

Math behind Neural Network Training

Geena Kim



review: Learning in Perceptron

Learning in perceptron

- Perceptron rule
- Delta rule (Gradient Descent)

review: Learning in Perceptron

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)

$$W = W - X_i$$

$$w'_j \leftarrow w_j - 1 \left(\hat{y}_i - y_i \right) X_{ij}$$

$$\begin{array}{l} \hat{y}_i = 1 \\ y_i = -1 \quad \underline{w_j = w_j - 2} \\ \hat{y}_i = -1 \quad \swarrow w_j \\ y_i = +1 \quad \underline{w_j = w_j + 2} \end{array}$$

Perceptron

Delta rule (Gradient Descent)

$$\omega_j \leftarrow \omega_j - \alpha \frac{\partial \mathcal{L}}{\partial \omega_j}$$

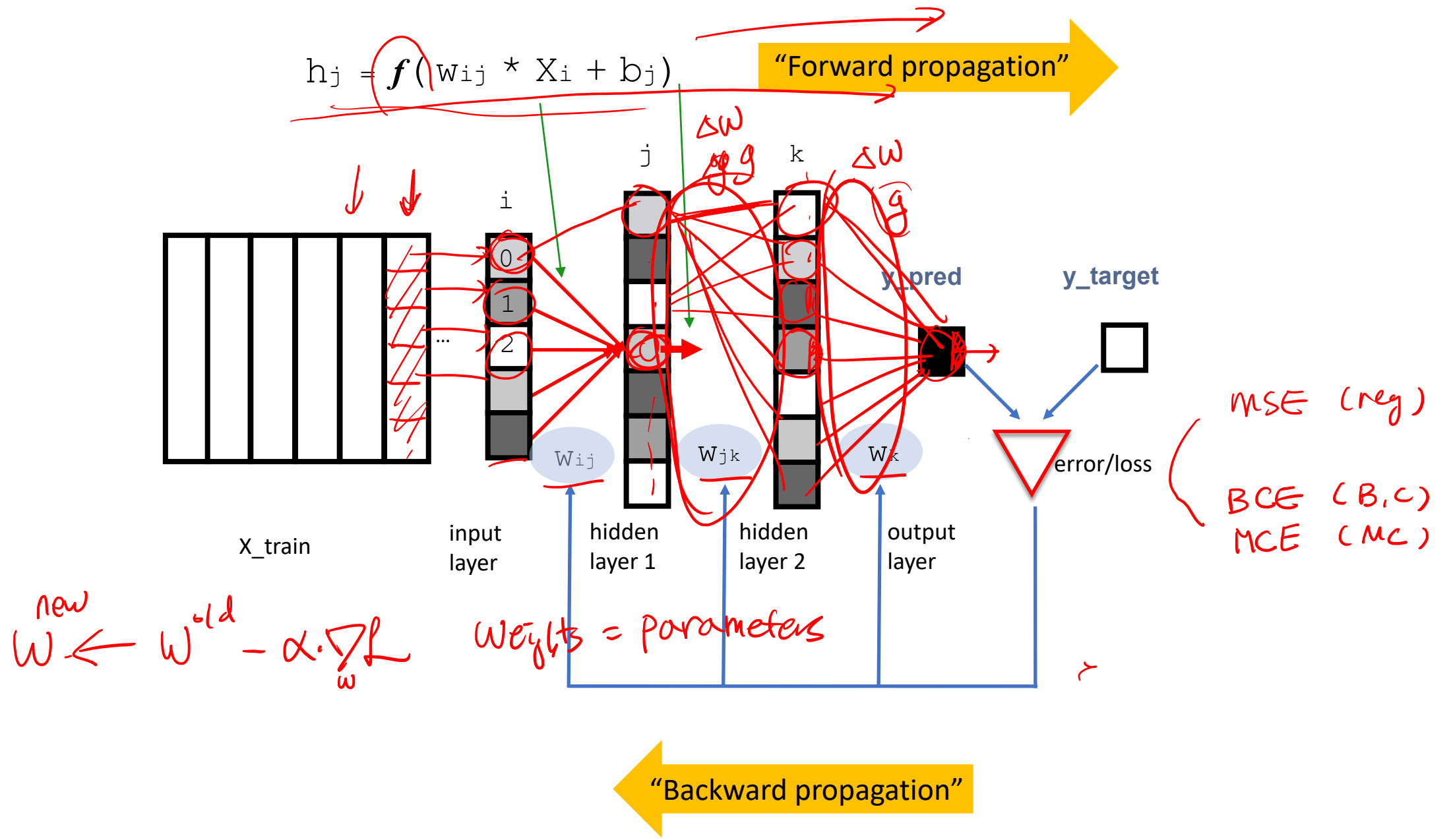
$\mathcal{L}(\omega_j, x_{ji})$
 $MSE = \sum_i \frac{1}{2} (y_i - \hat{y}_i)^2$
 $\frac{\partial}{\partial w} (y - (wx + b))^2$
 $(y - (wx + b)) \cdot (-x)$

