

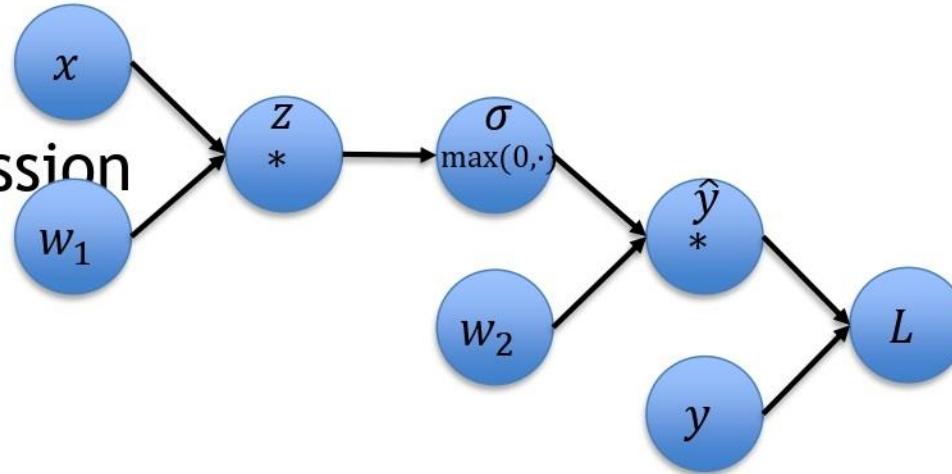
# Lecture 6

# Optimization

# Back to Compute Graphs & NNs

- Inputs  $x$  and targets  $y$
- Two-layer NN for regression with ReLU activation
- Function we want to optimize:

$$\sum_{i=1}^n \|w_2 \max(0, w_1 x_i) - y_i\|_2^2$$



# Gradient Descent for Neural Networks

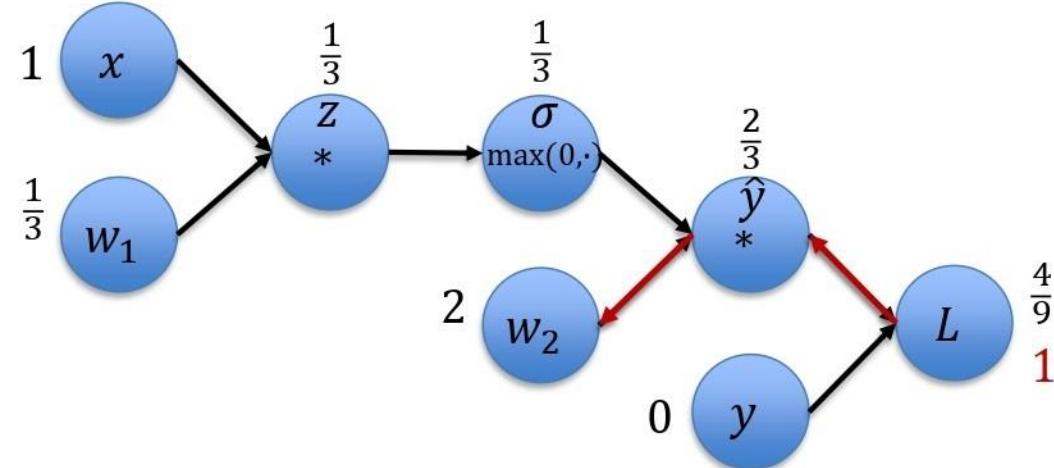
Initialize  $x = 1$ ,  $y = 0$ ,  
 $w_1 = \frac{1}{3}$ ,  $w_2 = 2$

$$L(\mathbf{y}, \hat{\mathbf{y}}; \boldsymbol{\theta}) = \frac{1}{n} \sum_i^n ||\hat{y}_i - y_i||_2^2$$

In our case  $n, d = 1$ :

$$L = (\hat{y} - y)^2 \Rightarrow \frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)$$

$$\hat{y} = w_2 \cdot \sigma \quad \Rightarrow \frac{\partial \hat{y}}{\partial w_2} = \sigma$$



Backpropagation

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2}$$

# Gradient Descent for Neural Networks

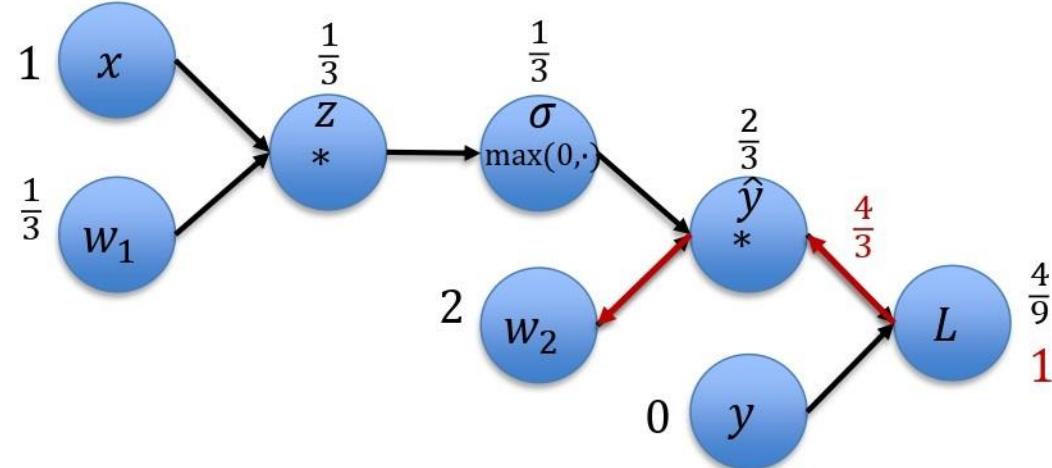
Initialize  $x = 1$ ,  $y = 0$ ,  
 $w_1 = \frac{1}{3}$ ,  $w_2 = 2$

$$L(\mathbf{y}, \hat{\mathbf{y}}; \boldsymbol{\theta}) = \frac{1}{n} \sum_i^n \|\hat{y}_i - y_i\|_2^2$$

In our case  $n, d = 1$ :

$$L = (\hat{y} - y)^2 \Rightarrow \boxed{\frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)}$$

$$\hat{y} = w_2 \cdot \sigma \quad \Rightarrow \frac{\partial \hat{y}}{\partial w_2} = \sigma$$



Backpropagation

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2}$$

$$2 \cdot \frac{2}{3}$$

# Gradient Descent for Neural Networks

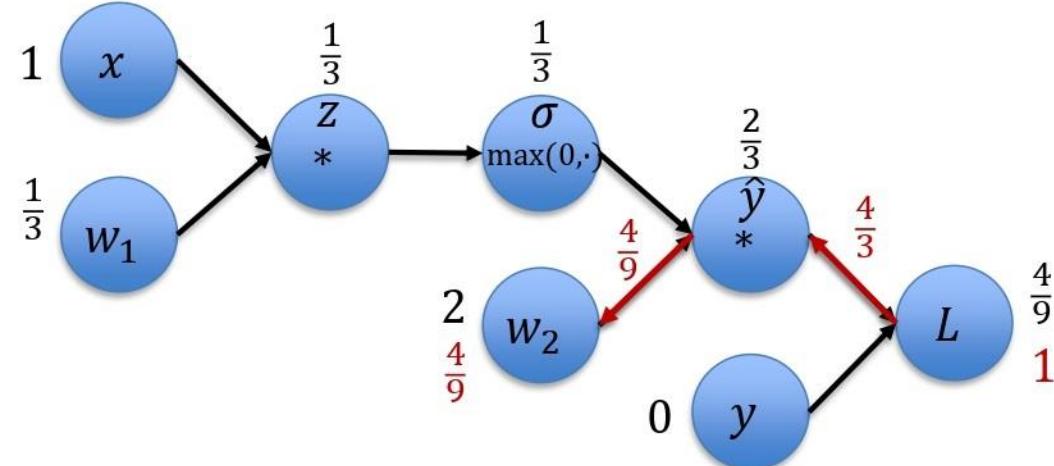
Initialize  $x = 1$ ,  $y = 0$ ,  
 $w_1 = \frac{1}{3}$ ,  $w_2 = 2$

$$L(\mathbf{y}, \hat{\mathbf{y}}; \boldsymbol{\theta}) = \frac{1}{n} \sum_i^n \|\hat{y}_i - y_i\|_2^2$$

In our case  $n, d = 1$ :

$$L = (\hat{y} - y)^2 \Rightarrow \frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)$$

$$\hat{y} = w_2 \cdot \sigma \quad \Rightarrow \boxed{\frac{\partial \hat{y}}{\partial w_2} = \sigma}$$



Backpropagation

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2}$$

$$2 \cdot \frac{2}{3} \cdot \frac{1}{3}$$

# Gradient Descent for Neural Networks

Initialize  $x = 1$ ,  $y = 0$ ,  
 $w_1 = \frac{1}{3}$ ,  $w_2 = 2$

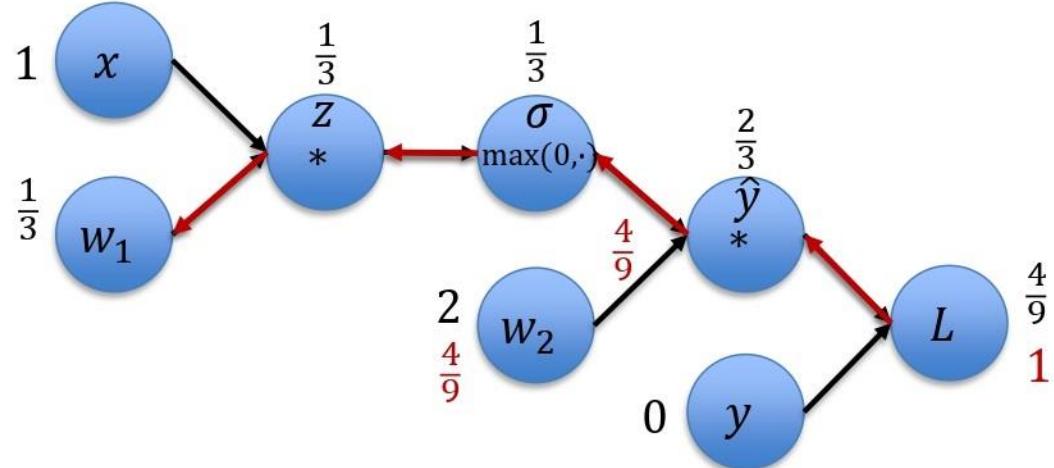
In our case  $n, d = 1$ :

$$L = (\hat{y} - y)^2 \Rightarrow \frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)$$

$$\hat{y} = w_2 \cdot \sigma \Rightarrow \frac{\partial \hat{y}}{\partial \sigma} = w_2$$

$$\sigma = \max(0, z) \Rightarrow \frac{\partial \sigma}{\partial z} = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{else} \end{cases}$$

