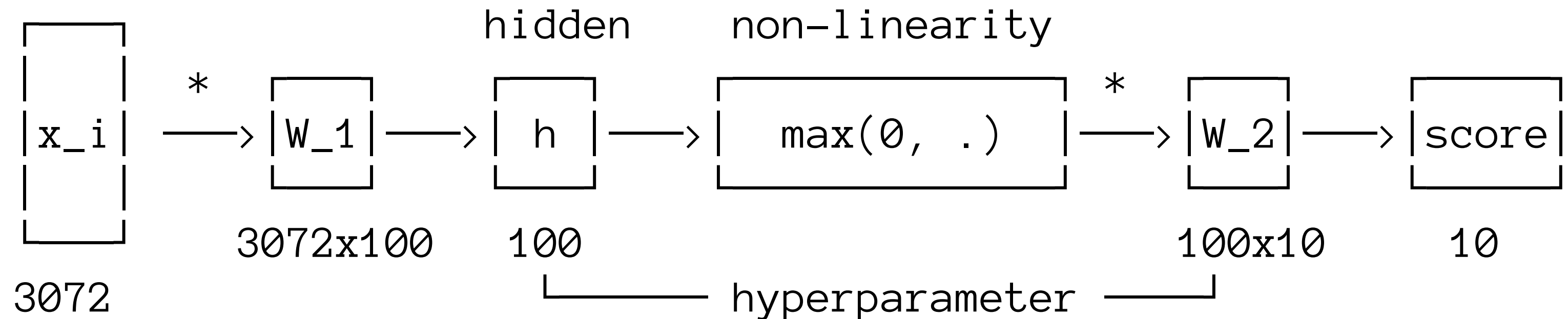
- If one of the inputs to the multiply gate is very small and the other is very big, then multiply gate will do something unintuitive
- It will assign a relatively huge gradient to the small input and a tiny gradient to the large input

# Gradients of vectorized operators

- All concepts extend in a straight-forward manner to matrix and vector operations
- One must pay closer attention to dimensions and transpose operations
- Matrix-Matrix multiply gradient: most tricky operation

# Neural Networks

- Linear score:  $f = W\vec{x}$
- 2-Layer Neural Network:  $f = W_2 \max(0, W_1 \vec{x})$

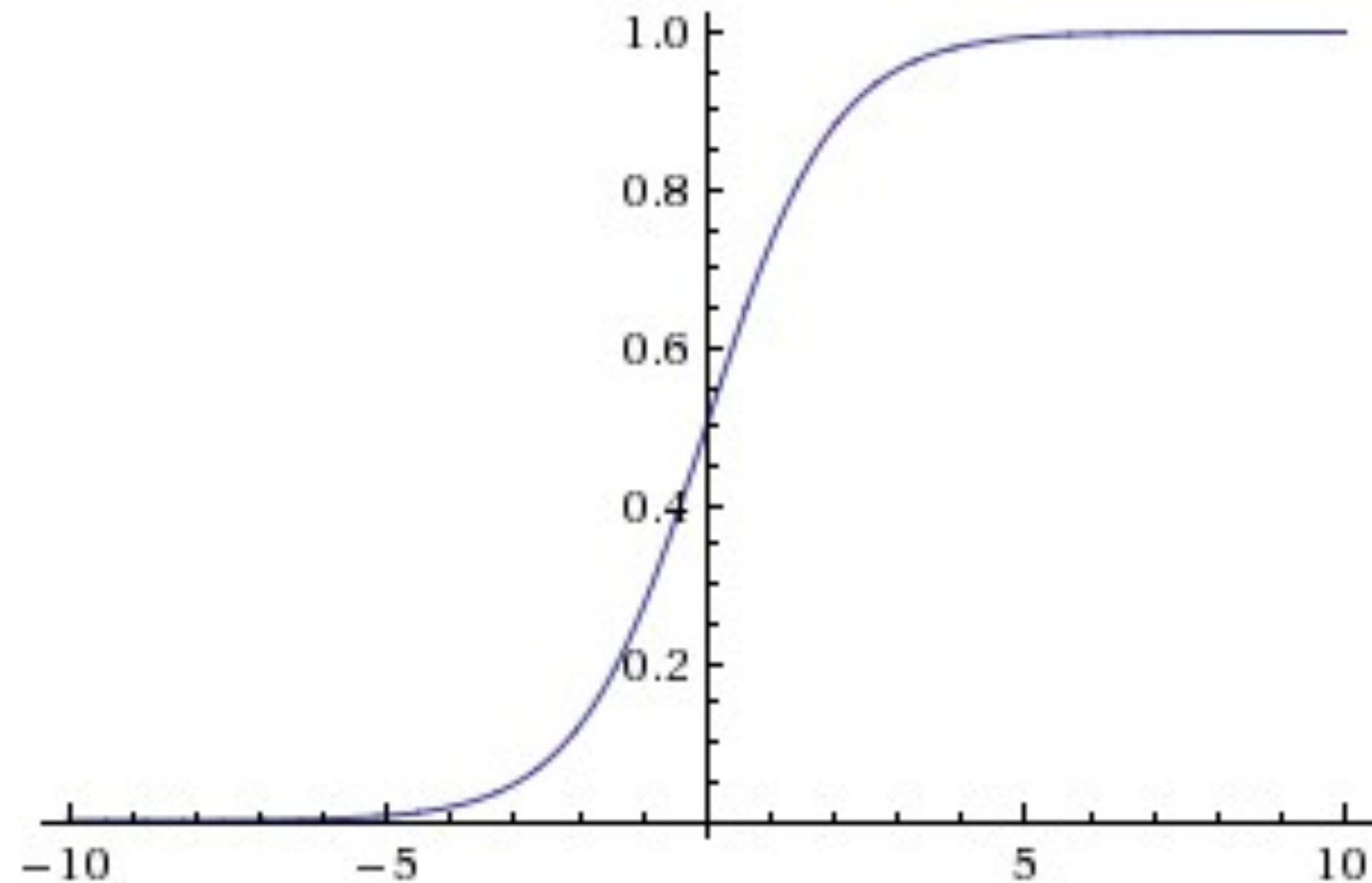


# Activation functions

- Sigmoid:  $\sigma(x) = 1/(1 + \exp(-x))$
- Tanh:  $\tanh(x)$
- ReLU (hinge):  $\max(0, x)$
- Leaky ReLU:  $\max(0.1x, x)$
- Others: Maxout, PReLU, ELU, etc. (active area of research)