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

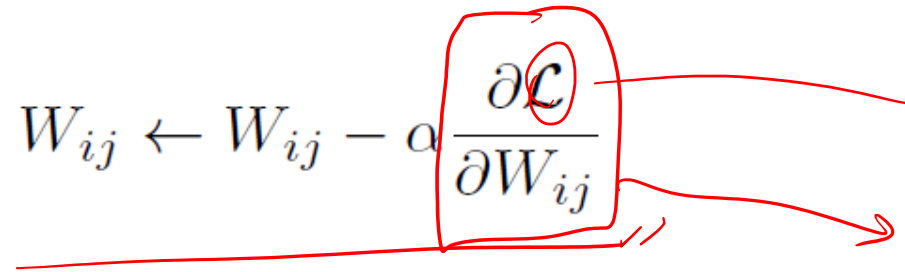
$$\hat{y}_i = \sum_j \omega_j X_{ij}$$

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

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

$$\frac{\partial}{\partial x} \frac{1}{x} = -\frac{1}{x^2}$$


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

Example: gradient of sigmoid

$$\sigma(z) = \frac{1}{1 + e^{-z}} = \frac{1}{g}$$

$z = \sum w_i x_i$

$g = 1 + e^{-z}$



$$\frac{\partial \sigma}{\partial w_i} = + \frac{1}{g^2} \cdot e^{-z} \cdot 1 \cdot x_i$$

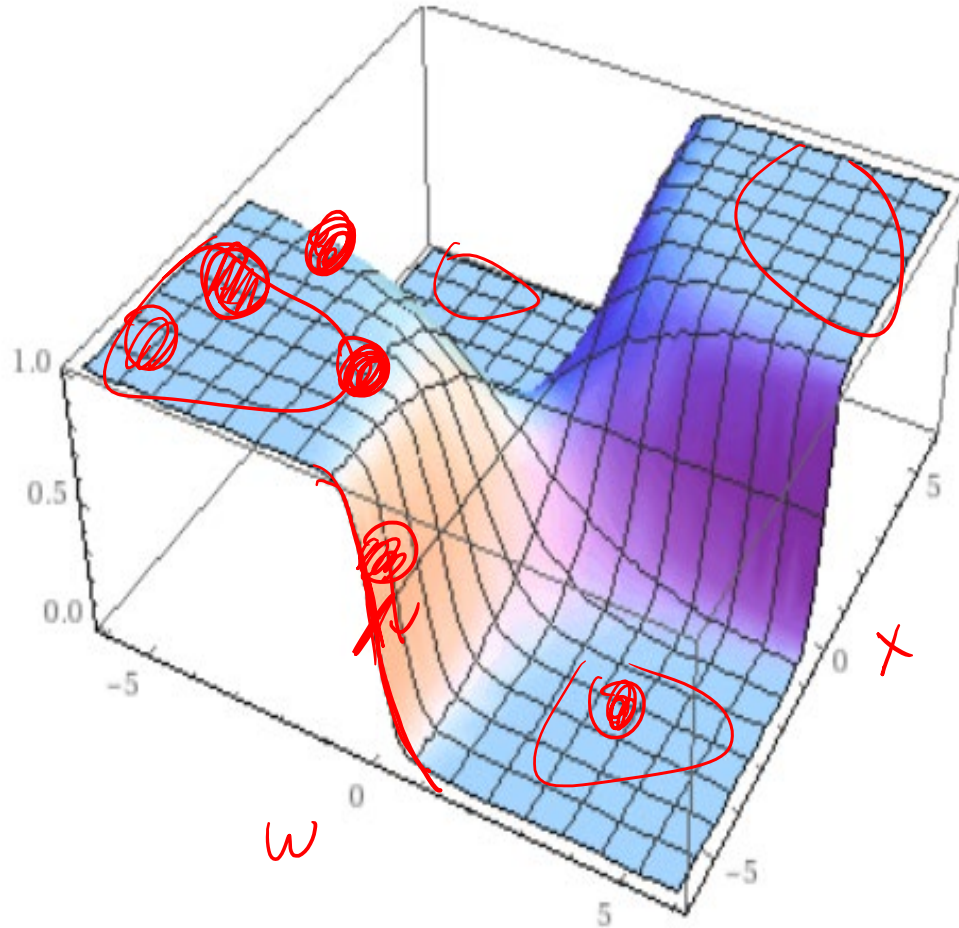
$\frac{1}{1 + e^{-z}}$ (labeled σ)
 $\frac{e^{-z}}{1 + e^{-z}}$ (labeled $1 - \sigma$)
 $\cdot x_i$