$$z = x \cdot w_1 \Rightarrow \frac{\partial z}{\partial w_1} = x$$



Backpropagation

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

# Gradient Descent for Neural Networks

Initialize  $x = 1$ ,  $y = 0$ ,  
 $w_1 = \frac{1}{3}$ ,  $w_2 = 2$

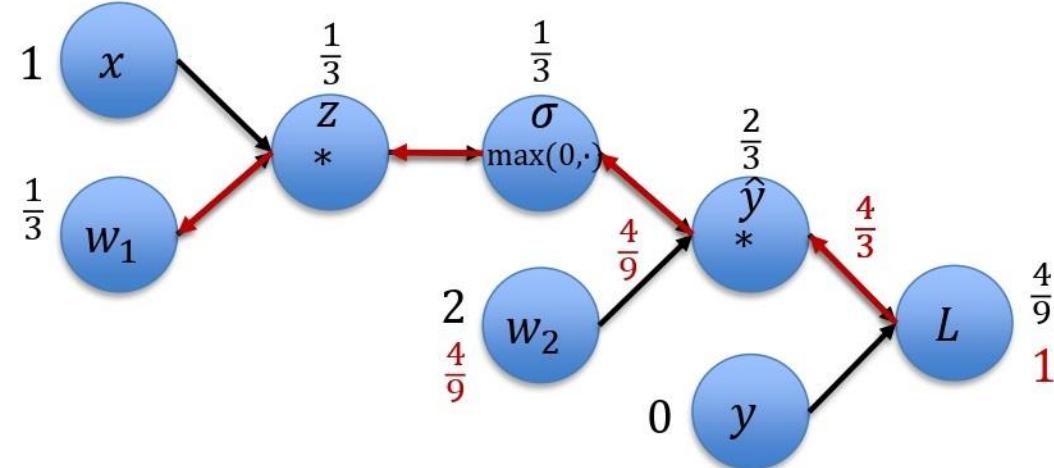
In our case  $n, d = 1$ :

$$L = (\hat{y} - y)^2 \Rightarrow \boxed{\frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)}$$

$$\hat{y} = w_2 \cdot \sigma \Rightarrow \frac{\partial \hat{y}}{\partial \sigma} = w_2$$

$$\sigma = \max(0, z) \Rightarrow \frac{\partial \sigma}{\partial z} = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{else} \end{cases}$$

$$z = x \cdot w_1 \Rightarrow \frac{\partial z}{\partial w_1} = x$$



Backpropagation

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

$$2 \cdot \frac{2}{3}$$

# Gradient Descent for Neural Networks

Initialize  $x = 1$ ,  $y = 0$ ,  
 $w_1 = \frac{1}{3}$ ,  $w_2 = 2$

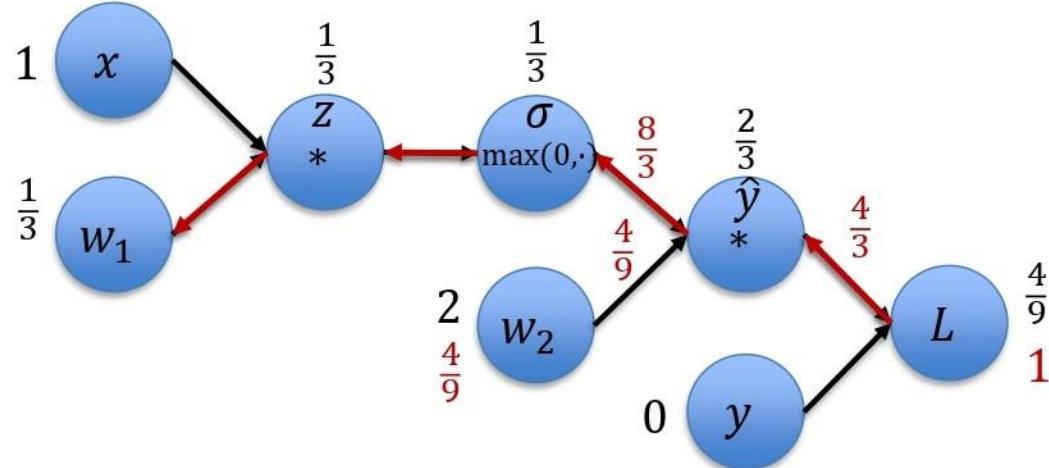
In our case  $n, d = 1$ :

$$L = (\hat{y} - y)^2 \Rightarrow \frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)$$

$$\hat{y} = w_2 \cdot \sigma \Rightarrow \boxed{\frac{\partial \hat{y}}{\partial \sigma} = w_2}$$

$$\sigma = \max(0, z) \Rightarrow \frac{\partial \sigma}{\partial z} = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{else} \end{cases}$$

$$z = x \cdot w_1 \Rightarrow \frac{\partial z}{\partial w_1} = x$$



Backpropagation

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

$$2 \cdot \frac{2}{3} \cdot 2$$

# Gradient Descent for Neural Networks

Initialize  $x = 1$ ,  $y = 0$ ,  
 $w_1 = \frac{1}{3}$ ,  $w_2 = 2$

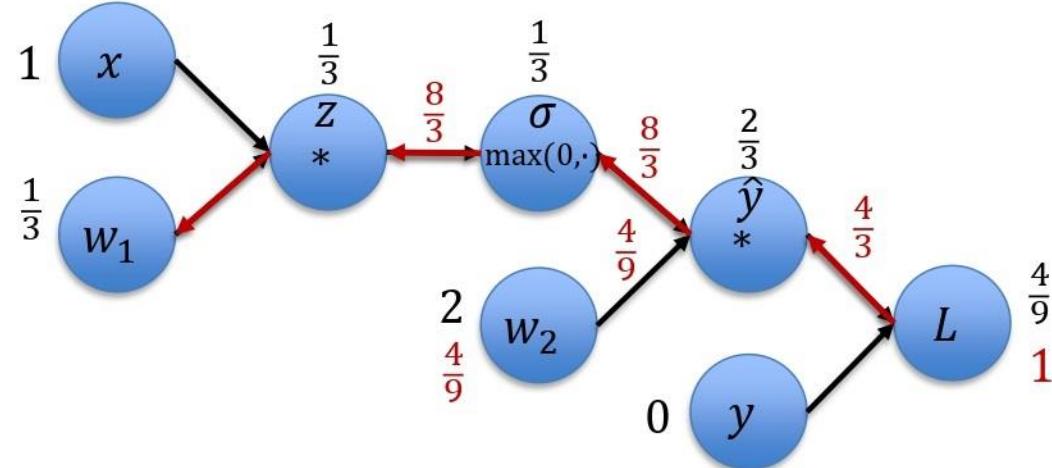
In our case  $n, d = 1$ :

$$L = (\hat{y} - y)^2 \Rightarrow \frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)$$

$$\hat{y} = w_2 \cdot \sigma \Rightarrow \frac{\partial \hat{y}}{\partial \sigma} = w_2$$

$$\sigma = \max(0, z) \Rightarrow \boxed{\frac{\partial \sigma}{\partial z} = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{else} \end{cases}}$$

$$z = x \cdot w_1 \Rightarrow \frac{\partial z}{\partial w_1} = x$$



Backpropagation

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

$$2 \cdot \frac{2}{3} \cdot 2 \cdot 1$$

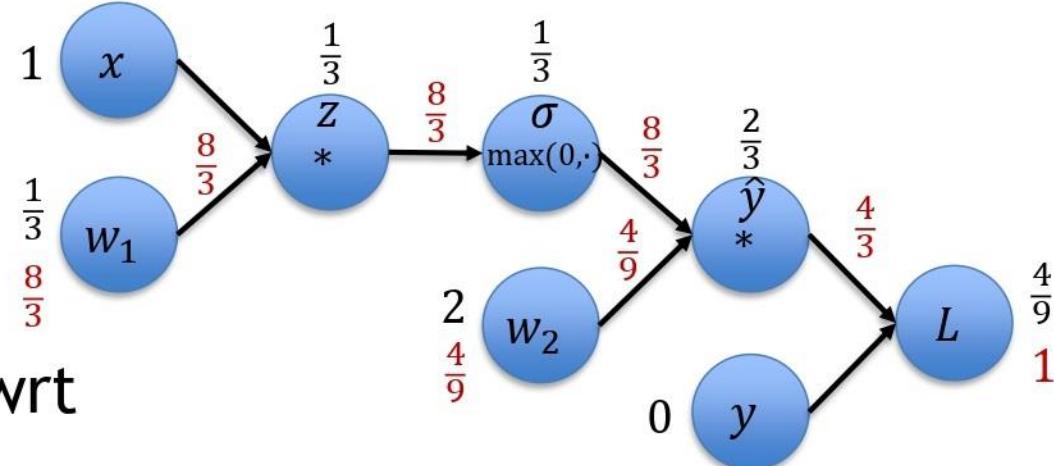
# Gradient Descent for Neural Networks

- Function we want to optimize:

$$f(x, \mathbf{w}) = \sum_{i=1}^n \|w_2 \max(0, w_1 x_i) - y_i\|_2^2$$

- Computed gradients wrt to weights  $w_1$  and  $w_2$
- Now: update the weights

$$\begin{aligned}\mathbf{w}' &= \mathbf{w} - \alpha \cdot \nabla_{\mathbf{w}} f = \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} - \alpha \cdot \begin{pmatrix} \nabla_{w_1} f \\ \nabla_{w_2} f \end{pmatrix} \\ &= \begin{pmatrix} \frac{1}{3} \\ 2 \end{pmatrix} - \alpha \cdot \begin{pmatrix} \frac{8}{3} \\ \frac{4}{9} \end{pmatrix}\end{aligned}$$



But: how to choose a good learning rate  $\alpha$  ?

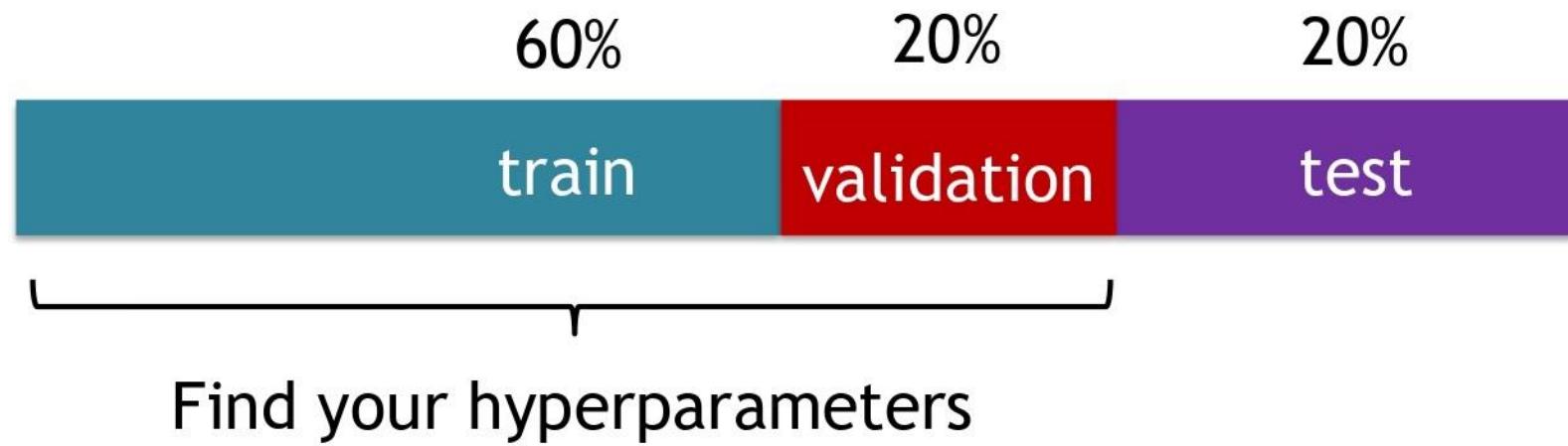
# Gradient Descent

- How to pick good learning rate?
- How to compute gradient for single training pair?
- How to compute gradient for large training set?
- How to speed things up? More to see in next lectures...