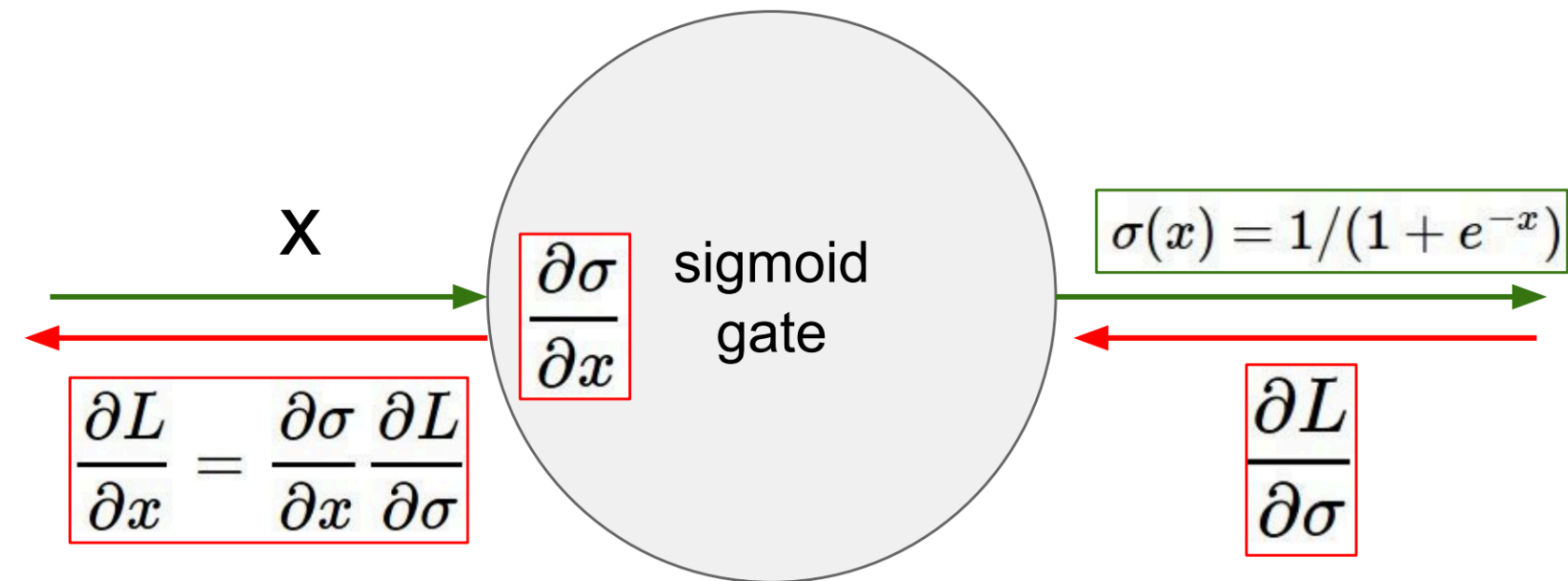
# Sigmoid activation

- $\sigma(x) = 1/(1 + \exp(-x))$
- Squashes numbers to the range (0, 1)
- Issues:
  1. Saturated neurons "kill" the gradient
  2. Outputs not zero centered
  3.  $\exp()$  computationally expensive



# 1. Saturation

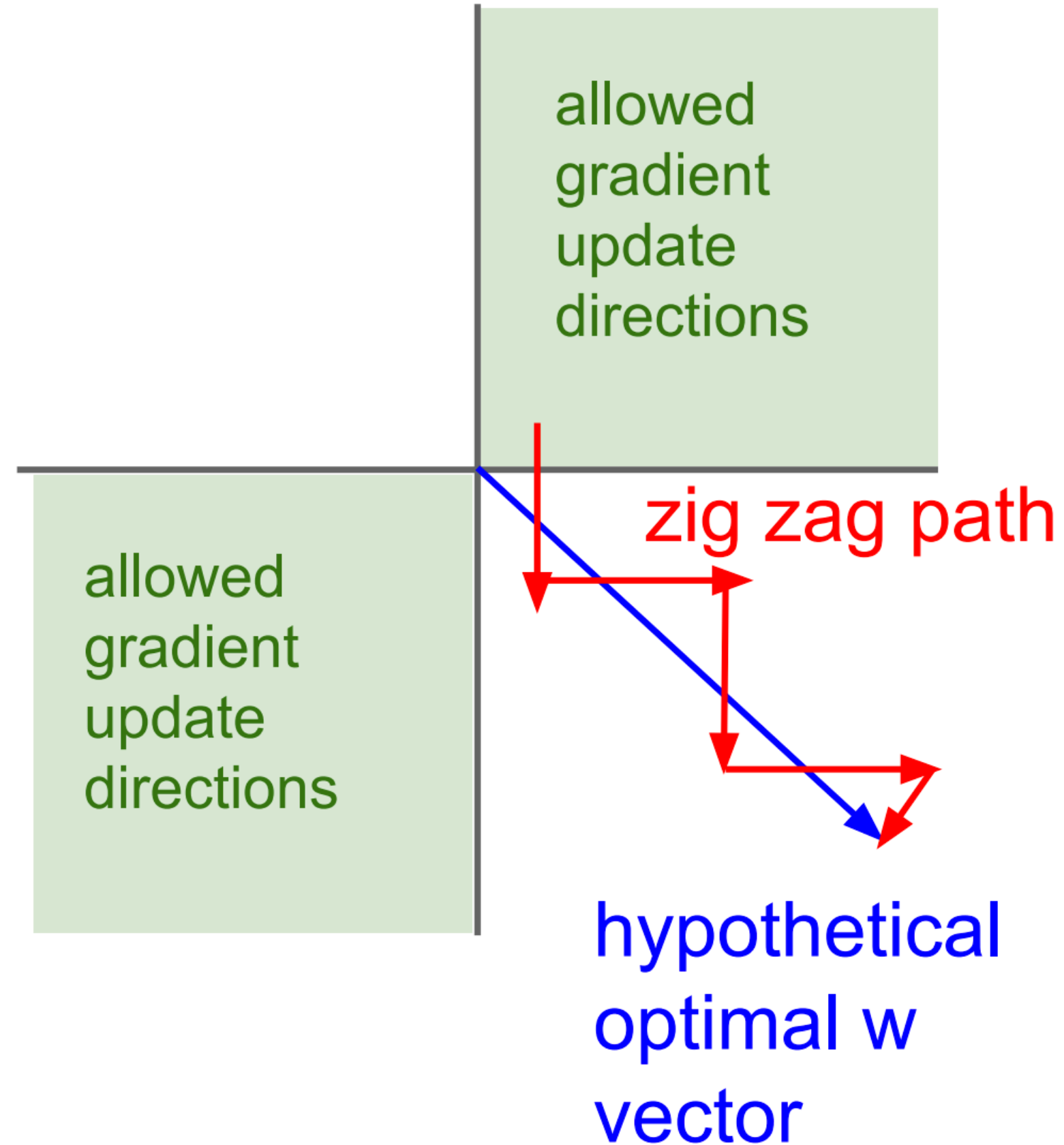
- $\frac{\partial \sigma}{\partial x} = \sigma(x)(1 - \sigma(x))$
- What happens to backward gradient when:
  - $x = -10$ ?
  - $x = 0$ ?
  - $x = 10$ ?