$\sigma \cdot (1 - \sigma) \cdot x_i$

Chain Rule Reminder

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

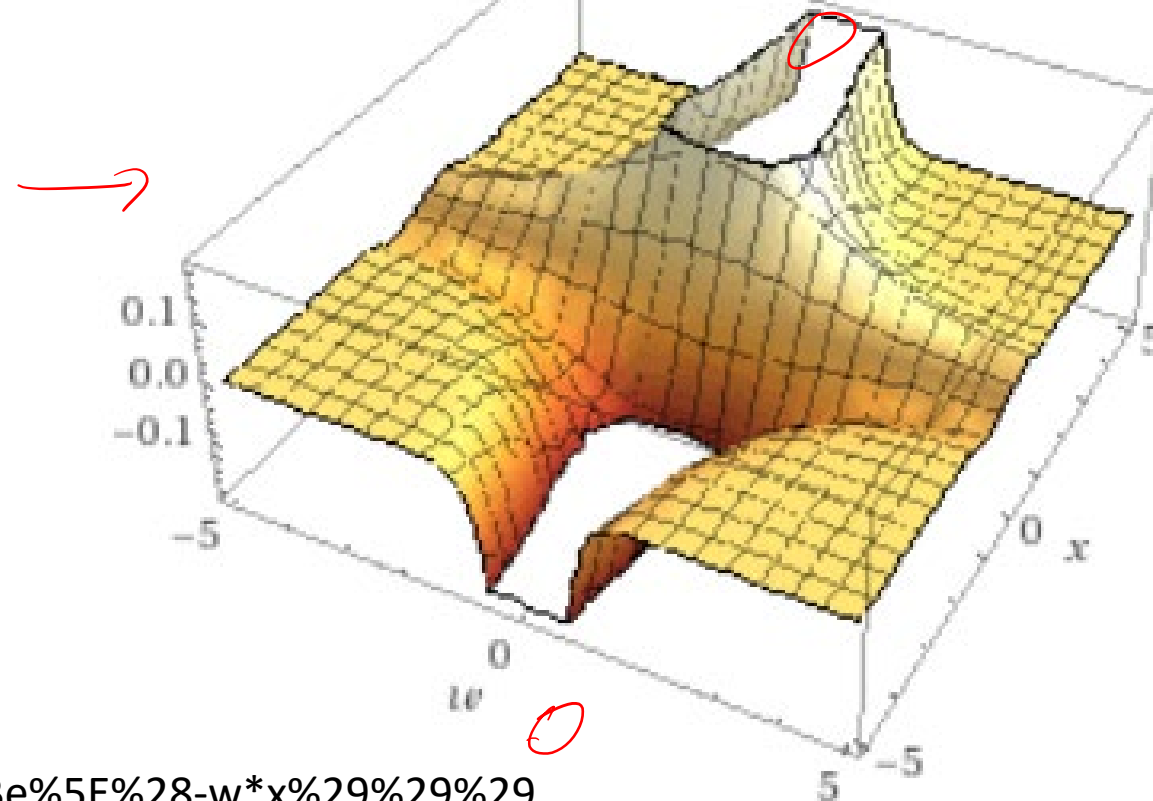
Handwritten red annotations: a red box around the fraction, a red arrow pointing from the box to the right, and the handwritten text $w x$ below the denominator.



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

$$w \leftarrow w - \alpha \nabla L$$

Handwritten red arrow pointing from this equation to the second 3D plot.