# Regularization

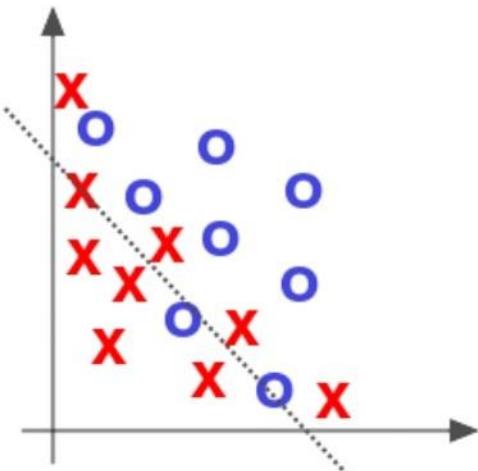
# Recap: Basic Recipe for ML

- Split your data

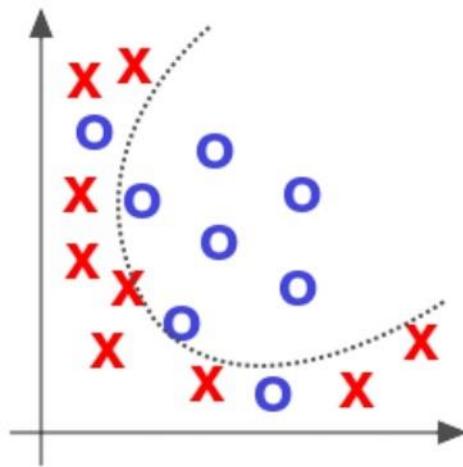


Other splits are also possible (e.g., 80%/10%/10%)

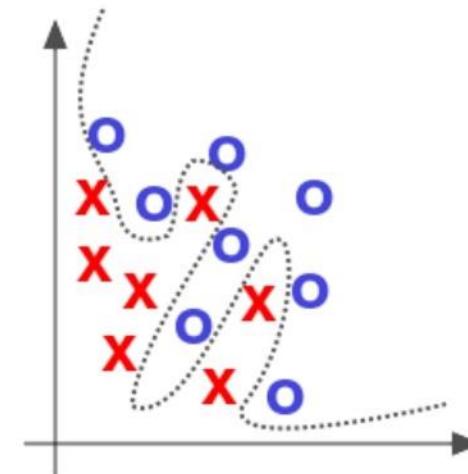
# Over- and Underfitting



**Underfitted**



**Appropriate**

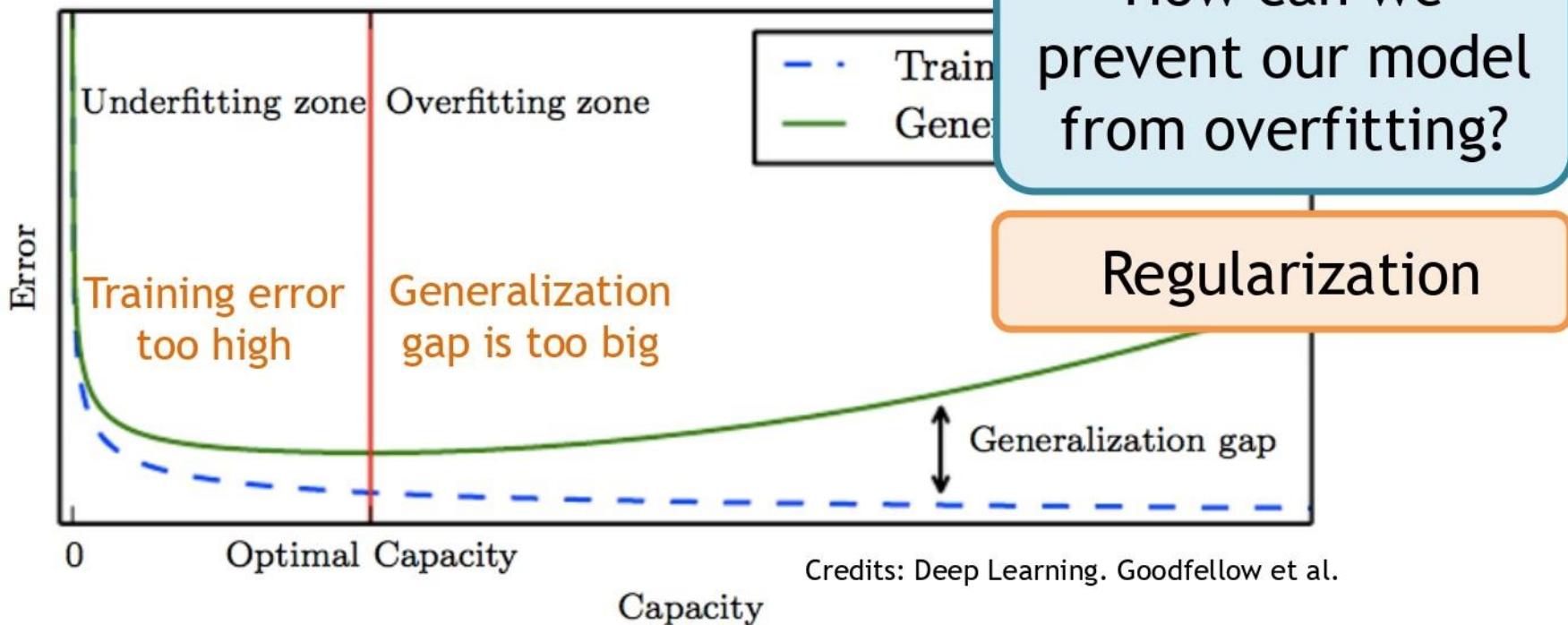


**Overfitted**

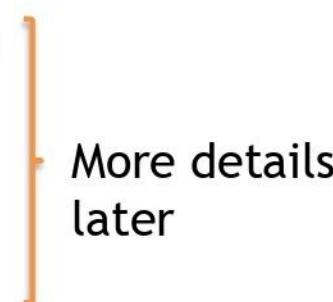
Source: Deep Learning by Adam Gibson, Josh Patterson, O'Reilly Media Inc., 2017

# Training a Neural Network

- Training/ Validation curve



# Regularization

- Loss function  $L(\mathbf{y}, \hat{\mathbf{y}}, \boldsymbol{\theta}) = \sum_{i=1}^n (\hat{y}_i - y_i)^2$
  - Regularization techniques
    - L2 regularization
    - L1 regularization
    - Max norm regularization
    - Dropout
    - Early stopping
    - ...
- 
- 

# Regularization: Example

- Input: 3 features  $x = [1, 2, 1]$
- Two linear classifiers that give the same result:
  - $\theta_1 = [0, 0.75, 0]$   Ignores 2 features
  - $\theta_2 = [0.25, 0.5, 0.25]$   Takes information from all features

# Regularization: Example

- Loss  $L(\mathbf{y}, \hat{\mathbf{y}}, \boldsymbol{\theta}) = \sum_{i=1}^n (x_i \theta_{ji} - y_i)^2 + \lambda R(\boldsymbol{\theta})$
- L2 regularization  $R(\boldsymbol{\theta}) = \sum_{i=1}^n \theta_i^2$   
 $\theta_1 \rightarrow 0 + 0.75^2 + 0 = 0.5625$   
 $\theta_2 \rightarrow 0.25^2 + 0.5^2 + 0.25^2 = 0.375$  Minimization

$$x = [1, 2, 1], \theta_1 = [0, 0.75, 0], \theta_2 = [0.25, 0.5, 0.25]$$

# Regularization: Example

- Loss  $L(\mathbf{y}, \hat{\mathbf{y}}, \boldsymbol{\theta}) = \sum_{i=1}^n (x_i \theta_{ji} - y_i)^2 + \lambda R(\boldsymbol{\theta})$

- L1 regularization  $R(\boldsymbol{\theta}) = \sum_{i=1}^n |\theta_i|$

$$\theta_1 \longrightarrow 0 + 0.75 + 0 = 0.75$$

$$\theta_2 \longrightarrow 0.25 + 0.5 + 0.25 = 1$$

Minimization

$$x = [1, 2, 1], \theta_1 = [0, 0.75, 0], \theta_2 = [0.25, 0.5, 0.25]$$

# Regularization: Example

- Input: 3 features  $x = [1, 2, 1]$
- Two linear classifiers that give the same result:

$\theta_1 = [0, 0.75, 0]$   Ignores 2 features

$\theta_2 = [0.25, 0.5, 0.25]$   Takes information  
from all features

# Regularization: Example

- Input: 3 features  $x = [1, 2, 1]$
- Two linear classifiers that give the same result:

$$\theta_1 = [0, 0.75, 0]$$

→ L1 regularization  
enforces **sparsity**

$$\theta_2 = [0.25, 0.5, 0.25]$$

→ Takes information  
from all features

# Regularization: Example

- Input: 3 features  $x = [1, 2, 1]$
- Two linear classifiers that give the same result:

$$\theta_1 = [0, 0.75, 0]$$



L1 regularization  
enforces **sparsity**

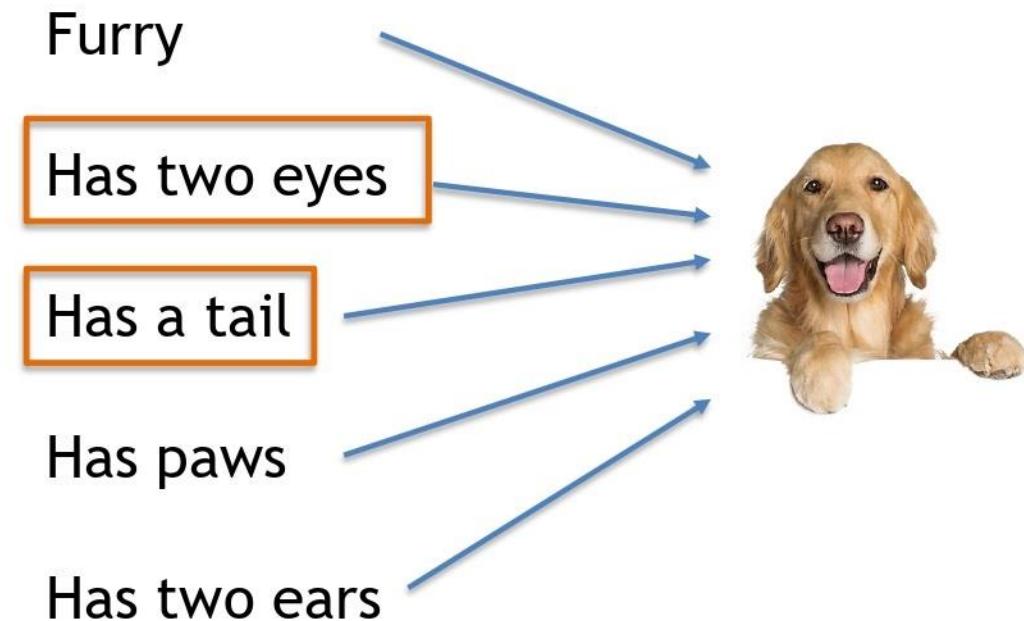
$$\theta_2 = [0.25, 0.5, 0.25]$$



L2 regularization  
enforces that the weights  
have **similar values**

# Regularization: Effect

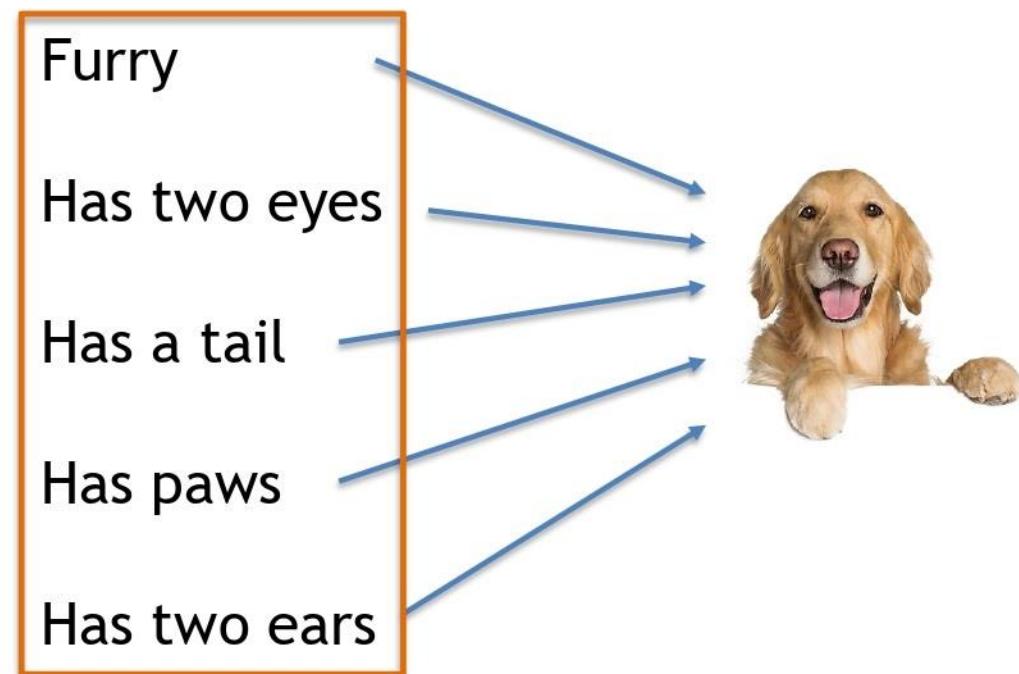
- Dog classifier takes different inputs



**L1 regularization**  
will focus all the  
attention to a  
few key features

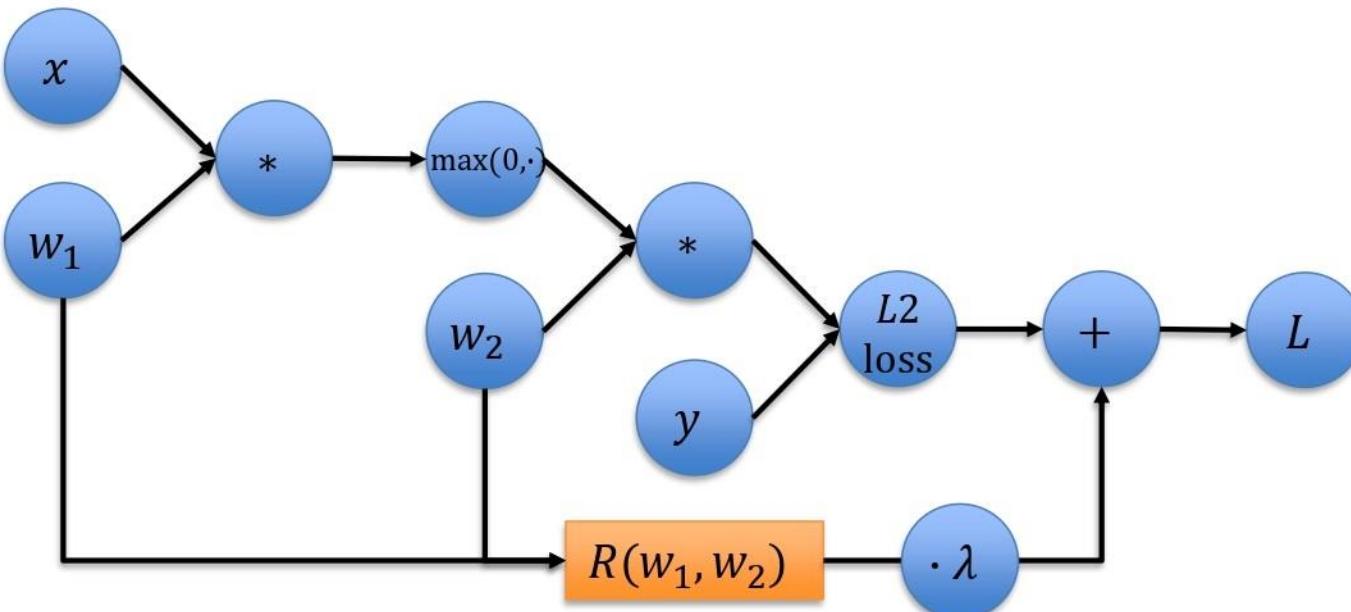
# Regularization: Effect

- Dog classifier takes different inputs



**L2 regularization**  
will take all  
information into  
account to make  
decisions

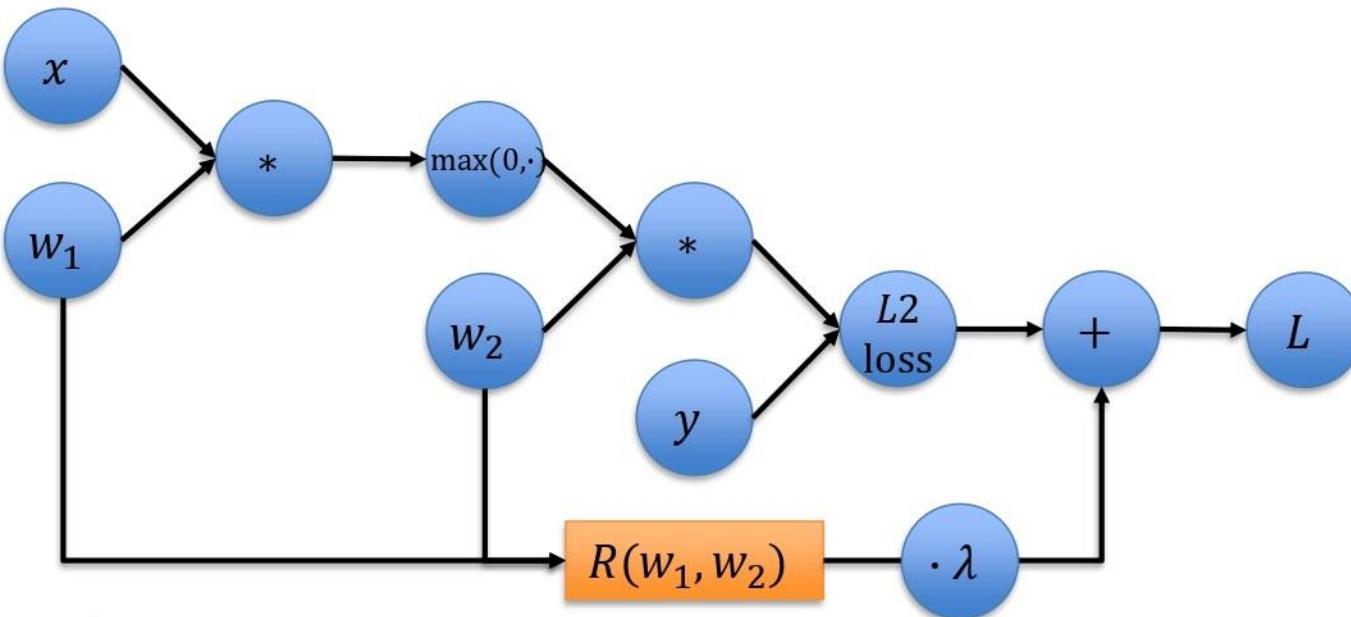
# Regularization for Neural Networks



Combining nodes:  
Network output + L2-loss +  
regularization

$$\sum_{i=1}^n \|w_2 \max(0, w_1 x_i) - y_i\|_2^2 + \lambda R(w_1, w_2)$$

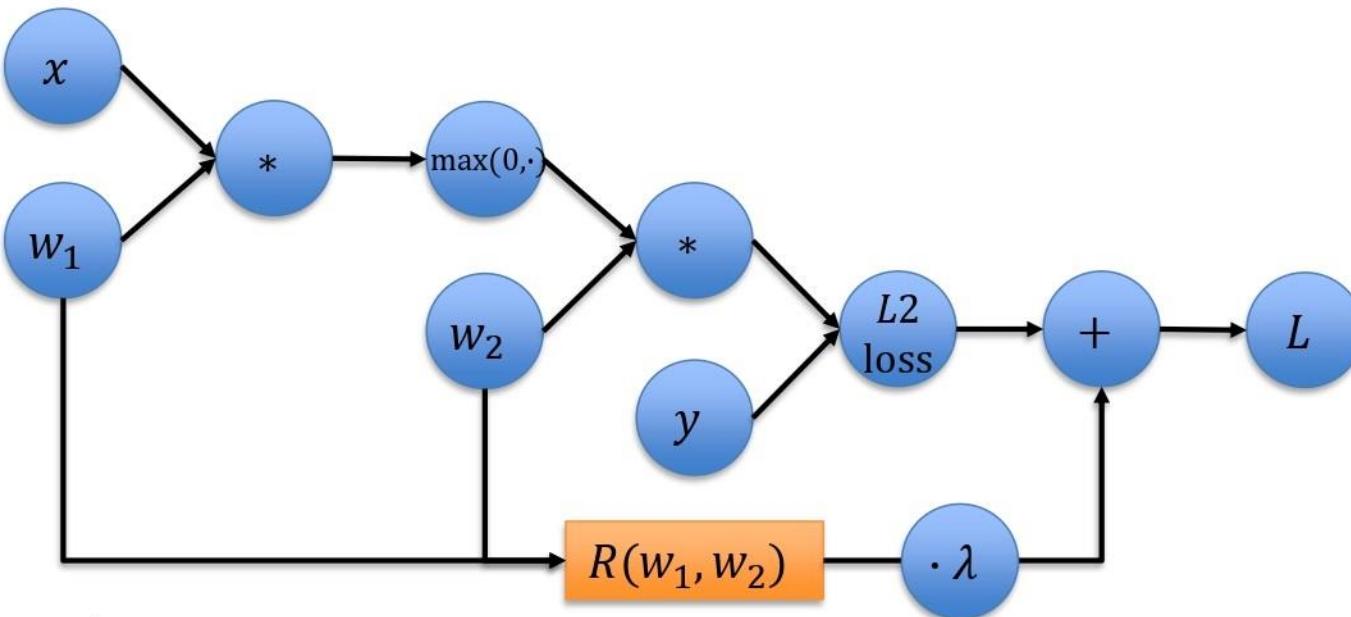
# Regularization for Neural Networks



Combining nodes:  
Network output + L2-loss +  
regularization

$$\sum_{i=1}^n \|w_2 \max(0, w_1 x_i) - y_i\|_2^2 + \lambda \left\| \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} \right\|_2^2$$

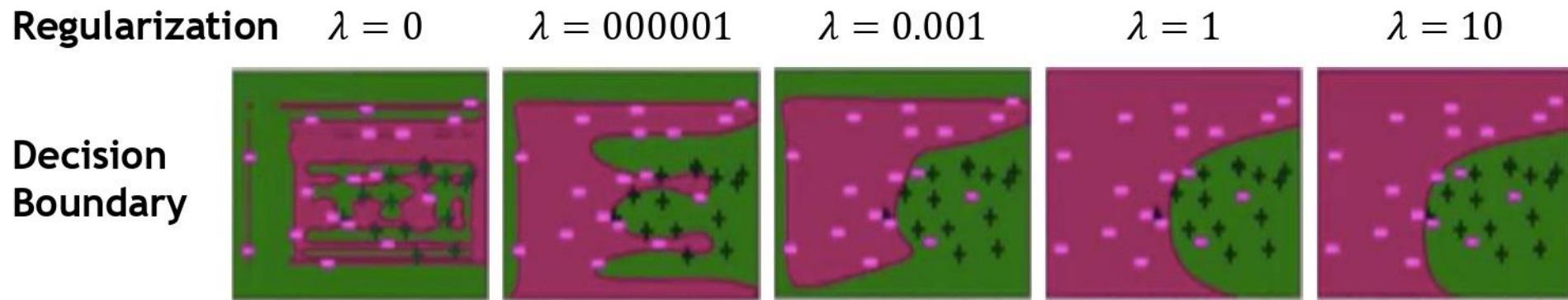
# Regularization for Neural Networks



Combining nodes:  
Network output + L2-loss +  
regularization

$$\sum_{i=1}^n \|w_2 \max(0, w_1 x_i) - y_i\|_2^2 + \lambda(w_1^2 + w_2^2)$$

# Regularization



Credits: University of Washington

What is the goal of regularization?