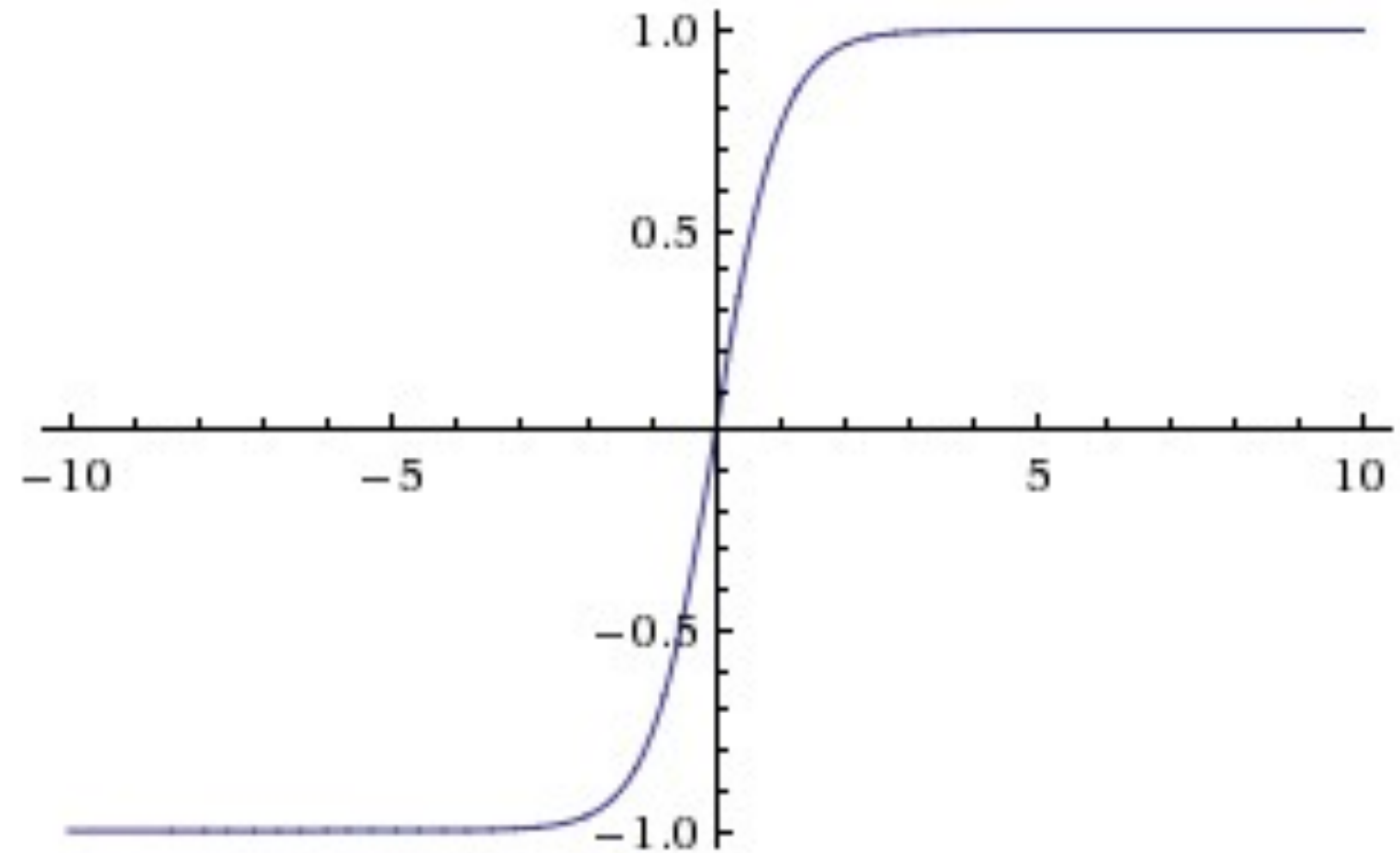
## 2. Positive output

- Suppose  $f = \sum_i w_i x_i + b, x_i > 0$
- $\frac{\partial f}{\partial w_i} = x_i$
- $\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial f} \frac{\partial f}{\partial w_i} = x_i \frac{\partial L}{\partial f}$
- So all the backward gradients **always** have the same sign as  $\frac{\partial L}{\partial f}$



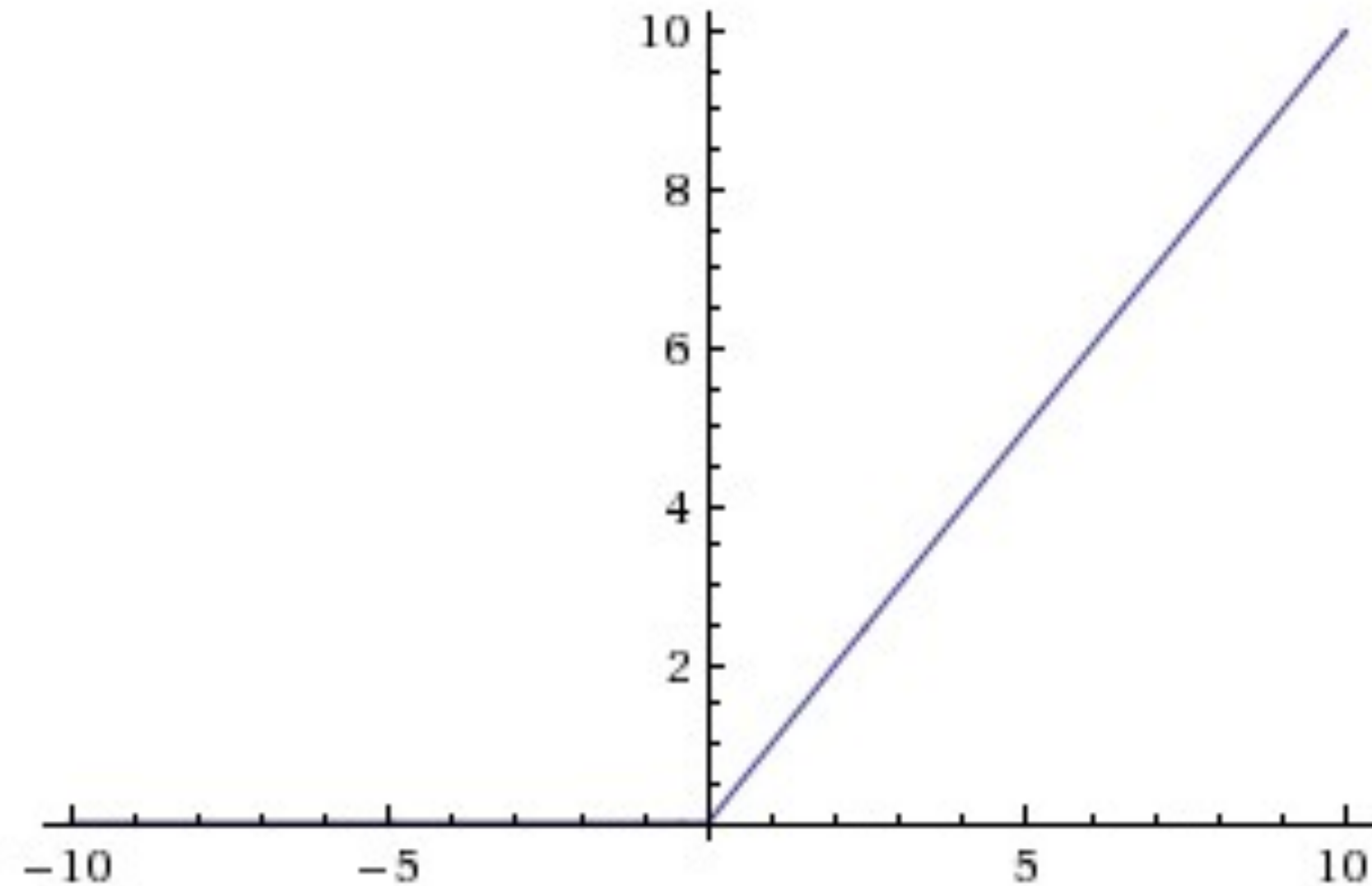
# Tanh

- Squashes numbers to  $(-1, 1)$
- Zero centered
- Still kills gradients when saturated