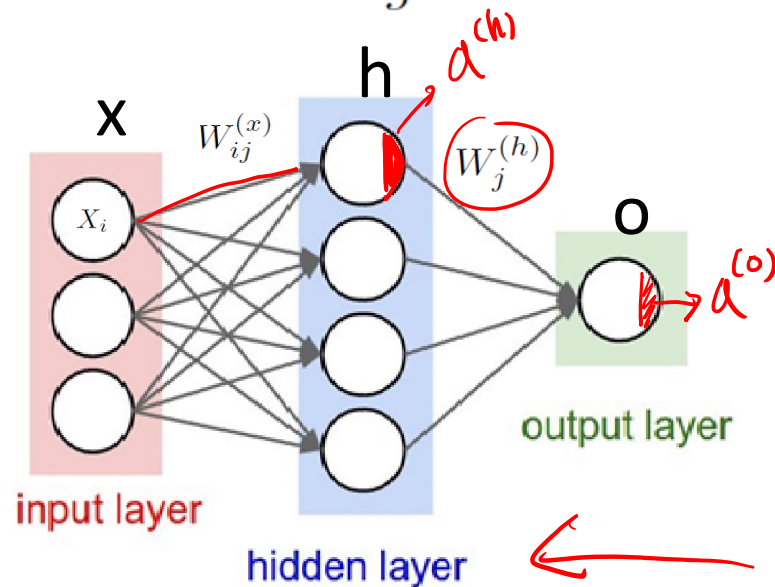
Calculating Gradient- Chain Rule

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

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

$$\frac{\partial \log x}{\partial x} = \frac{1}{x}$$

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



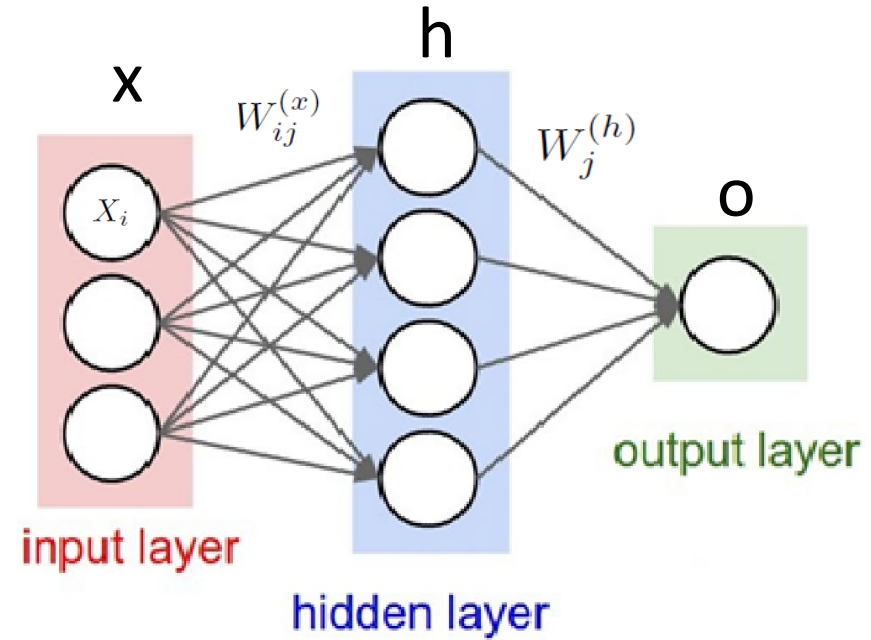
$$\frac{\partial \mathcal{L}}{\partial W_{ij}^{(x)}} = \frac{\partial \mathcal{L}(a^{(o)})}{\partial a^{(o)}} \frac{\partial a^{(o)}}{\partial W_{ij}^{(x)}}$$

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