What happens to the training error?

# Regularization

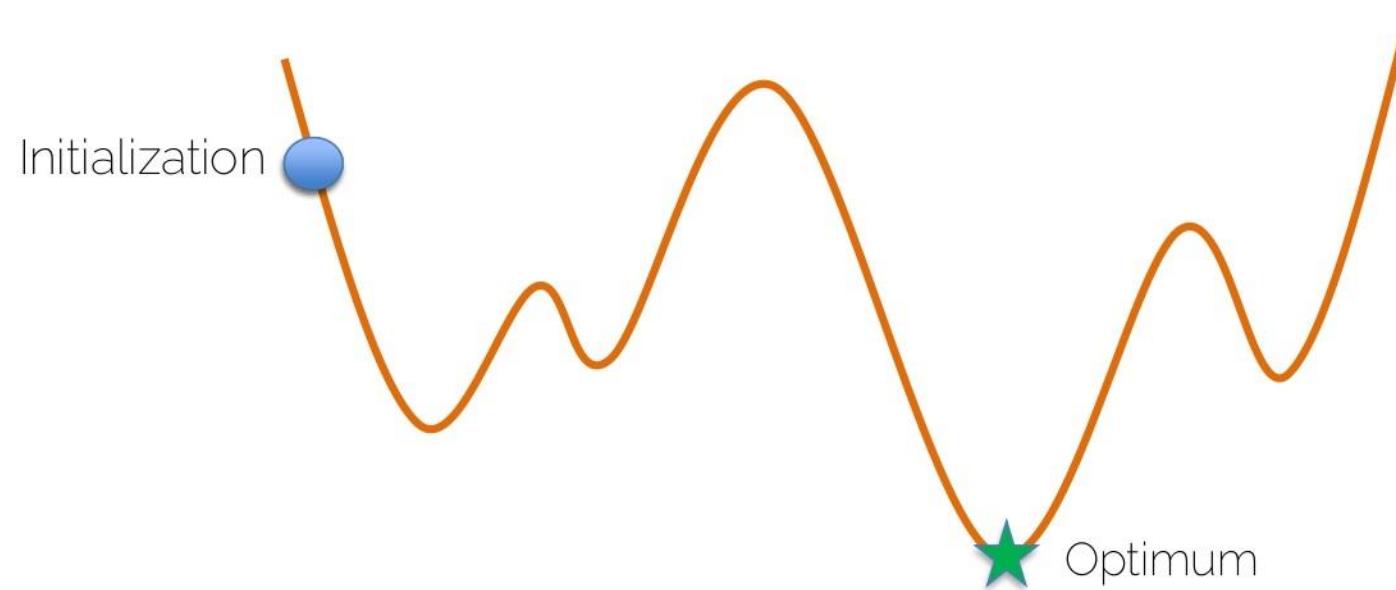
- Any strategy that aims to

Lower  
validation  
error

Increasing  
training error

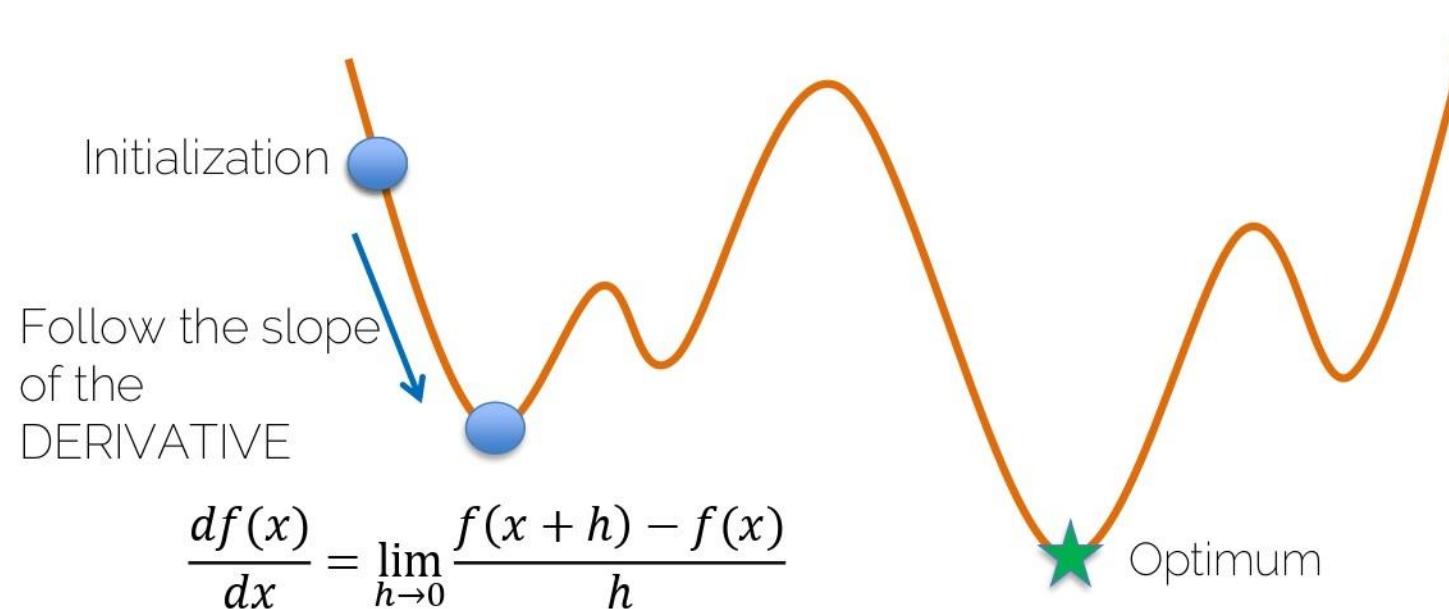
# Gradient Descent

$$x^* = \arg \min f(x)$$



# Gradient Descent

$$x^* = \arg \min f(x)$$



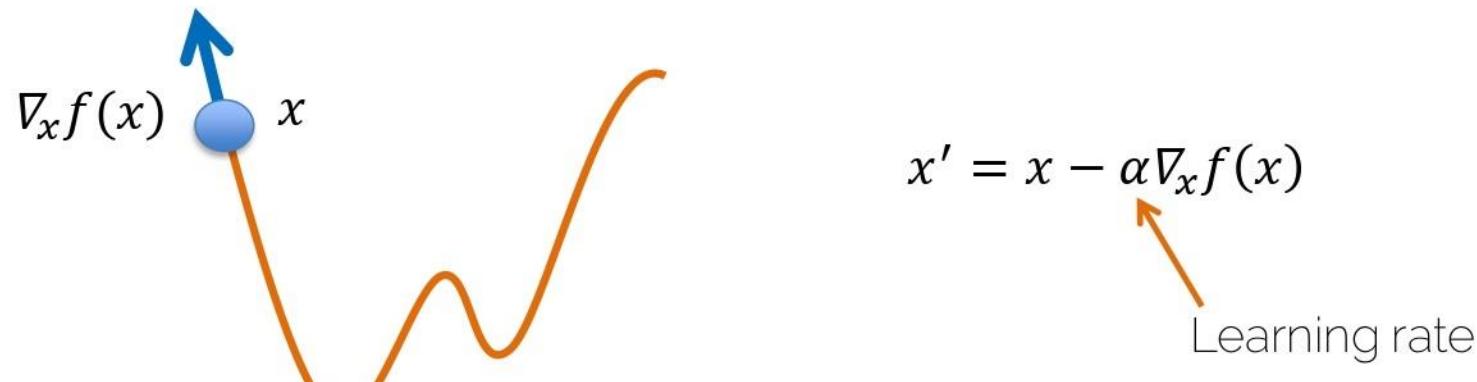
# Gradient Descent

- From derivative to gradient

$$\frac{df(x)}{dx} \longrightarrow \nabla_x f(x)$$

Direction of greatest increase of the function

- Gradient steps in direction of negative gradient



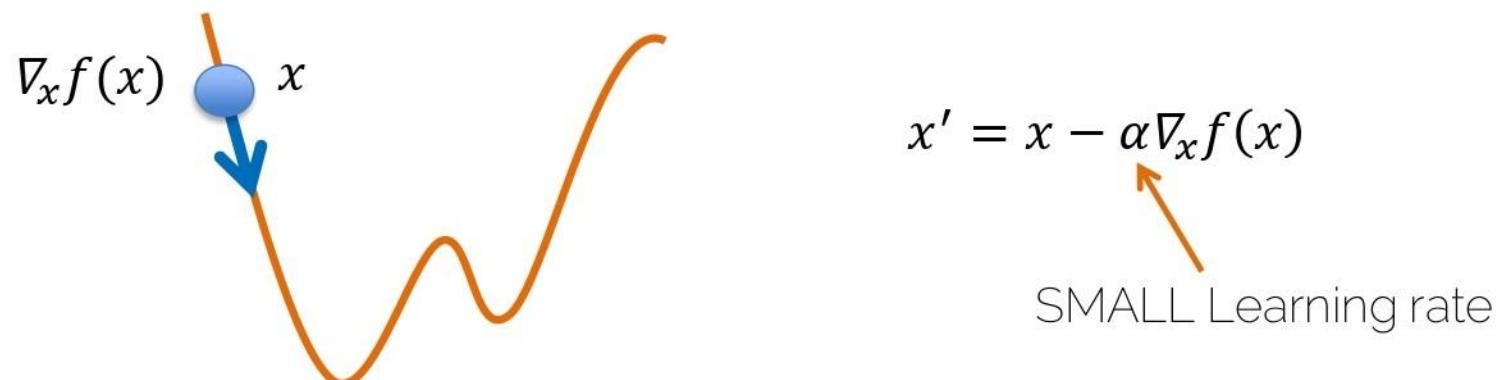
# Gradient Descent

- From derivative to gradient

$$\frac{df(x)}{dx} \longrightarrow \nabla_x f(x)$$

Direction of greatest increase of the function

- Gradient steps in direction of negative gradient



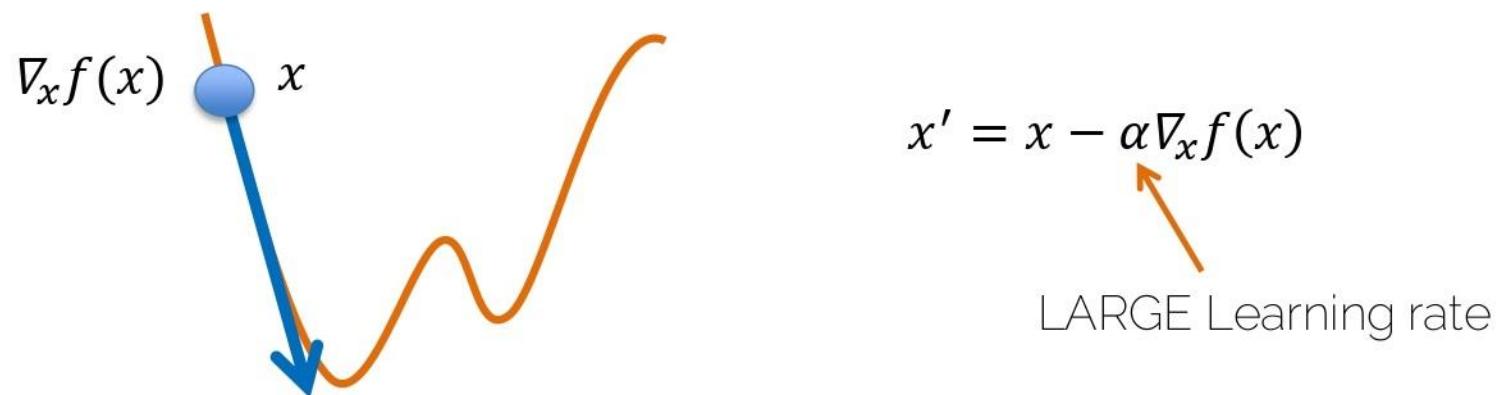
# Gradient Descent

- From derivative to gradient

$$\frac{df(x)}{dx} \longrightarrow \nabla_x f(x)$$

Direction of greatest increase of the function