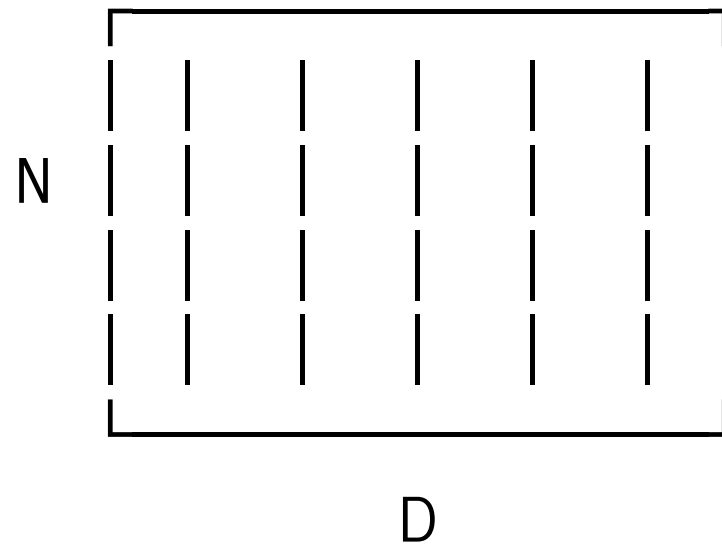
# ReLU (Rectified Linear Unit)

- $f(x) = \max(0, x)$
- Does not saturate
- Very computationally efficient
- Converges faster than sigmoid or tanh in practice
- Issues:
  - Non-zero centered
  - Not smooth



# Data Preprocessing: $X$ ( $[N, D]$ )

- Mean subtraction:  $X -= \text{np.mean}(X, \text{axis}=0)$
- Normalization:  $X /= \text{np.std}(X, \text{axis}=0)$



- PCA: decorrelates the data
- Whitening

# Preprocessing for images: center only

- Subtract the mean image: e.g. CIFAR-10 needs a  $[32, 32, 3]$  array
- Subtract per-channel mean: 3 numbers, e.g. with  $[N, C, W, H]$  format

$X[:, 0, :, :] == 103.939$  # R mean

$X[:, 1, :, :] == 116.779$  # G mean

$X[:, 2, :, :] == 123.68$  # B mean

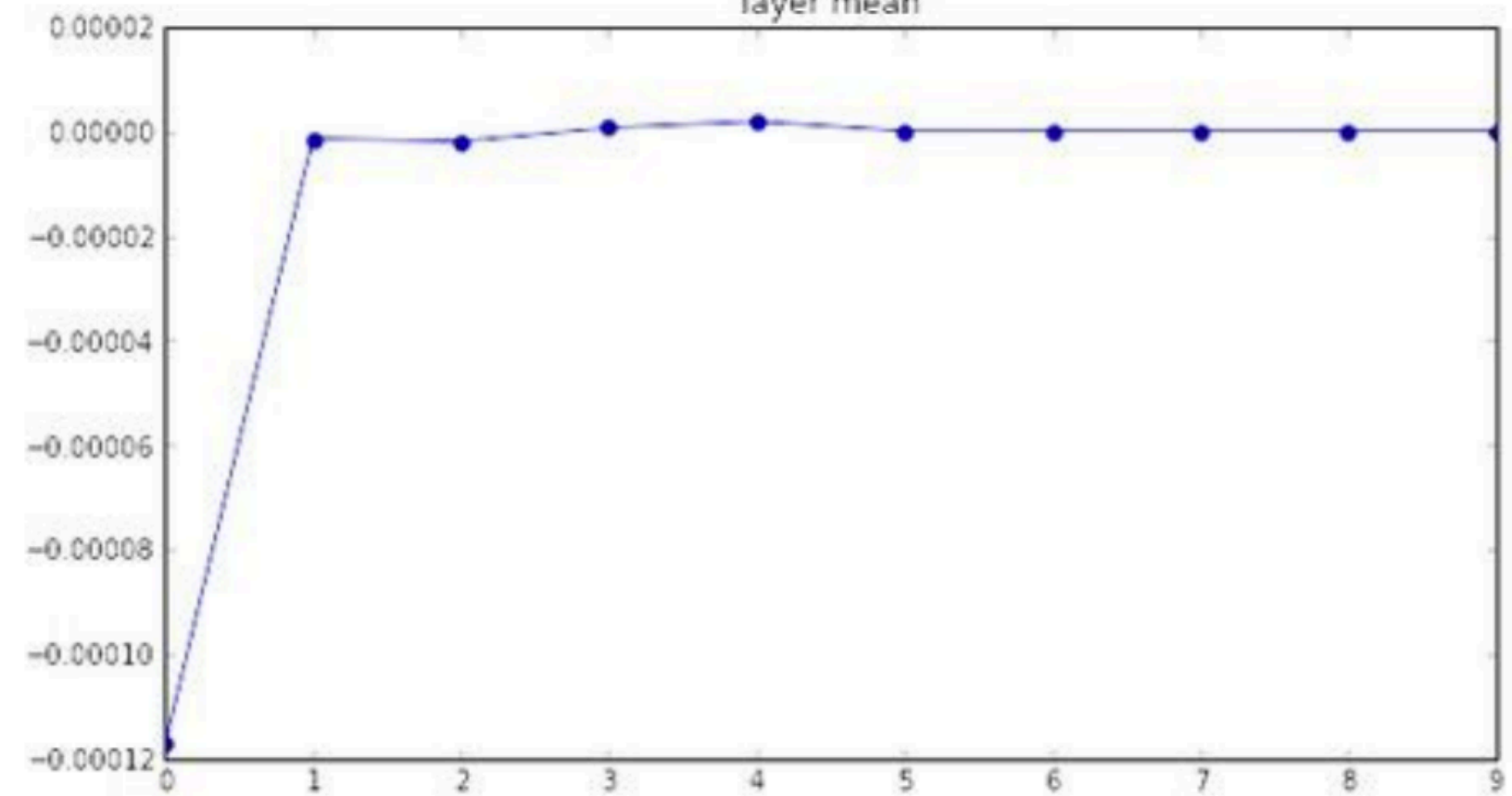
# Weight initialization

- What happens when  $W = 0$  initialization is used?
- First idea: small random numbers