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

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

Calculating Gradient- Chain Rule



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

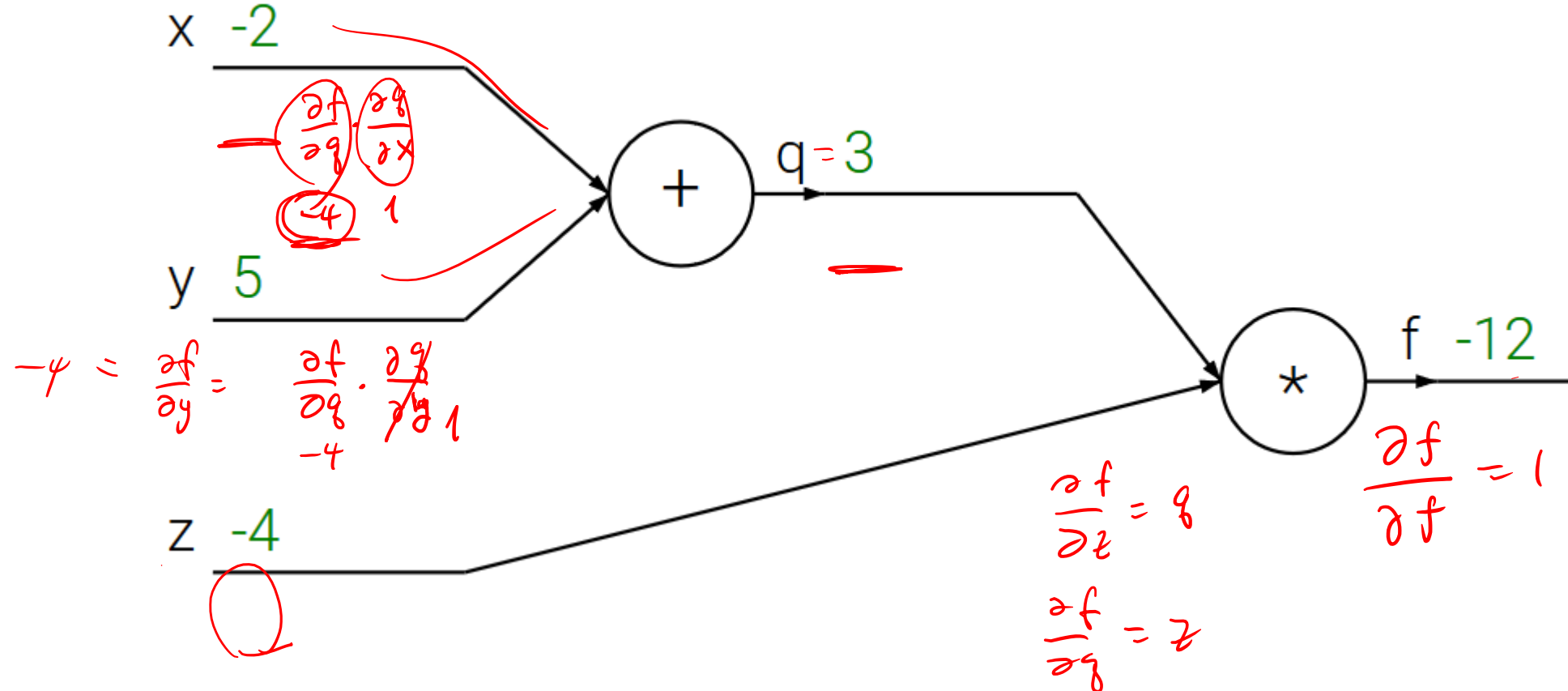
$$= \sum_j^{n^{(h)}} W_j^{(h)} \frac{\partial \sigma(z_j^{(h)})}{\partial W_{ij}^{(x)}} = \sum_j^{n^{(h)}} W_j^{(h)} \sigma(z_j^{(h)}) (1 - \sigma(z_j^{(h)})) \frac{\partial (z_j^{(h)})}{\partial W_{ij}^{(x)}} = \sum_j^{n^{(h)}} W_j^{(h)} \sigma(z_j^{(h)}) (1 - \sigma(z_j^{(h)})) X_i$$