- Gradient steps in direction of negative gradient



# Gradient Descent

$$\mathbf{x}^* = \arg \min f(\mathbf{x})$$

Initialization

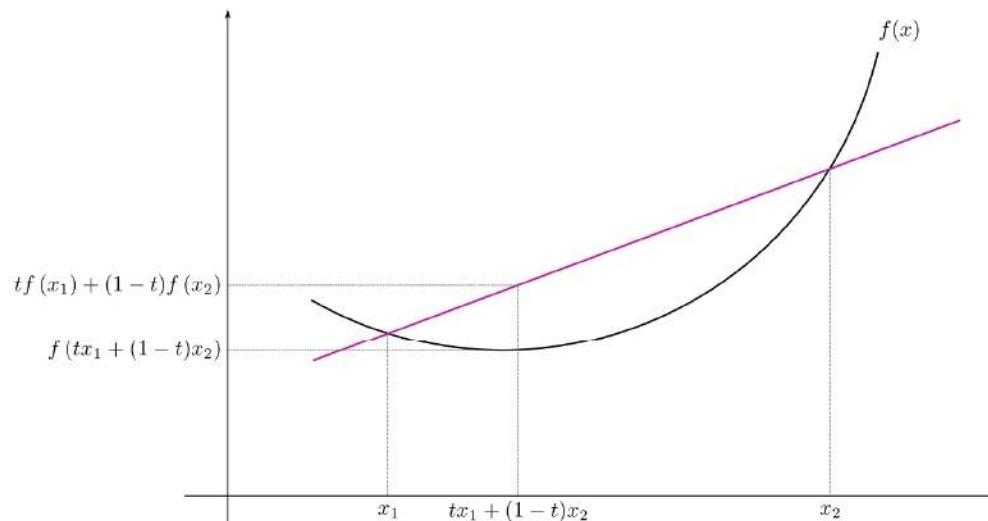
What is the gradient when we reach this point?

Optimum

Not guaranteed to reach the global optimum

# Convergence of Gradient Descent

- Convex function: all local minima are global minima

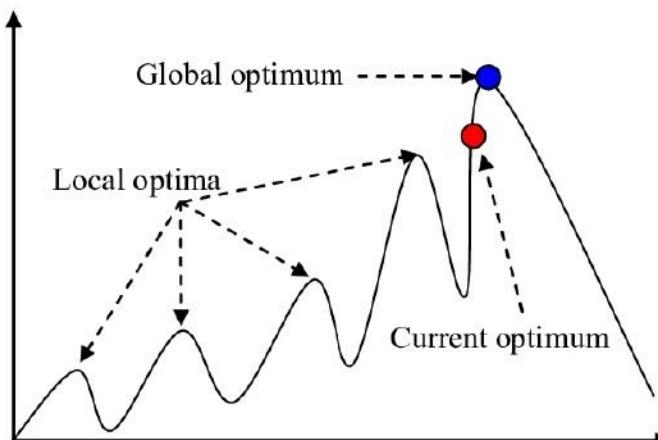


Source: [https://en.wikipedia.org/wiki/Convex\\_function#/media/File:ConvexFunction.svg](https://en.wikipedia.org/wiki/Convex_function#/media/File:ConvexFunction.svg)

If line/plane segment between any two points lies above or on the graph

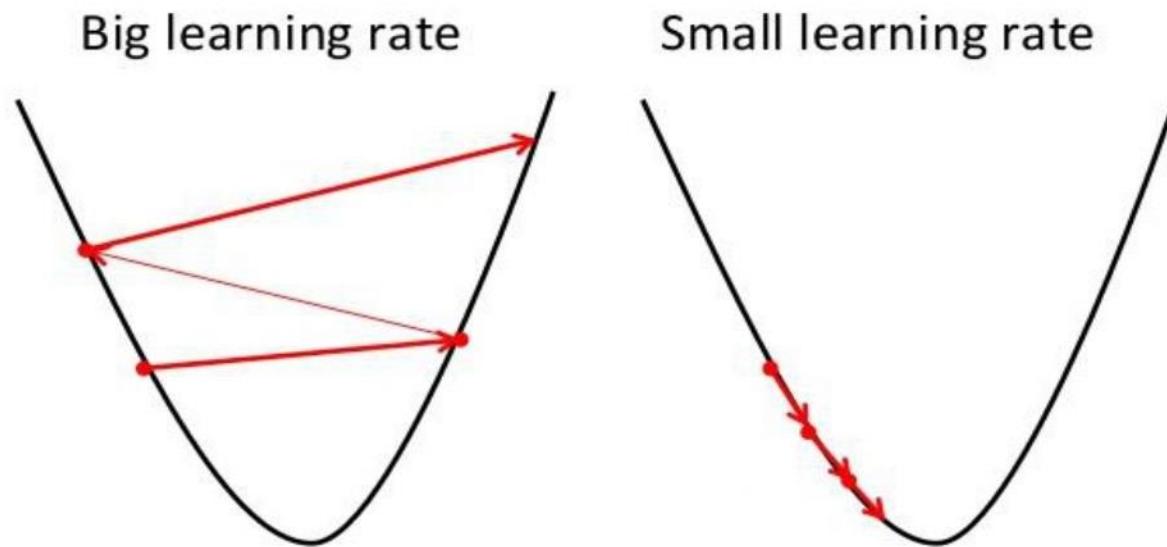
# Convergence of Gradient Descent

- Neural networks are non-convex
  - many (different) local minima
  - no (practical) way to say which is globally optimal



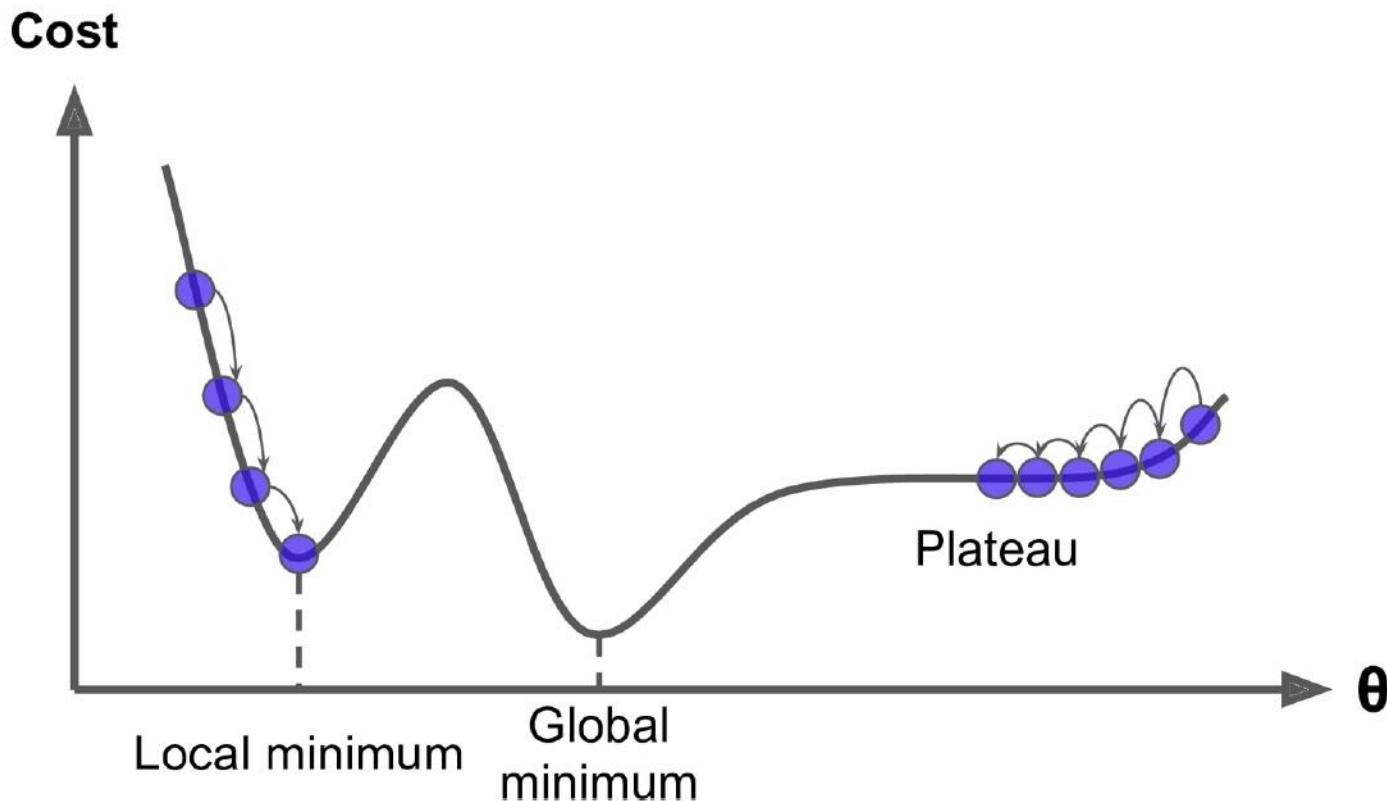
Source: Li, Qi. (2006). Challenging Registration of Geologic Image Data

