# Neural Network Training

Geena Kim

## Learning in perceptron

- Perceptron rule
- Delta rule (Gradient Descent)

## Perceptron algorithm

- Cycle through the training instances

- Only update $W$ on misclassified instances

- If instance misclassified:

  - If instance is positive class (positive misclassified as negative)

  $$W = W + X_i$$

  - If instance is negative class (negative misclassified as positive)

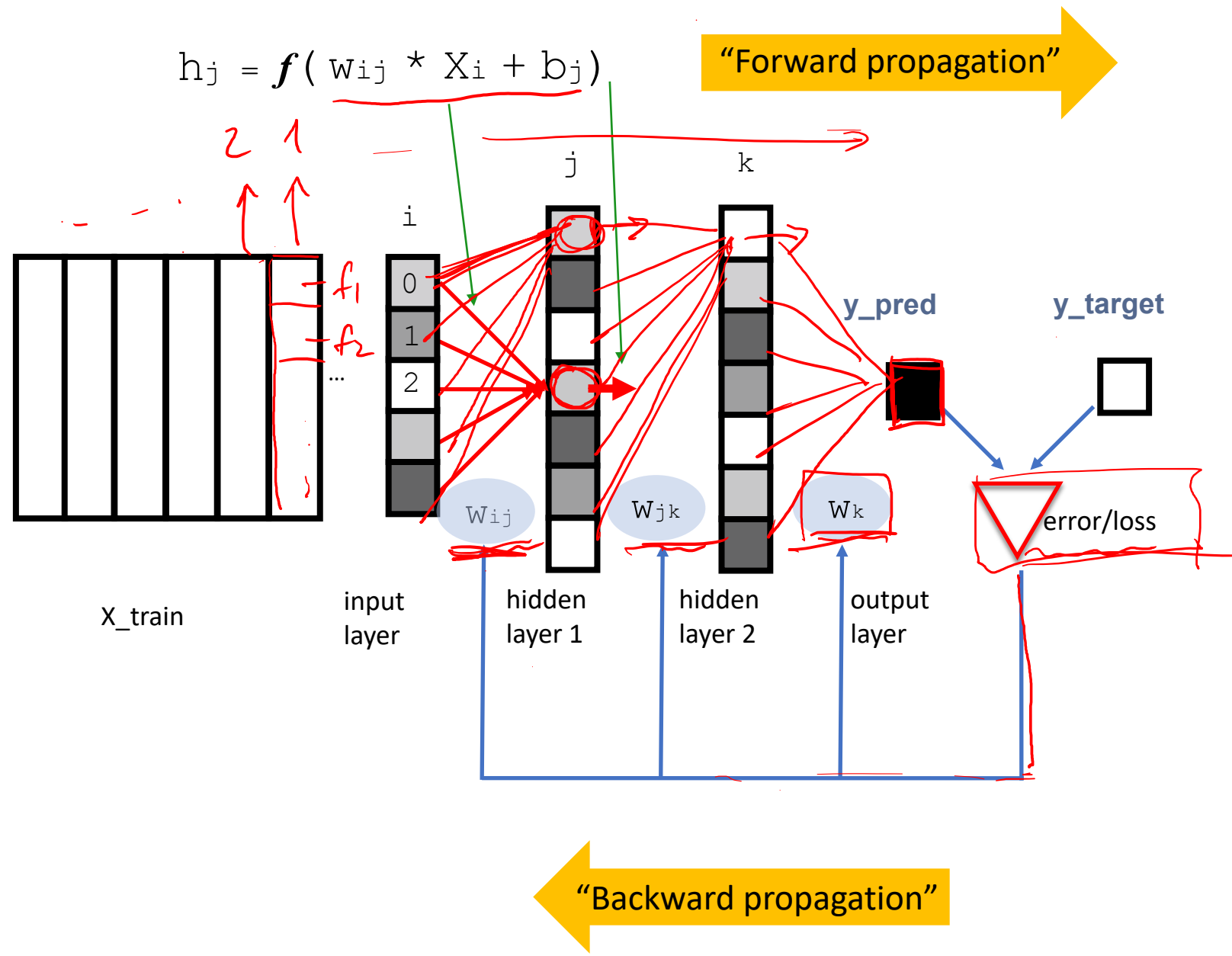  $$W = W - X_i$$

# Perceptron

Delta rule (Gradient Descent)