$$\frac{\partial \mathcal{L}}{\partial W_{ij}^{(x)}} = \left(\frac{y}{a^{(o)}} - \frac{1-y}{1-a^{(o)}} \right) \sigma(z^{(o)}) (1 - \sigma(z^{(o)})) \sum_j^{n^{(h)}} W_j^{(h)} \sigma(z_j^{(h)}) (1 - \sigma(z_j^{(h)})) X_i$$

Back Propagation- Computation Graph

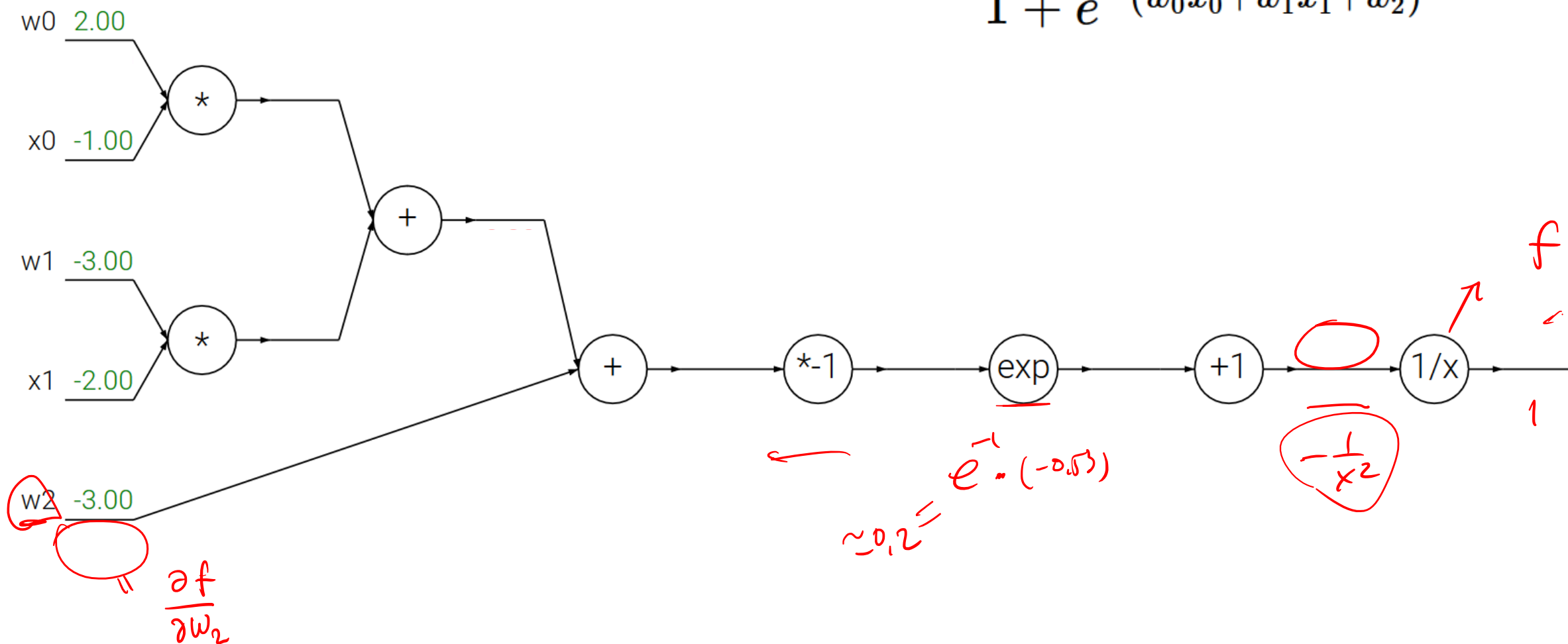
$$q = x + y$$

$$f = q * z$$

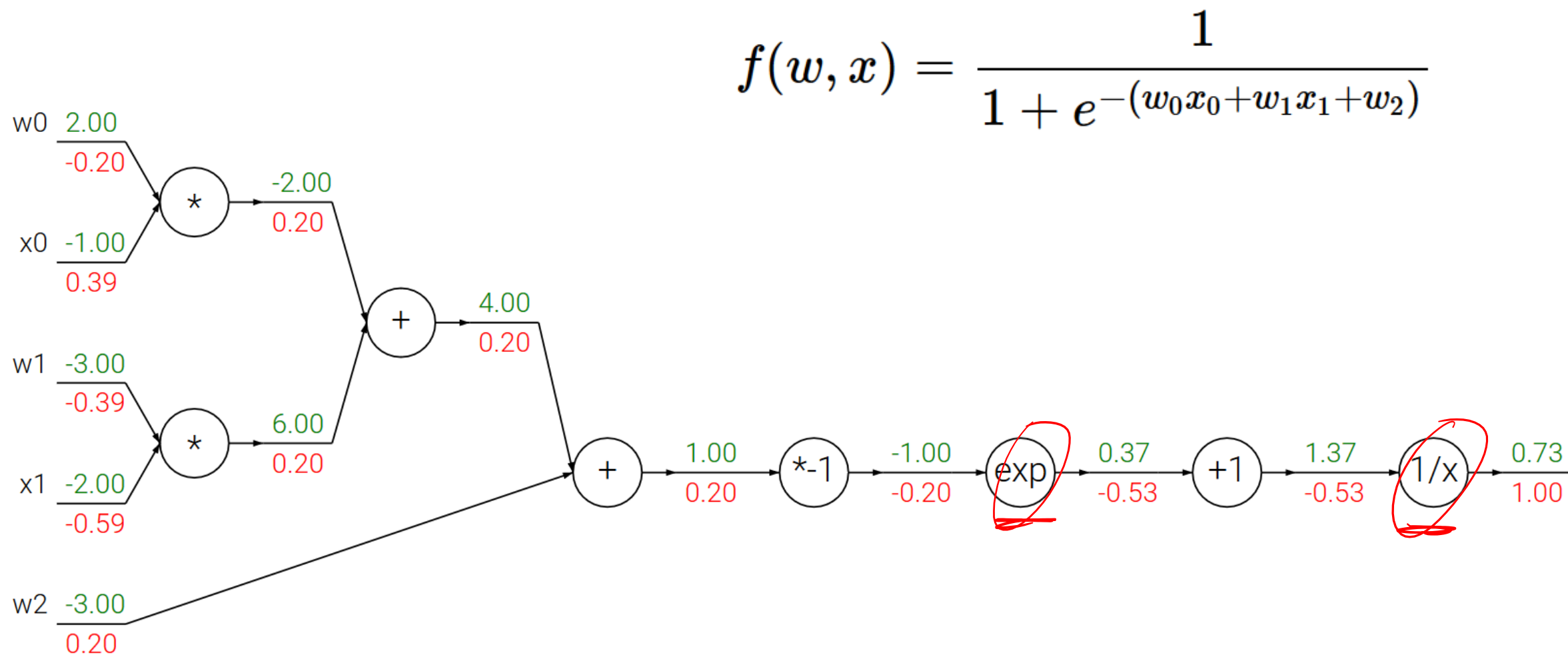


Back Propagation- Computation Graph

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$



Back Propagation- Computation Graph



How does the computer perform differentiation?