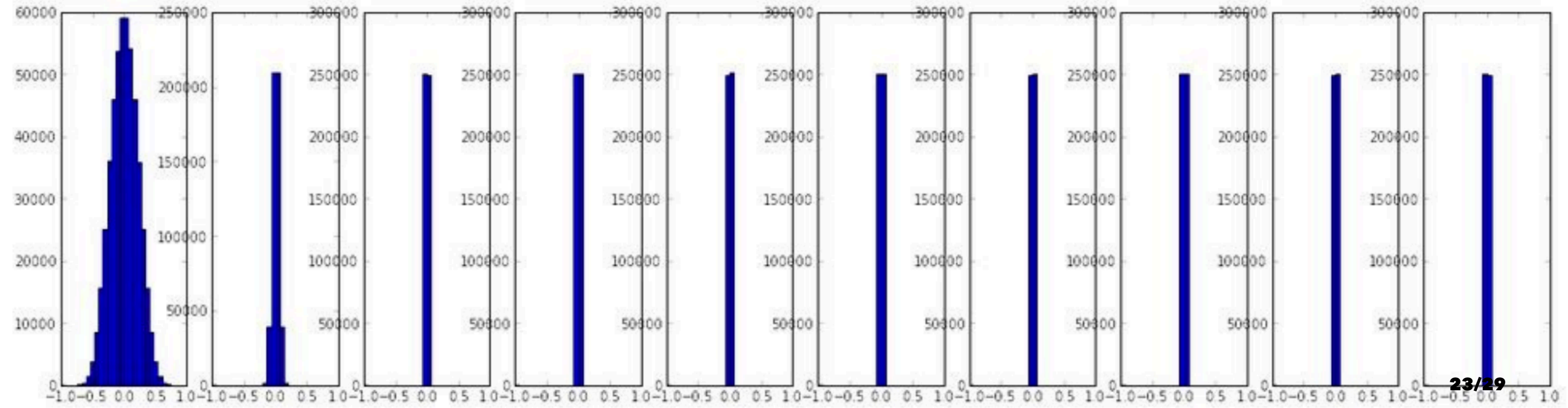
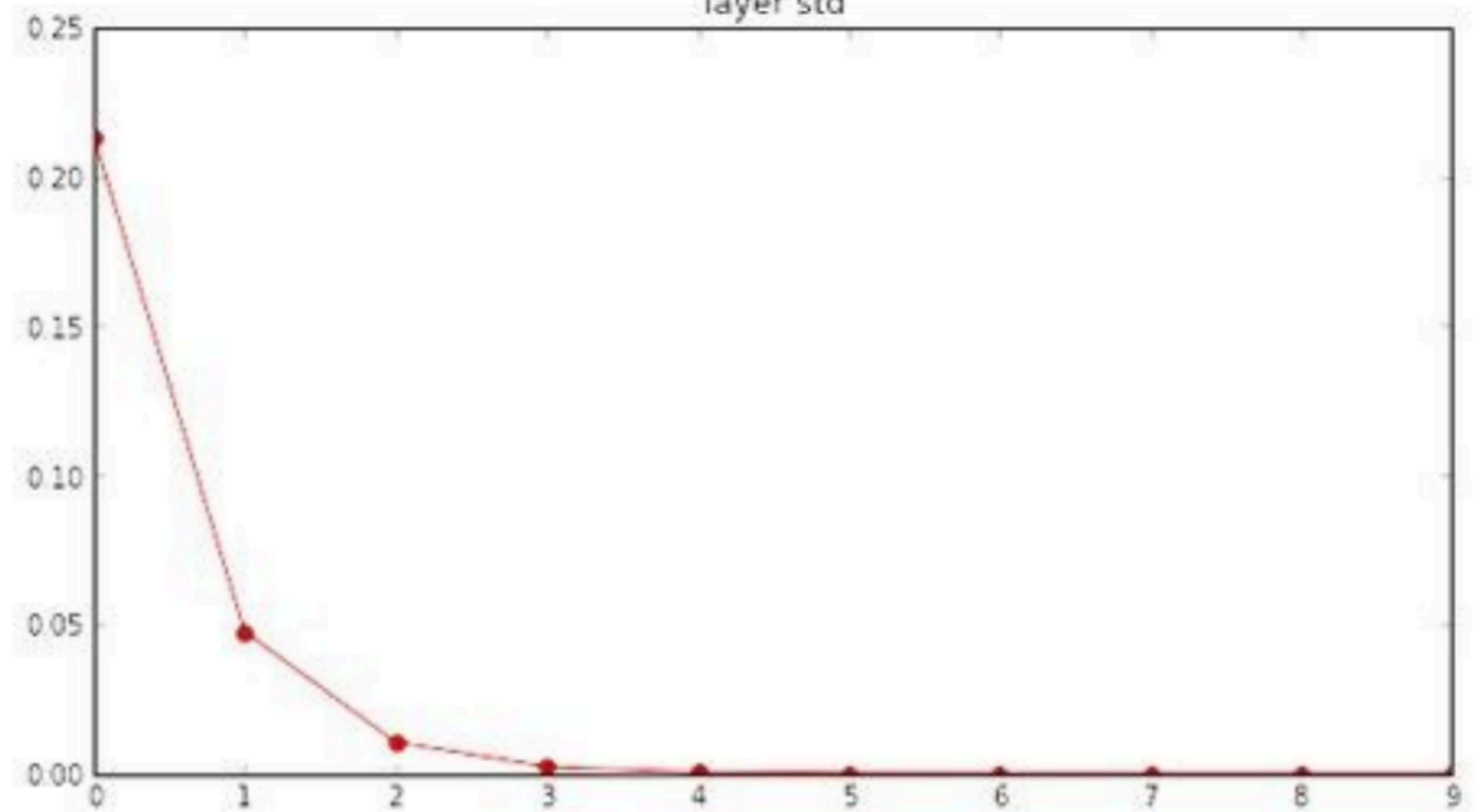
$$W = 0.01 * \text{np.random.randn}(D, H)$$