$$\omega_j \leftarrow \omega_j - \alpha \frac{\partial \mathcal{L}}{\partial \omega_j}$$
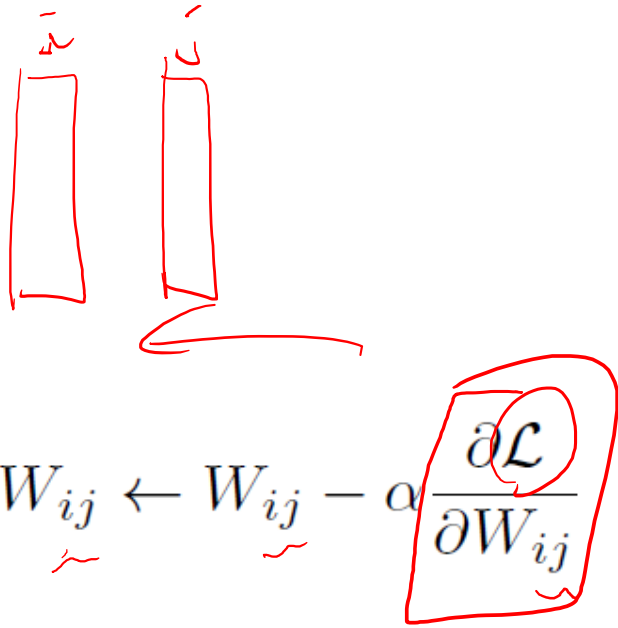
$$\mathcal{L} = \frac{1}{2}(\hat{y}_i - y_i)^2$$

$$\hat{y}_i = \sum_j \omega_j X_{ij}$$

$$\omega_j \leftarrow \omega_j - \alpha(\hat{y}_i - y_i)X_{ij}$$

# How Neural Network Training Works

# Weight Update in deep neural nets

$$W_{ij} \leftarrow W_{ij} - \alpha \frac{\partial \mathcal{L}}{\partial W_{ij}}$$

**Weight Update Rule**

**Loss Function**

**Gradient (Chain Rule)**

**Back Propagation**

# Chain Rule Reminder

$$\frac{\partial f(g(x))}{\partial x} = \frac{\partial f}{\partial g}\frac{\partial g}{\partial x}$$

$$f(g(h(\cdots)))$$

$$\frac{\partial \sigma}{\partial z} = \frac{\partial}{\partial z}\left(\frac{1}{f}\right) \qquad f = \frac{1}{1+e^{-z}}$$