# Convergence of Gradient Descent



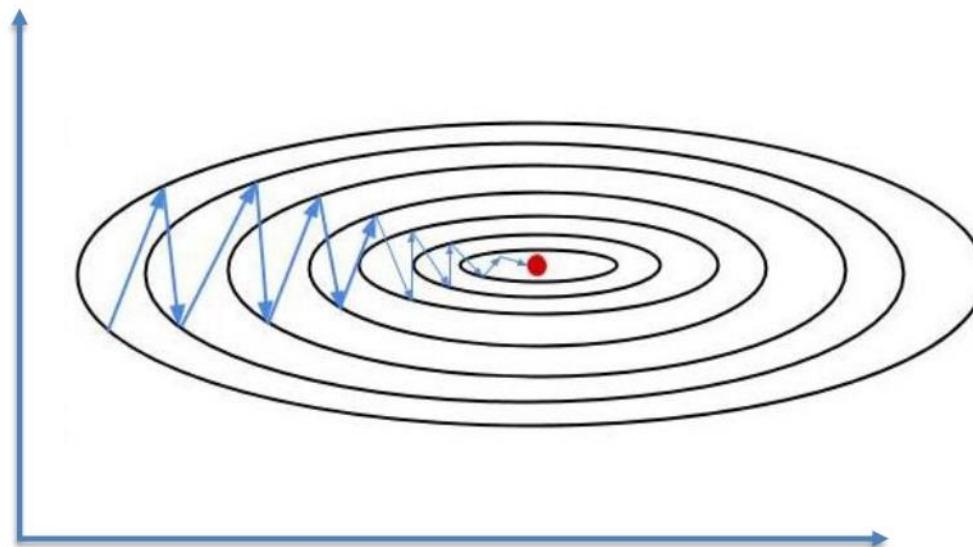
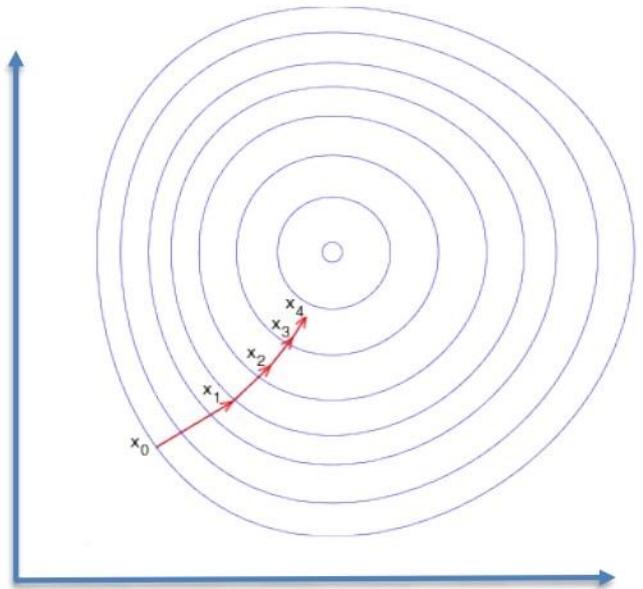
Source: <https://builtin.com/data-science/gradient-descent>

# Convergence of Gradient Descent



Source: A. Geron

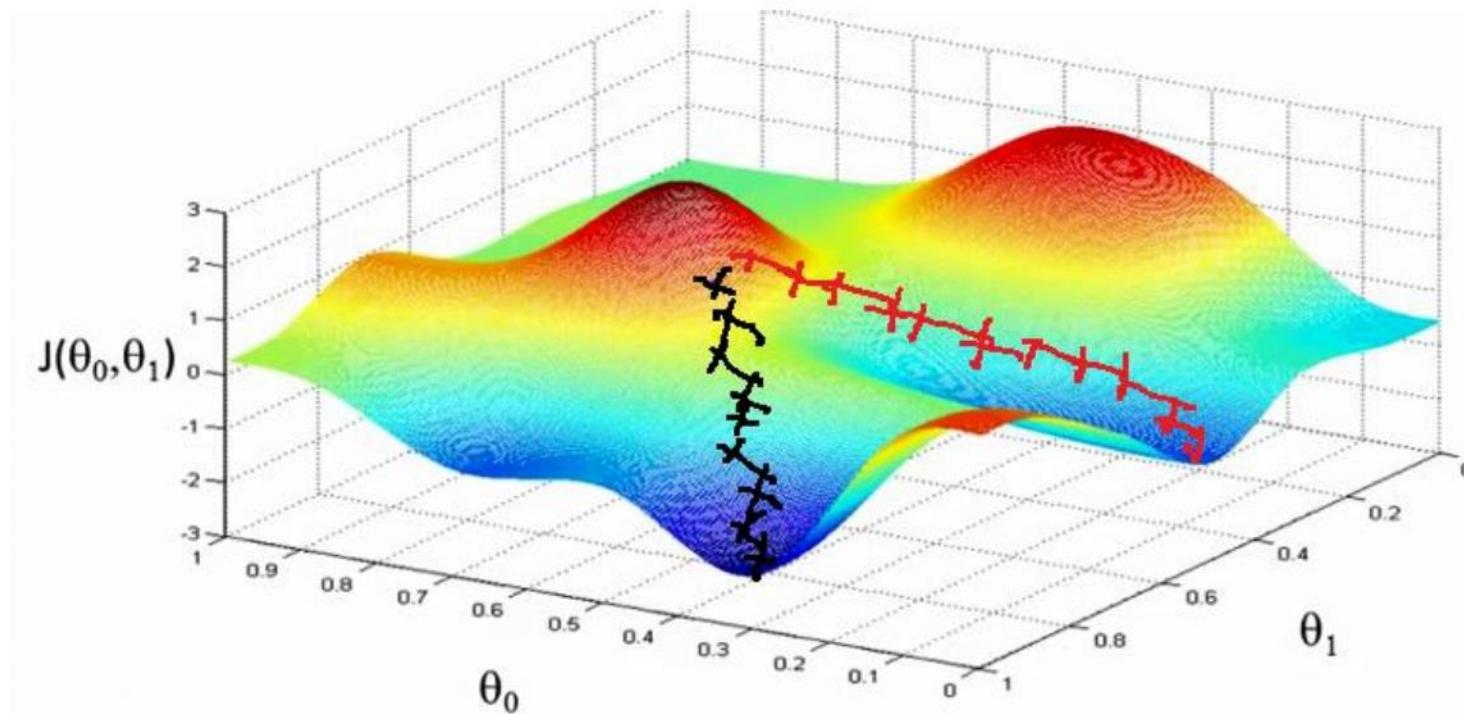
# Gradient Descent: Multiple Dimensions



Source: [builtin.com/data-science/gradient-descent](https://builtin.com/data-science/gradient-descent)

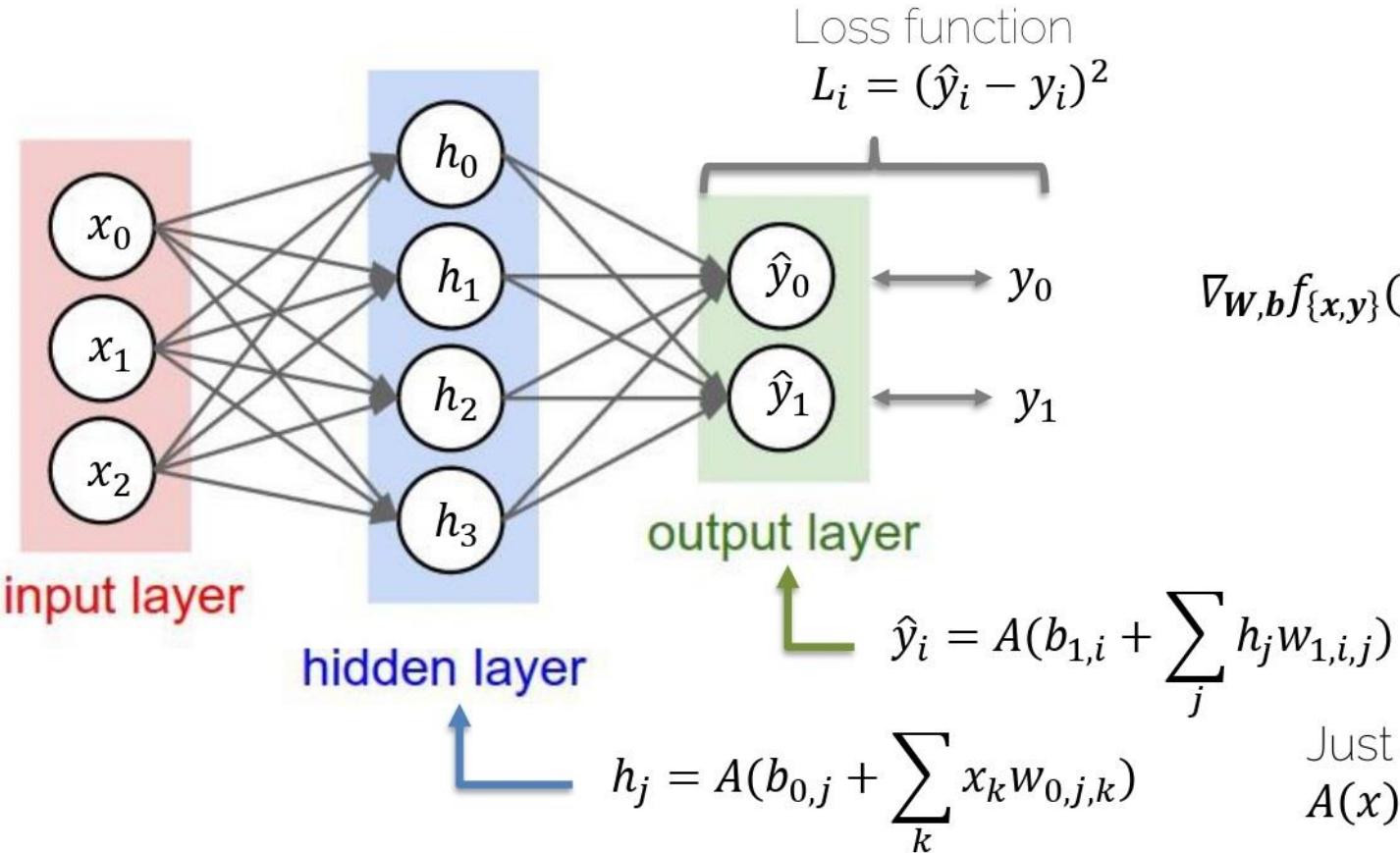
Various ways to visualize...