Automatic Differentiation (Autodiff)

$$f(x,y) \left(\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right)$$

```
10  def vjp(anp.add,          lambda g, ans, x, y : unbroadcast(x, g),
11                                     lambda g, ans, x, y : unbroadcast(y, g))
12  def vjp(anp.multiply,    lambda g, ans, x, y : unbroadcast(x, y * g),
13                                     lambda g, ans, x, y : unbroadcast(y, x * g))
14  def vjp(anp.subtract,    lambda g, ans, x, y : unbroadcast(x, g),
15                                     lambda g, ans, x, y : unbroadcast(y, -g))
16  def vjp(anp.divide,       lambda g, ans, x, y : unbroadcast(x,  g / y),
17                                     lambda g, ans, x, y : unbroadcast(y, - g * x / y**2))
18  def vjp(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  def vjp(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

Gradient Descent

Optimization Goal

Find a set of (optimized) weights which minimize the error (or loss function) at the output

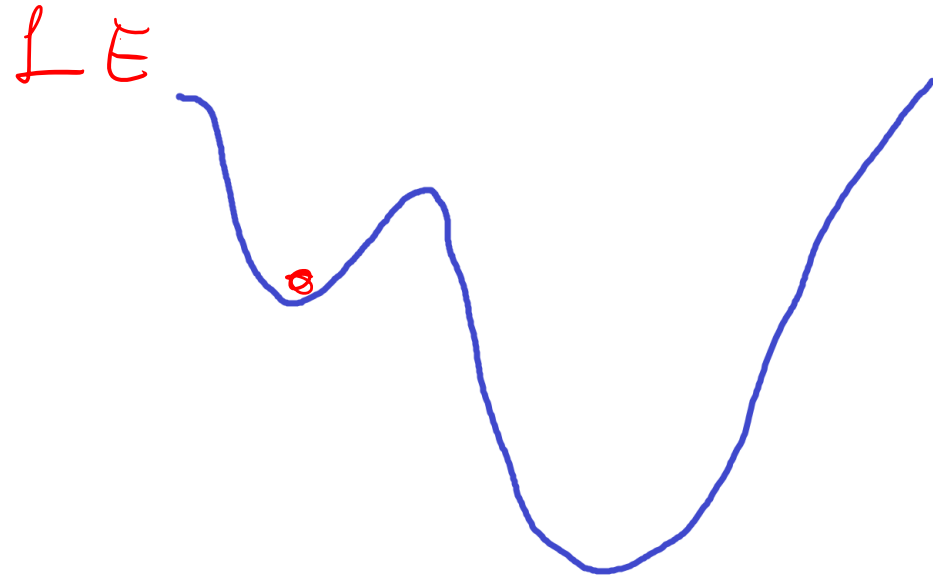
Weight update rule

$$\dot{W}_{nm}^L \leftarrow W_{nm}^L - \alpha * \delta W_{nm}^L$$

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

ΔW

Global minimum vs. local minimum

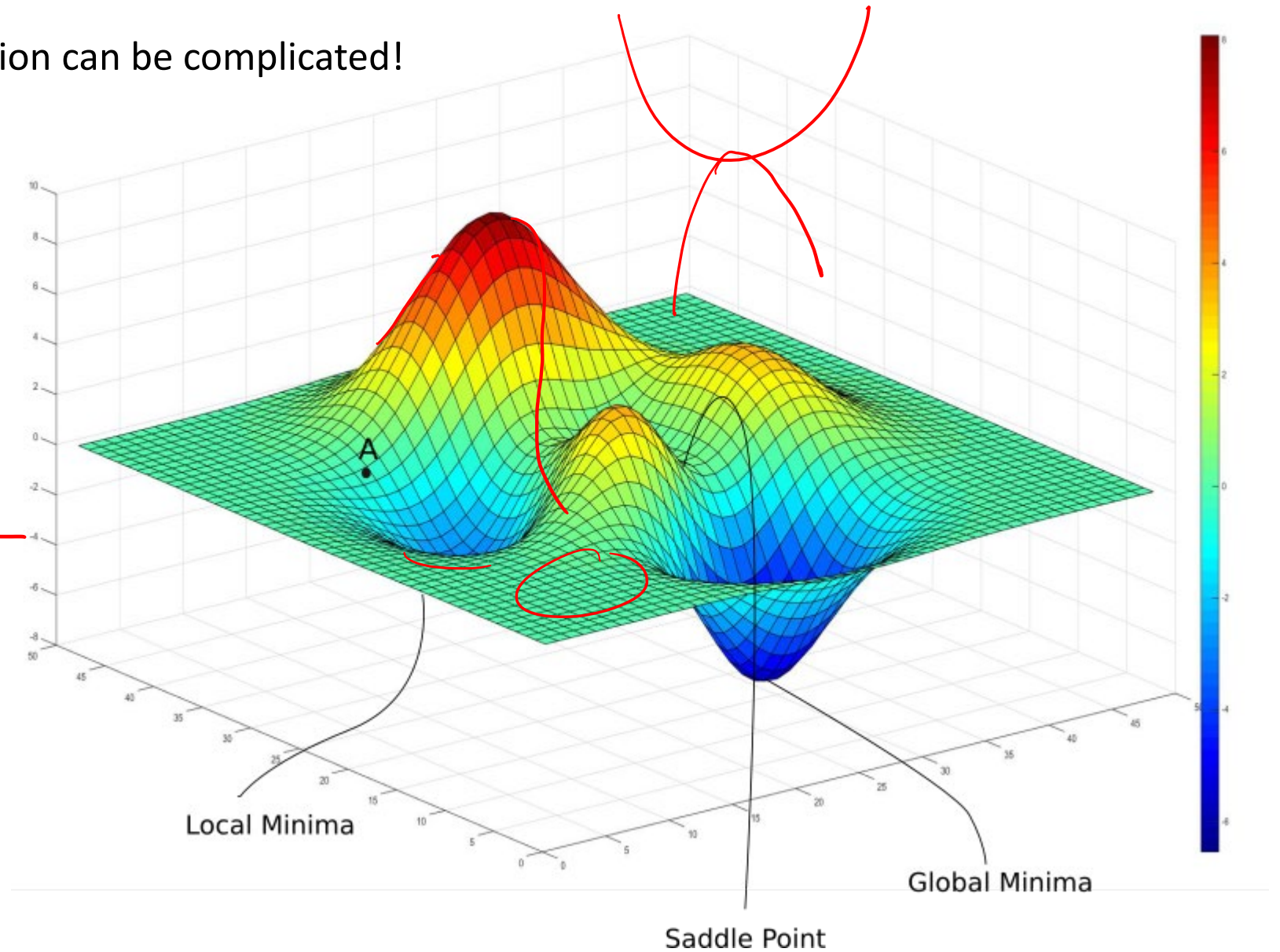


Gradient Descent

Error surface in the multi dimension can be complicated!

WT W-~~2L~~

also... cliffs and plateaus.



Stochastic Gradient Descent