$$+\frac{1}{f^2}e^{-z}\cdot(+1)$$

Example: gradient of sigmoid

$$\frac{\partial \sigma}{\partial w_j} \qquad \sigma(z) = \frac{1}{1+e^{-z}}$$

$$z = \sum_{j=0} w_j \cdot x_j$$

$$\frac{\partial \sigma}{\partial w_j} = \boxed{\frac{\partial \sigma}{\partial z}} \cdot \boxed{\frac{\partial z}{\partial w_j}} \qquad z = \sum_j w_j \cdot x_j$$
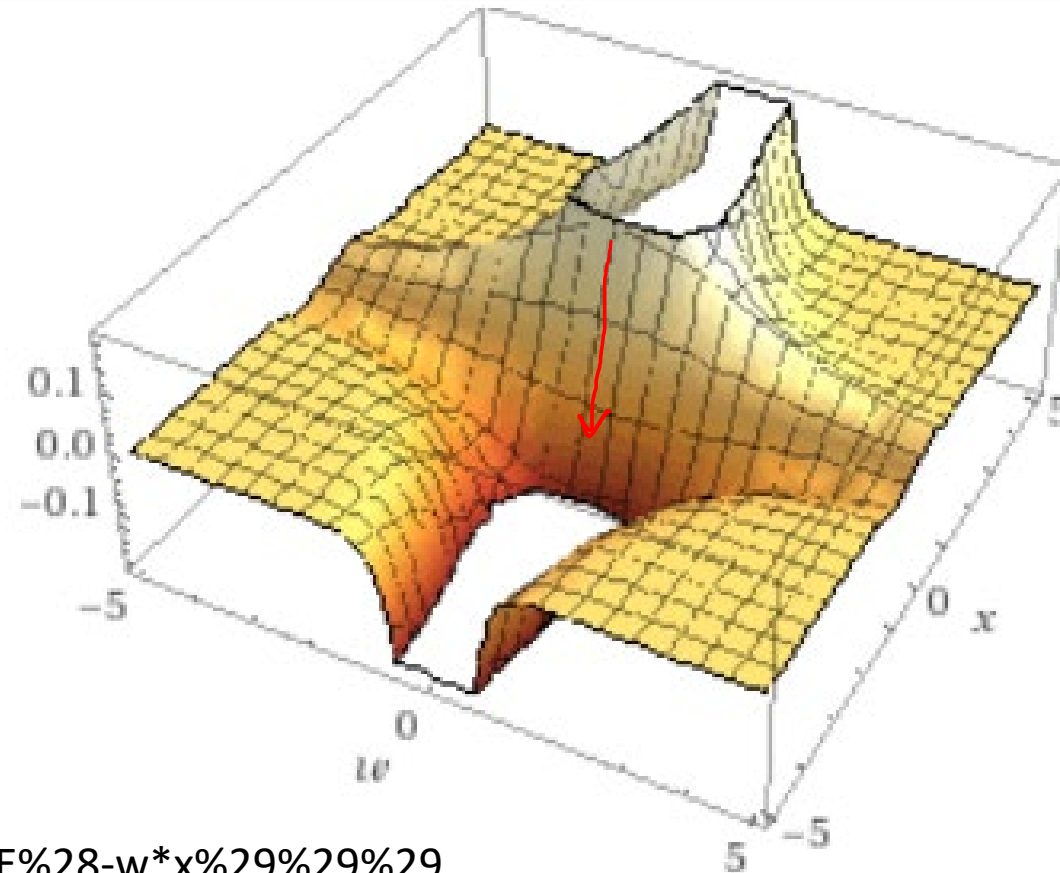
$$\sigma\cdot(1-\sigma) \qquad\qquad x_j$$

$$\frac{\partial \sigma}{\partial w_j} = \sigma\cdot(1-\sigma)\cdot x_j$$

# Chain Rule Reminder



$$\frac{1}{1+e^{-wx}}$$

$$\frac{\partial}{\partial w}\left(\frac{1}{1+e^{-wx}}\right) = \frac{x\,e^{-wx}}{\left(e^{-wx}+1\right)^2}$$