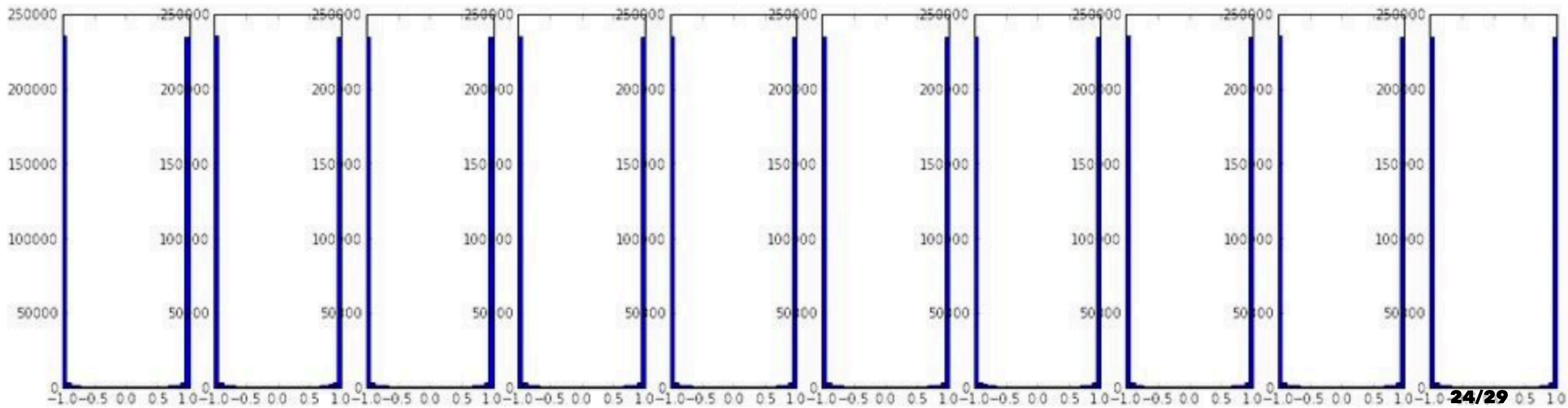
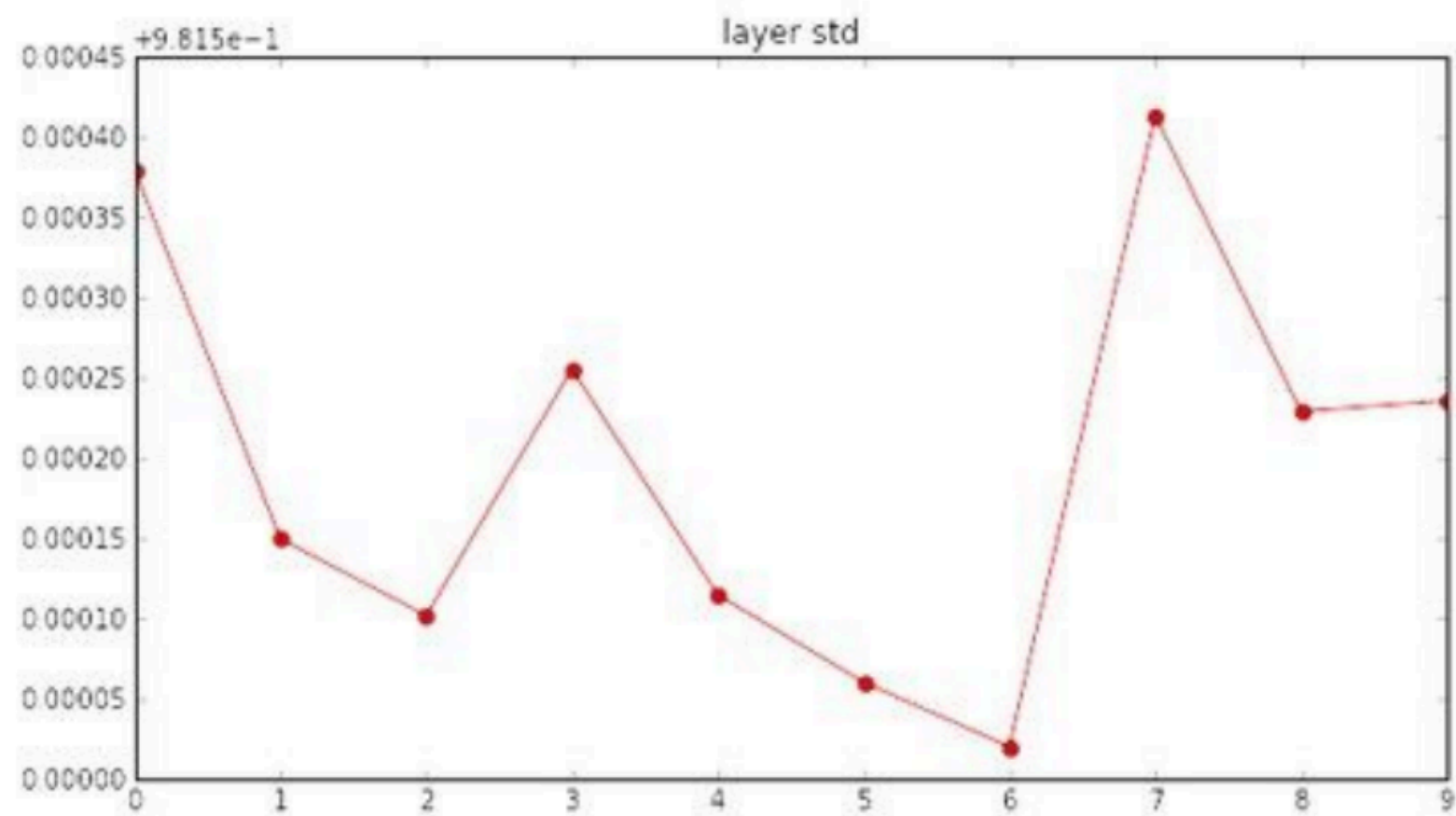
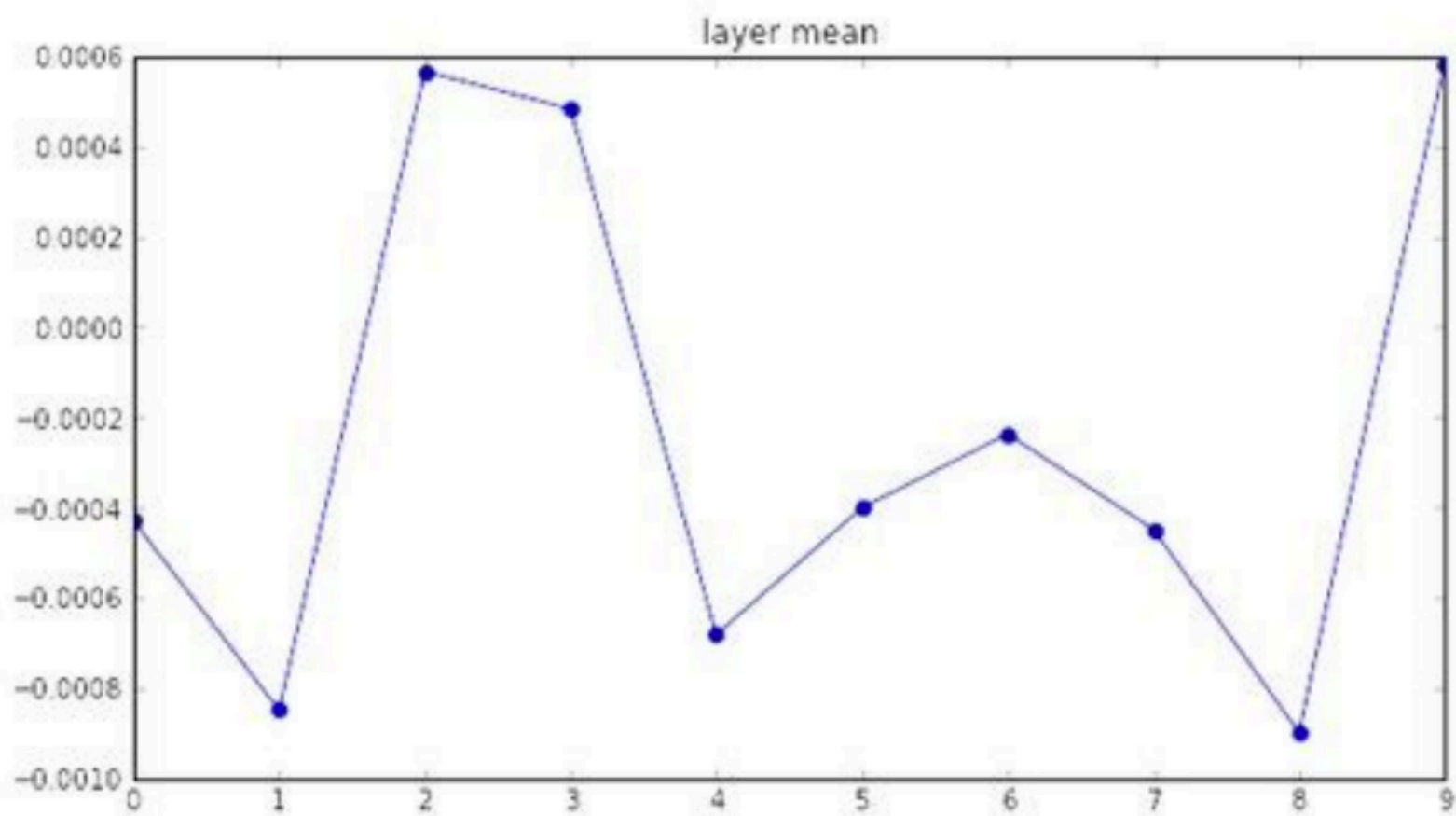
- To see the behavior need activation statistics:
  - mean and variance of values of the weight matrix
  - histogram of values of weight matrix