How many training samples at a time do we include to calculate the error?

W

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

Practically we use mini batches

1 ep

show 10,000

1 ep =

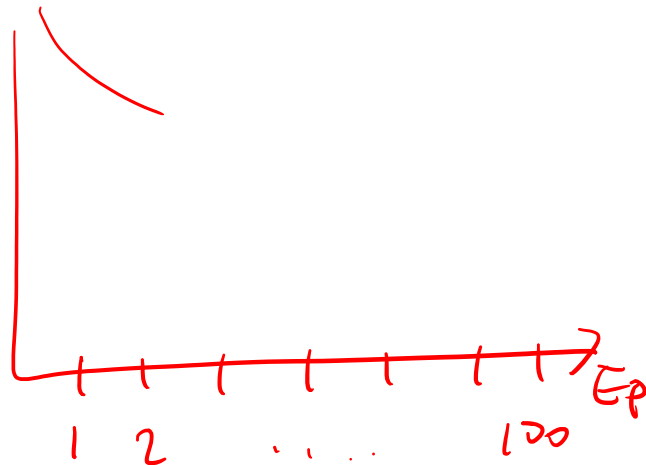
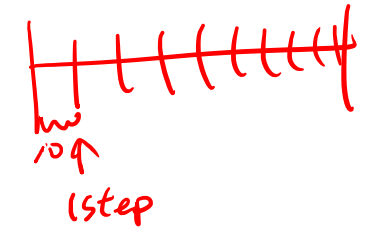
minibatch bs=10

1000 batches

→ 1 step



1000



Training speed and accuracy vs. minibatch size

Stochastic Gradient Descent

With decreasing learning rate (Learning rate scheduling)

Algorithm 8.1 Stochastic gradient descent (SGD) update

Require: Learning rate schedule $\epsilon_1, \epsilon_2, \dots$

Require: Initial parameter θ

$k \leftarrow 1$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Apply update: $\theta \leftarrow \theta - \epsilon_k \hat{\mathbf{g}}$

$k \leftarrow k + 1$

end while

Stochastic Gradient Descent with momentum

SGD with learning rate alone is slow to converge

Adding a momentum (moving average) can make it faster

$$\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \nabla_{\boldsymbol{\theta}} \left(\frac{1}{m} \sum_{i=1}^m L(\mathbf{f}(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)}) \right),$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mathbf{v}.$$

warning- different notations used (from deeplearningbook.org)

```
tf.keras.optimizers.SGD(  
    learning_rate=0.01, momentum=0.0, nesterov=False, name='SGD', **kwargs  
)
```

Next time

- More optimization methods
- Regularization techniques
- Tips for training NN