https://www.wolframalpha.com/input/?i=d%2Fdw+%281%2F%281%2Be%5E%28-w*x%29%29%29
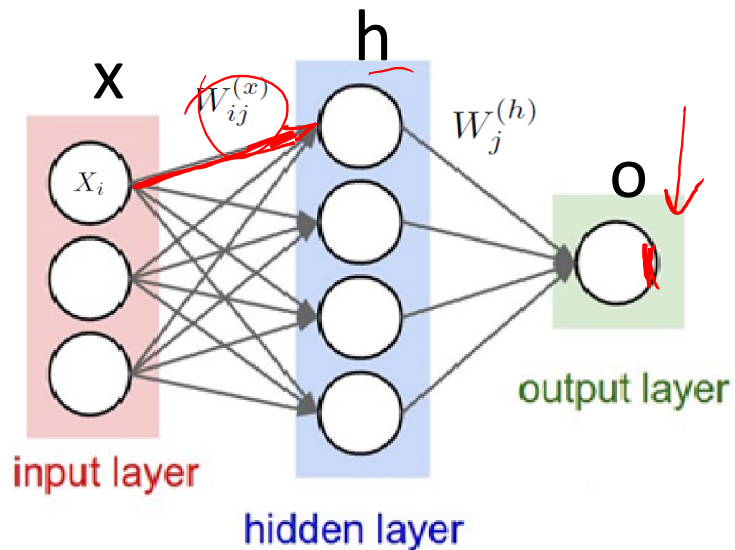
# Calculating Gradient- Chain Rule

$$\mathcal{L} = -\frac{1}{m} \sum_{s}^{m} y_s \log a_s + (1 - y_s) \log(1 - a_s)$$

$$W_{ij} \leftarrow W_{ij} - \alpha \frac{\partial \mathcal{L}}{\partial W_{ij}}$$

$$\frac{\partial f(g(x))}{\partial x} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x}$$

h

x



$W_{ij}^{(x)}$

$X_i$

$W_j^{(h)}$

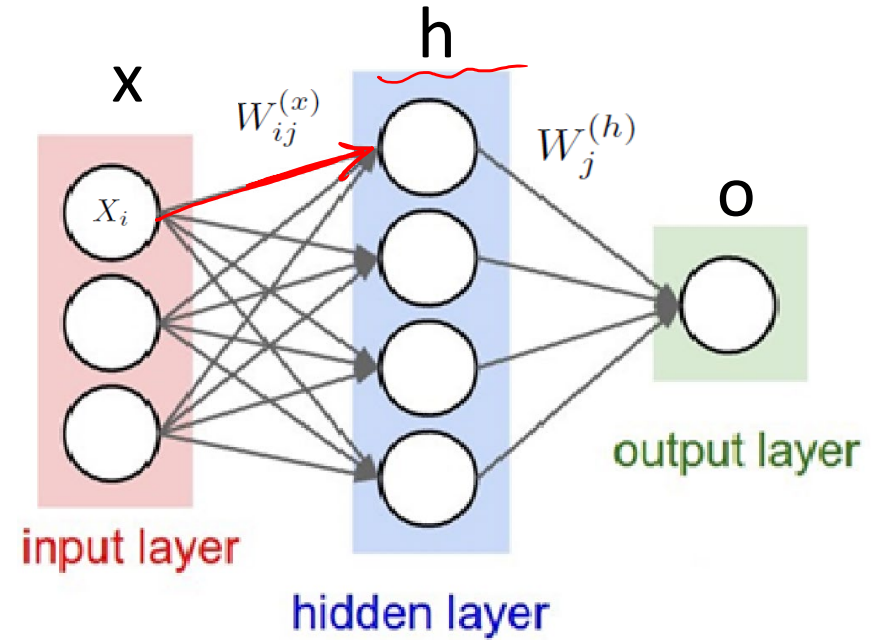o

output layer

input layer

hidden layer

$$\frac{\partial \mathcal{L}}{\partial W_{ij}^{(x)}} = \frac{\partial \mathcal{L}(a^{(o)})}{\partial a^{(o)}} \frac{\partial a^{(o)}}{\partial W_{ij}^{(x)}}$$

$$= \left( \frac{y}{a^{(o)}} - \frac{1 - y}{1 - a^{(o)}} \right) \frac{\partial a^{(o)}}{\partial W_{ij}^{(x)}}$$

$$\frac{\partial a^{(o)}}{\partial W_{ij}^{(x)}} = \frac{\partial \sigma(z^{(o)})}{\partial z^{(o)}} \frac{\partial z^{(h)}}{\partial W_{ij}^{(x)}}$$

$$= \frac{\partial \sigma(z^{(o)})}{\partial z^{(o)}} \frac{\partial z^{(o)}}{\partial W_{ij}^{(x)}} = \sigma(z^{(o)})(1 - \sigma(z^{(o)})) \frac{\partial z^{(o)}}{\partial W_{ij}^{(x)}}$$