layer mean



layer std







# Nice weights?

- For  $f_j = W_j \vec{x} = \sum_{i=1}^{n_{\text{in}}} w_{ji} x_i$  we have:  
$$\text{Var}(w_{ji} x_i) = [\mathbb{E}(x_i)]^2 \text{Var}(w_{ji}) + [\mathbb{E}(w_{ji})]^2 \text{Var}(x_i) + \text{Var}(x_i) \text{Var}(w_{ji})$$
- If  $x_i$  and  $w_{ji}$  have mean 0 and are independent random then  $\text{Var}(f_j) = n_{\text{in}} \text{Var}(x_i) \text{Var}(w_{ji})$
- For the variance of input and output to be the same set  $\text{Var}(w_{ji}) = 1/n_{\text{in}}$  for all  $i$  and  $j$

- Gradient to have the same variance:  $\text{Var}(w_{ji}) = 1/n_{\text{out}}$
- Compromise:  $\text{Var}(W) = 2/(n_{\text{in}} + n_{\text{out}})$

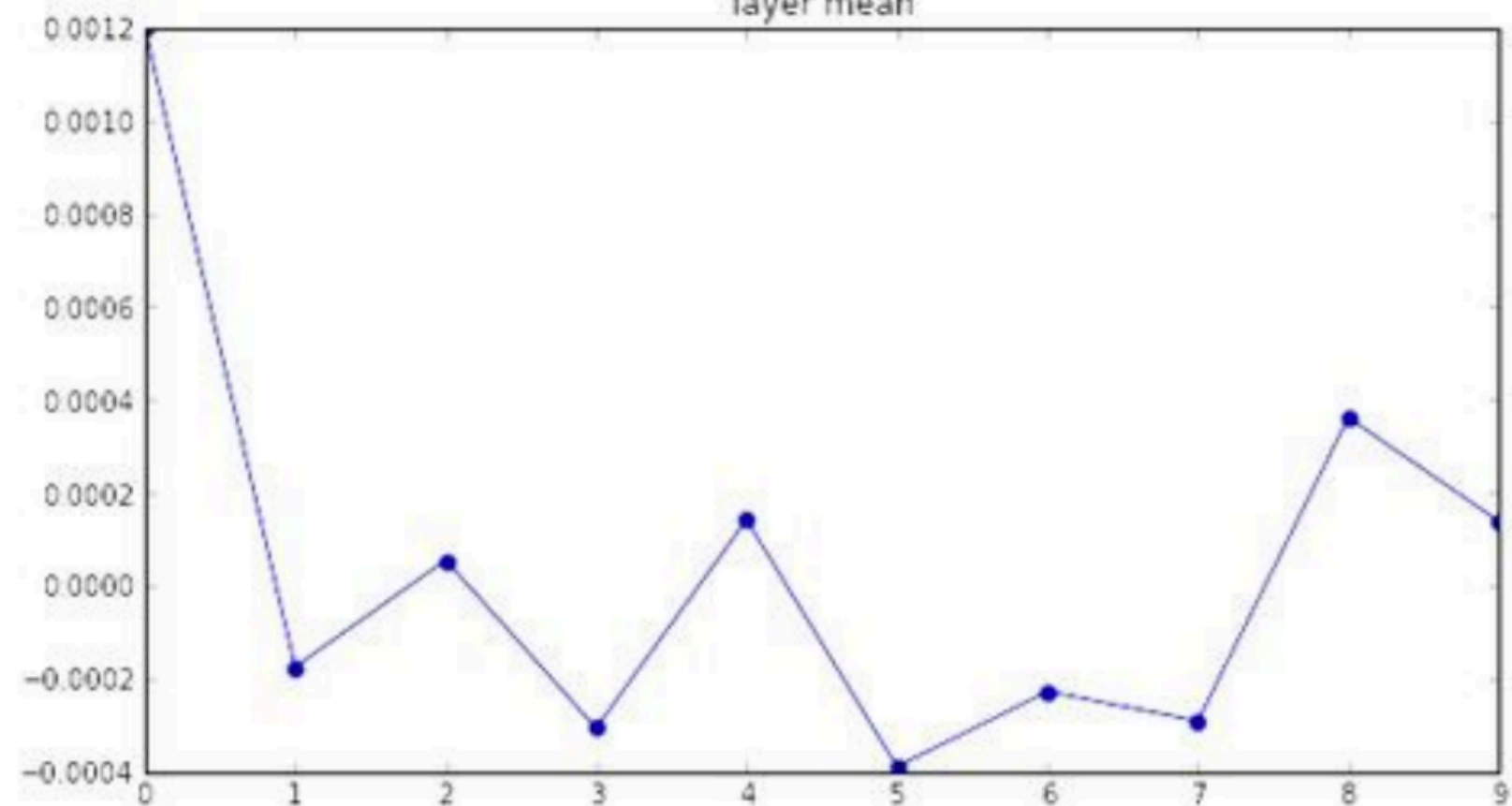
```
from keras import initializers
```

```
model.add(Dense(64,  
                kernel_initializer=initializers.glorot_normal()))
```