# Gradient Descent: Multiple Dimensions



Source: <http://blog.datumbox.com/wp-content/uploads/2013/10/gradient-descent.png>

# Gradient Descent for Neural Networks



$$\nabla_{W,b} f_{\{x,y\}}(W) = \begin{bmatrix} \frac{\partial f}{\partial w_{0,0,0}} \\ \dots \\ \frac{\partial f}{\partial w_{l,m,n}} \\ \dots \\ \frac{\partial f}{\partial b_{l,m}} \end{bmatrix}$$

Just simple:  
 $A(x) = \max(0, x)$

# Gradient Descent: Single Training Sample

- Given a loss function  $L$  and a single training sample  $\{\mathbf{x}_i, \mathbf{y}_i\}$
- Find best model parameters  $\boldsymbol{\theta} = \{\mathbf{W}, \mathbf{b}\}$
- Cost  $L_i(\boldsymbol{\theta}, \mathbf{x}_i, \mathbf{y}_i)$ 
  - $\boldsymbol{\theta} = \arg \min L_i(\mathbf{x}_i, \mathbf{y}_i)$
- Gradient Descent:
  - Initialize  $\boldsymbol{\theta}^1$  with 'random' values (more to that later)
  - $\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k - \alpha \nabla_{\boldsymbol{\theta}} L_i(\boldsymbol{\theta}^k, \mathbf{x}_i, \mathbf{y}_i)$
  - Iterate until convergence:  $|\boldsymbol{\theta}^{k+1} - \boldsymbol{\theta}^k| < \epsilon$

# Gradient Descent: Single Training Sample

- $\theta^{k+1} = \theta^k - \alpha \nabla_{\theta} L_i(\theta^k, x_i, y_i)$ 
  - Weights, biases after update step
  - Weights, biases at step k (current model)
- $\nabla_{\theta} L_i(\theta^k, x_i, y_i)$  computed via backpropagation
- Typically:  $\dim(\nabla_{\theta} L_i(\theta^k, x_i, y_i)) = \dim(\theta) \gg 1 \text{ million}$

## Gradient Descent: Multiple Training Samples

- Given a loss function  $L$  and multiple ( $n$ ) training samples  $\{\mathbf{x}_i, \mathbf{y}_i\}$
- Find best model parameters  $\boldsymbol{\theta} = \{\mathbf{W}, \mathbf{b}\}$
- Cost  $L = \frac{1}{n} \sum_{i=1}^n L_i(\boldsymbol{\theta}, \mathbf{x}_i, \mathbf{y}_i)$ 
  - $\boldsymbol{\theta} = \arg \min L$

# Gradient Descent: Multiple Training Samples

- Update step for multiple samples

$$\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k - \alpha \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^k, \mathbf{x}_{\{1..n\}}, \mathbf{y}_{\{1..n\}})$$

- Gradient is average / sum over residuals

$$\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^k, \mathbf{x}_{\{1..n\}}, \mathbf{y}_{\{1..n\}}) = \frac{1}{n} \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} L_i(\boldsymbol{\theta}^k, \mathbf{x}_i, \mathbf{y}_i)$$

Reminder: this comes from backprop.

- Often people are lazy and just write:  $\nabla L = \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} L_i$ 
  - omitting  $\frac{1}{n}$  is not 'wrong', it just means rescaling the learning rate

# Gradient Descent on Train Set

- Given large train set with  $n$  training samples  $\{\mathbf{x}_i, \mathbf{y}_i\}$ 
    - Let's say 1 million labeled images
    - Let's say our network has 500k parameters
  - Gradient has 500k dimensions
  - $n = 1 \text{ million}$
- Extremely expensive to compute

# Stochastic Gradient Descent (SGD)

- If we have  $n$  training samples we need to compute the gradient for all of them which is  $O(n)$
- If we consider the problem as empirical risk minimization, we can express the total loss over the training data as the expectation of all the samples

$$\frac{1}{n} \left( \sum_{i=1}^n L_i(\boldsymbol{\theta}, \mathbf{x}_i, \mathbf{y}_i) \right) = \mathbb{E}_{i \sim [1, \dots, n]} [L_i(\boldsymbol{\theta}, \mathbf{x}_i, \mathbf{y}_i)]$$

# Stochastic Gradient Descent (SGD)

- The expectation can be approximated with a small subset of the data

$$\mathbb{E}_{i \sim [1, \dots, n]}[L_i(\boldsymbol{\theta}, \mathbf{x}_i, \mathbf{y}_i)] \approx \frac{1}{|S|} \sum_{j \in S} (L_j(\boldsymbol{\theta}, \mathbf{x}_j, \mathbf{y}_j)) \text{ with } S \subseteq \{1, \dots, n\}$$

Minibatch  
choose subset of trainset  $m \ll n$

$$B_i = \{\{\mathbf{x}_1, \mathbf{y}_1\}, \{\mathbf{x}_2, \mathbf{y}_2\}, \dots, \{\mathbf{x}_m, \mathbf{y}_m\}\}$$
$$\{B_1, B_2, \dots, B_{n/m}\}$$

# Stochastic Gradient Descent (SGD)

- Minibatch size is hyperparameter
  - Typically power of 2 → 8, 16, 32, 64, 128...
  - Smaller batch size means greater variance in the gradients  
→ noisy updates
  - Mostly limited by GPU memory (in backward pass)
  - E.g.,
    - Train set has  $\mathbf{n} = 2^{20}$  (about 1 million) images
    - With batch size  $\mathbf{m} = 64$ :  $B_1 \dots n/m = B_1 \dots 16,384$  minibatches  
(Epoch = complete pass through training set)

# Stochastic Gradient Descent (SGD)

$$\theta^{k+1} = \theta^k - \alpha \nabla_{\theta} L(\theta^k, x_{\{1..m\}}, y_{\{1..m\}})$$

$\nabla_{\theta} L = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} L_i$

$m$  training samples in the current minibatch

Gradient for the  $k$ -th minibatch

$k$  now refers to  $k$ -th iteration

Note the terminology: iteration vs epoch

# Convergence of SGD

Suppose we want to minimize the function  $F(\theta)$  with the stochastic approximation

$$\theta^{k+1} = \theta^k - \alpha_k H(\theta^k, X)$$

where  $\alpha_1, \alpha_2 \dots \alpha_n$  is a sequence of positive step-sizes and  $H(\theta^k, X)$  is the unbiased estimate of  $\nabla F(\theta^k)$ , i.e.

$$\mathbb{E}[H(\theta^k, X)] = \nabla F(\theta^k)$$

Robbins, H. and Monro, S. "A Stochastic Approximation Method" 1951.

# Convergence of SGD

$$\theta^{k+1} = \theta^k - \alpha_k H(\theta^k, X)$$

converges to a local (**global**) minimum if the following conditions are met:

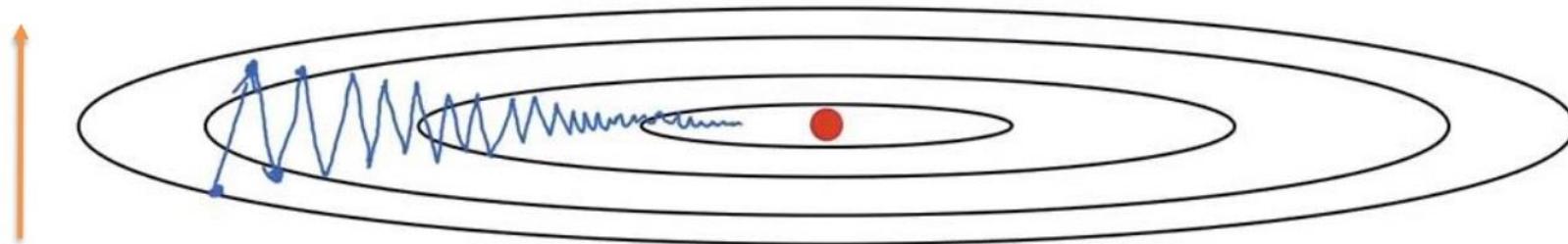
- 1)  $\alpha_n \geq 0, \forall n \geq 0$
- 2)  $\sum_{n=1}^{\infty} \alpha_n = \infty$
- 3)  $\sum_{n=1}^{\infty} \alpha_n^2 < \infty$
- 4)  **$F(\theta)$  is strictly convex**

The proposed sequence by Robbins and Monro is  $\alpha_n \propto \frac{\alpha}{n}$ , for  $n > 0$

# Problems of SGD

- Gradient is scaled equally across all dimensions
  - i.e., cannot independently scale directions
  - need to have conservative min learning rate to avoid divergence
  - Slower than 'necessary'
- Finding good learning rate is an art by itself
  - More next lecture

# Gradient Descent with Momentum



Source: A. Ng

We're making many steps back and forth along this dimension. Would love to track that this is averaging out over time.

Would love to go faster here...  
i.e., accumulated gradients over time

# Gradient Descent with Momentum

$$\boldsymbol{v}^{k+1} = \beta \cdot \boldsymbol{v}^k - \alpha \cdot \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^k)$$

accumulation rate ('friction', momentum)      velocity      learning rate      Gradient of current minibatch

$$\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k + \boldsymbol{v}^{k+1}$$

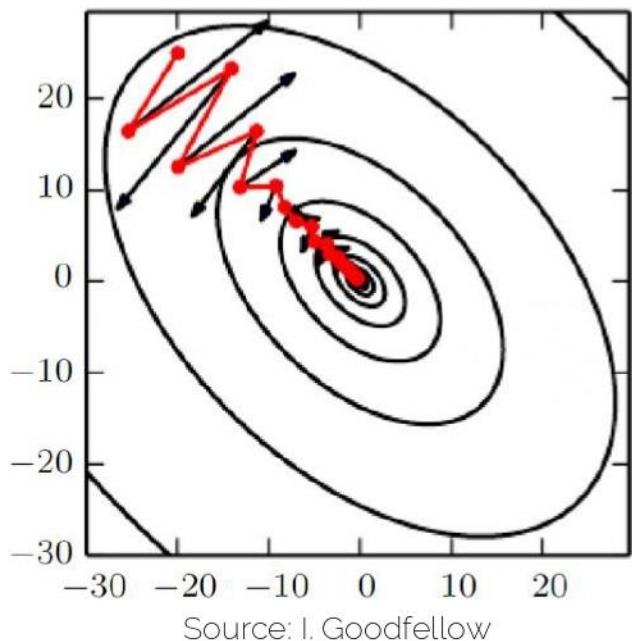
weights of model      velocity

Exponentially-weighted average of gradient

Important: velocity  $\boldsymbol{v}^k$  is vector-valued!

[Sutskever et al., ICML'13] On the importance of initialization and momentum in deep learning

# Gradient Descent with Momentum



Step will be largest when a sequence of gradients all point to the same direction

Hyperparameters are  $\alpha, \beta$   
 $\beta$  is often set to 0.9

$$\theta^{k+1} = \theta^k + v^{k+1}$$

# Gradient Descent with Momentum

- Can it overcome local minima?



$$\theta^{k+1} = \theta^k + v^{k+1}$$

# Nesterov Momentum

- Look-ahead momentum

$$\tilde{\theta}^{k+1} = \theta^k + \beta \cdot v^k$$

$$v^{k+1} = \beta \cdot v^k - \alpha \cdot \nabla_{\theta} L(\tilde{\theta}^{k+1})$$

$$\theta^{k+1} = \theta^k + v^{k+1}$$

Nesterov, Yurii E. "A method for solving the convex programming problem with convergence rate  $O(1/k^2)$ ." *Dokl. akad. nauk Sssr.* Vol. 269. 1983.

# Nesterov Momentum

- First make a big jump in the direction of the previous accumulated gradient.
- Then measure the gradient where you end up and make a correction.



brown vector = jump,      red vector = correction,      green vector = accumulated gradient

blue vectors = standard momentum

Source: G. Hinton

$$\begin{aligned}\tilde{\theta}^{k+1} &= \theta^k + \beta \cdot v^k \\ v^{k+1} &= \beta \cdot v^k - \alpha \cdot \nabla_{\theta} L(\tilde{\theta}^{k+1}) \\ \theta^{k+1} &= \theta^k + v^{k+1}\end{aligned}$$