# Calculating Gradient- Chain Rule



$$\frac{\partial z^{(o)}}{\partial W_{ij}^{(x)}} = \sum_{j}^{n^{(h)}} W_j^{(h)} \frac{\partial a_j^{(h)}}{\partial W_{ij}^{(x)}}$$

$$= \sum_{j}^{n^{(h)}} W_j^{(h)} \frac{\partial \sigma(z_j^{(h)})}{\partial W_{ij}^{(x)}} = \sum_{j}^{n^{(h)}} W_j^{(h)} \sigma(z_j^{(h)})(1-\sigma(z_j^{(h)})) \frac{\partial(z_j^{(h)})}{\partial W_{ij}^{(x)}} = \sum_{j}^{n^{(h)}} W_j^{(h)} \sigma(z_j^{(h)})(1-\sigma(z_j^{(h)})) X_i$$
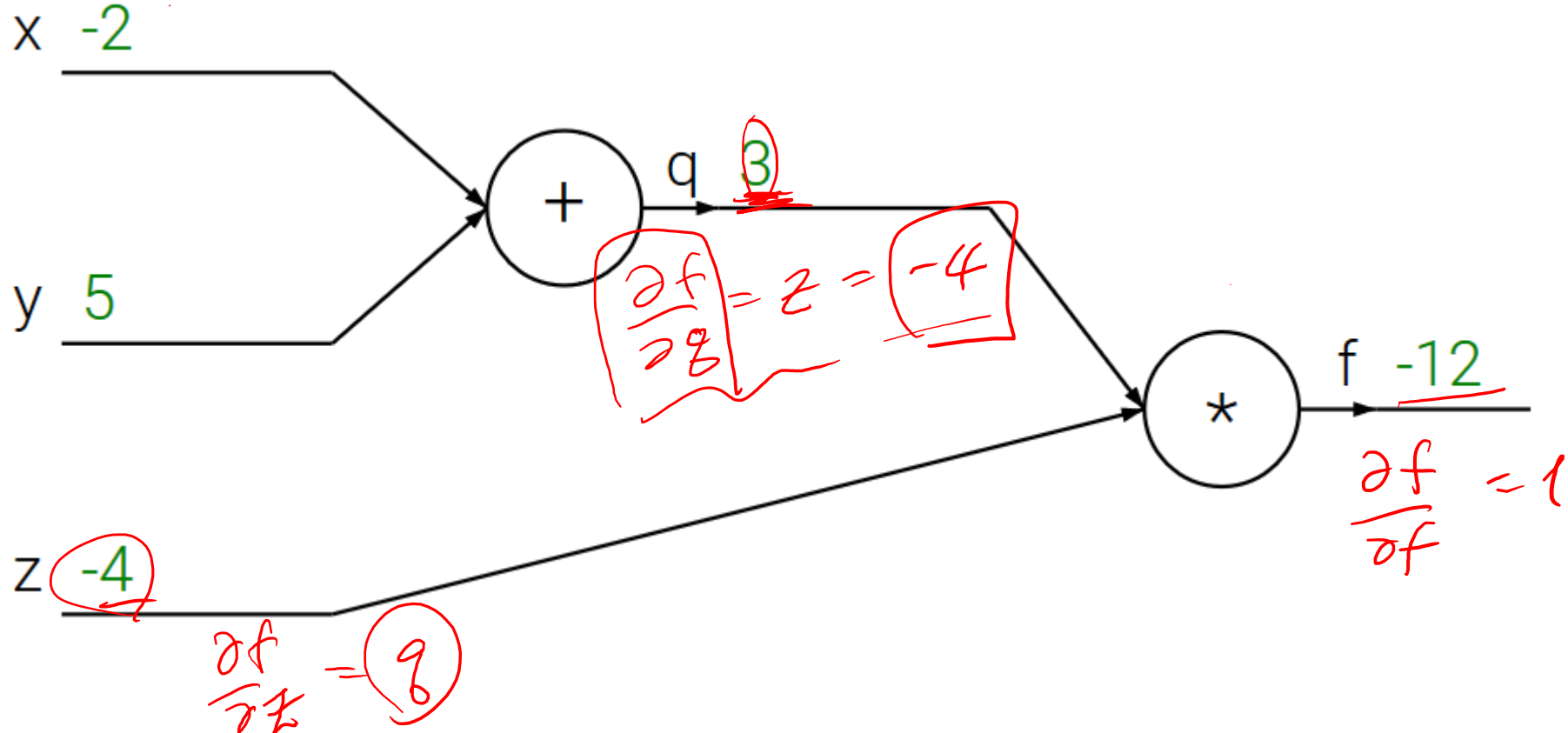
$$\frac{\partial \mathcal{L}}{\partial W_{ij}^{(x)}} = \left( \frac{y}{a^{(o)}} - \frac{1-y}{1-a^{(o)}} \right) \sigma(z^{(o)})(1-\sigma(z^{(o)})) \sum_{j}^{n^{(h)}} W_j^{(h)} \sigma(z_j^{(h)})(1-\sigma(z_j^{(h)})) X_i$$
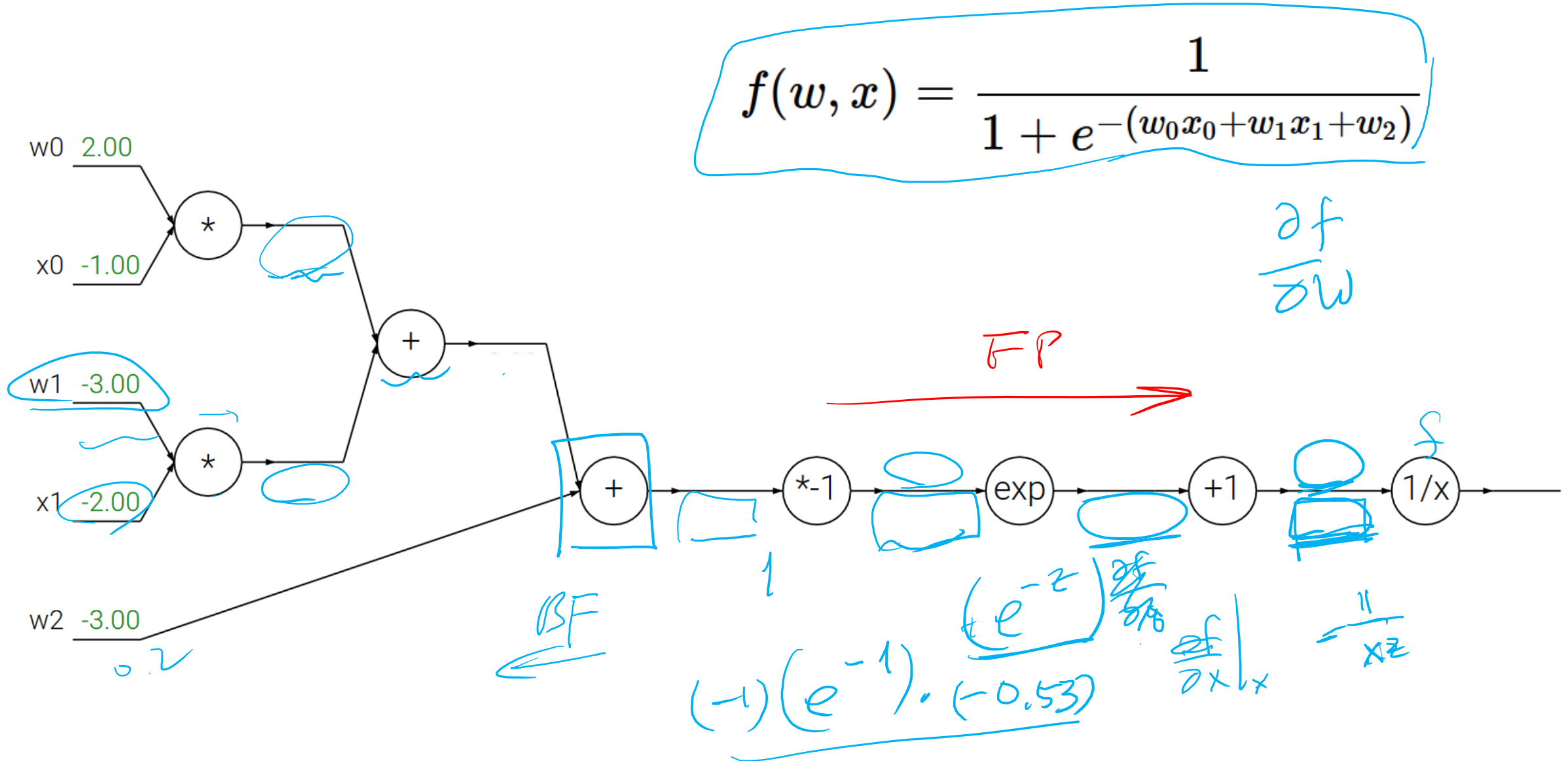
# Back Propagation- Computation Graph