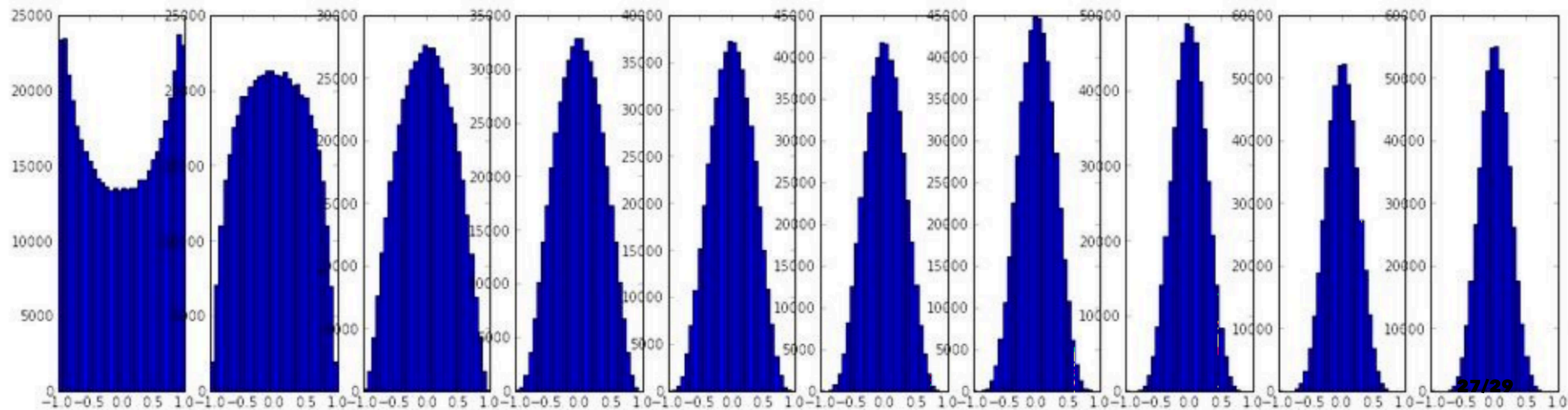
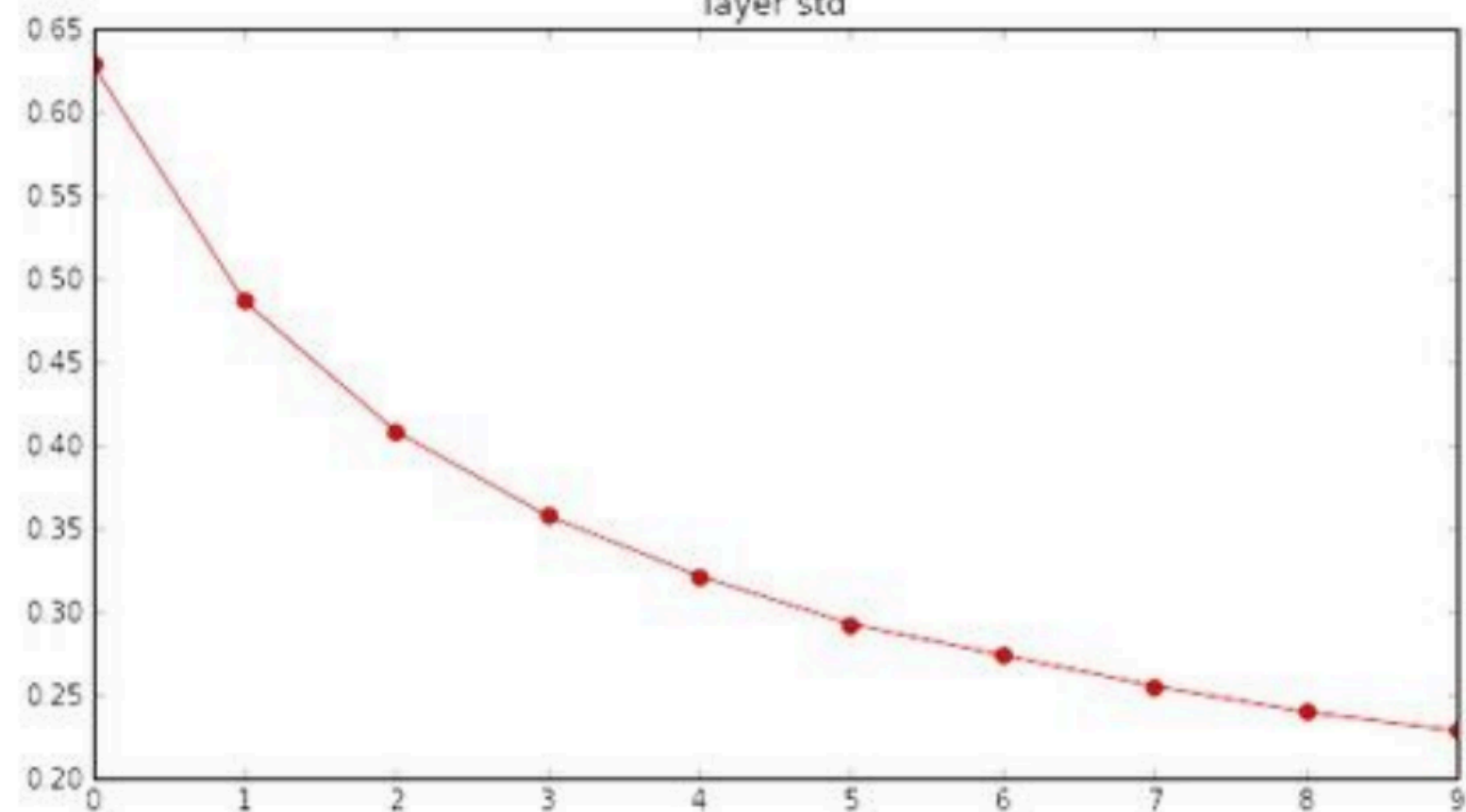
```
# also works; will use the default parameters.
```

```
model.add(Dense(64, kernel_initializer='glorot_normal'))
```

layer mean



layer std



# Delving Deep into Rectifiers

- Instead of identity consider ReLU non-linearities
- Then:  $\text{Var}(f_i) = n\mathbb{E}(x_i^2)\text{Var}(W_i)$  since mean of  $\max(0, x_i)$  is not 0
- So for variance of input and output to be the same we need:  
 $\text{Var}(W) = 2/n_{\text{in}}$

```
model.add(Dense(64, kernel_initializer='he_normal'))
```

# Optimal performance?\*

- May be using the wrong criteria
- The weight properties do not persist as learning starts
- Learning speed increases but inadvertently so does the generalization error
- For large layers (big  $n_{in}$ ), variance of  $2/n_{in}$  means extremely small weights

---

\* Deep Learning Book (2016)