q = x+y

f = q*z

forward pass

x  -2

y  5

+  q  3

$\frac{\partial f}{\partial q} = z = -4$

*  f  -12

$\frac{\partial f}{\partial f} = 1$

z  -4

$\frac{\partial f}{\partial z} = q$

# Back Propagation- Computation Graph



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

w0  2.00
x0  -1.00
w1  -3.00
x1  -2.00
w2  -3.00

*   +   *-1   exp   +1   1/x

$\frac{\partial f}{\partial w}$

FP

BF

1

$(-1)(e^{-1}) \cdot (-0.53)$

$(e^{-z}) \frac{\partial f}{\partial z}$

$\frac{\partial f}{\partial x}\big|_x$

$\frac{1}{xz}$

Image borrowed from Stanford cs231n

# Back Propagation- Computation Graph



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

Image borrowed from Stanford cs231n

# How does the computer perform differentiation?

## Automatic Differentiation (Autodiff)

```
10   defvjp(anp.add,          lambda g, ans, x, y : unbroadcast(x, g),
11                            lambda g, ans, x, y : unbroadcast(y, g))
12   defvjp(anp.multiply,     lambda g, ans, x, y : unbroadcast(x, y * g),
13                            lambda g, ans, x, y : unbroadcast(y, x * g))
14   defvjp(anp.subtract,     lambda g, ans, x, y : unbroadcast(x, g),
15                            lambda g, ans, x, y : unbroadcast(y, -g))
16   defvjp(anp.divide,       lambda g, ans, x, y : unbroadcast(x,    g / y),
17                            lambda g, ans, x, y : unbroadcast(y, - g * x / y**2))
18   defvjp(anp.true_divide,  lambda g, ans, x, y : unbroadcast(x,    g / y),
19                            lambda g, ans, x, y : unbroadcast(y, - g * x / y**2))
20   defvjp(anp.power,
21       lambda g, ans, x, y: unbroadcast(x, g * y * x ** anp.where(y, y - 1, 1.)),
22       lambda g, ans, x, y: unbroadcast(y, g * anp.log(replace_zero(x, 1.)) * x ** y))
```

https://github.com/mattjj/autodidact/blob/master/autograd/numpy/numpy_vjps.py

https://www.cs.toronto.edu/~rgrosse/courses/csc321_2018/slides/lec10.pdf