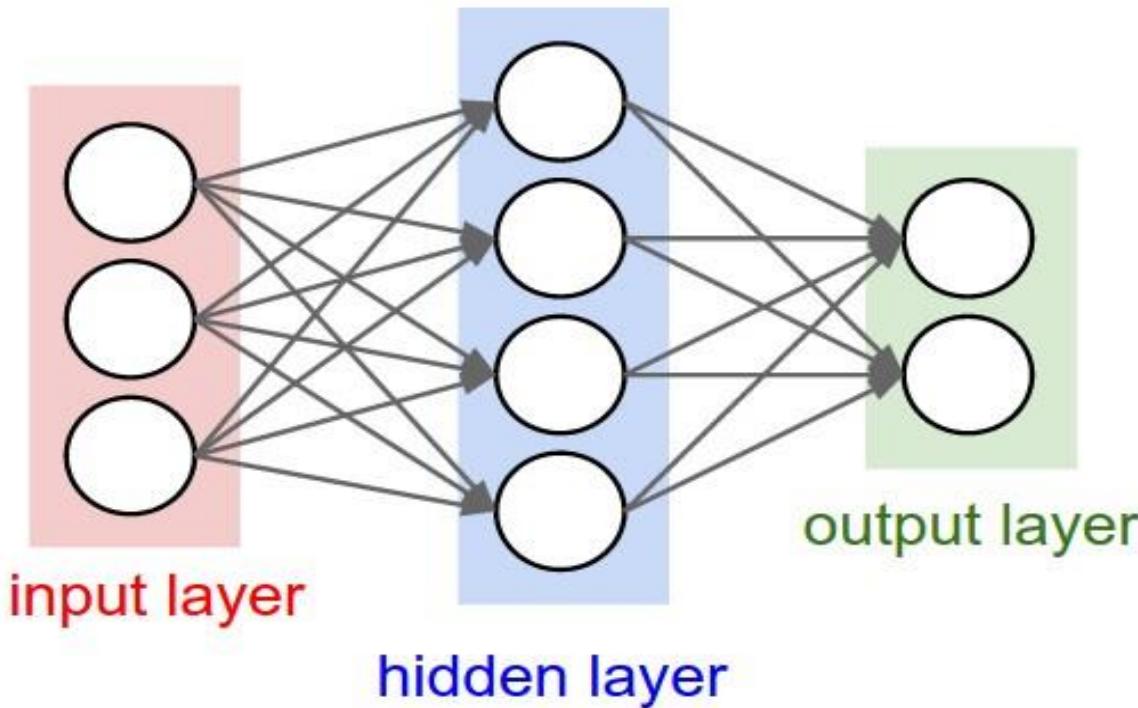


Scaling Optimization

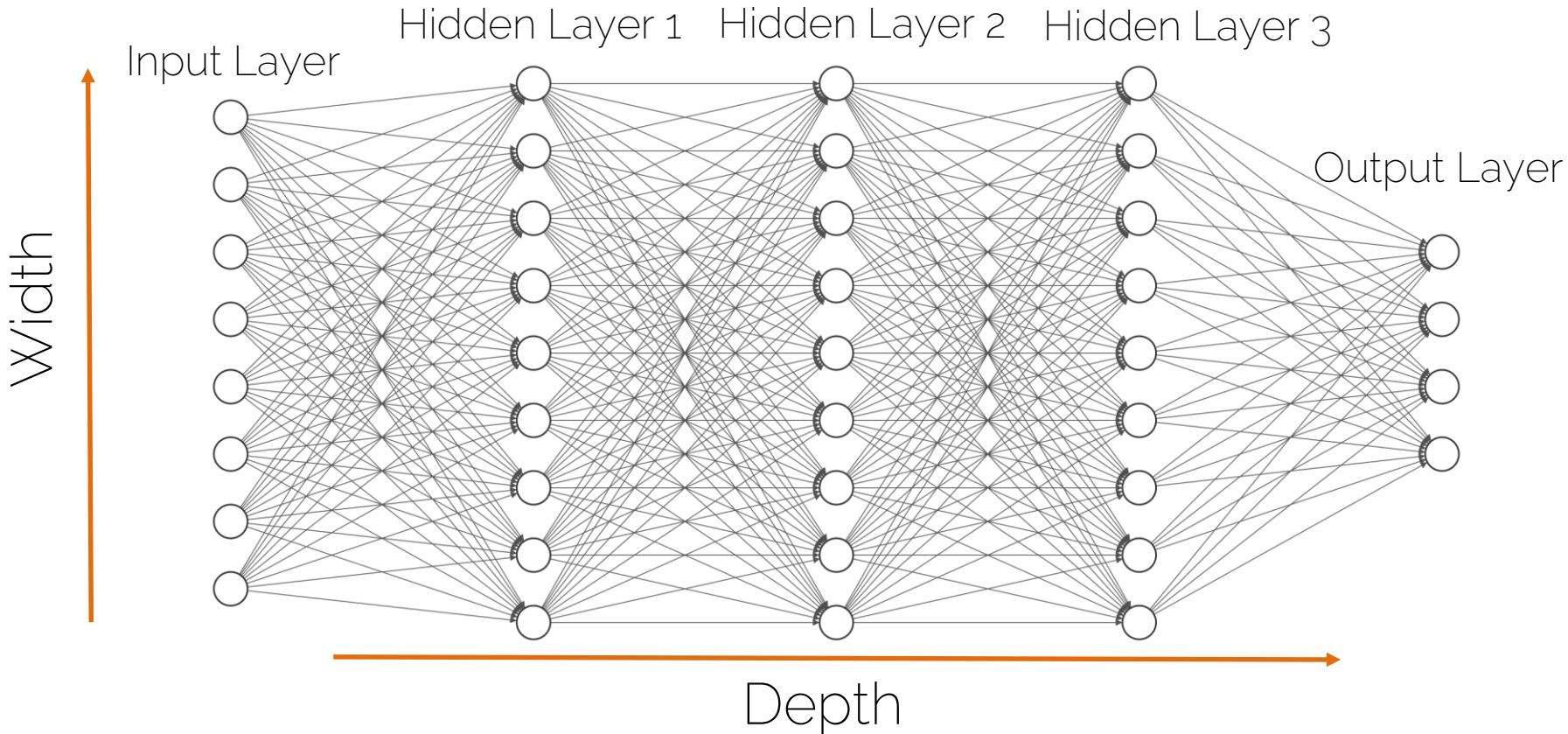
Lecture 4 Recap

Neural Network



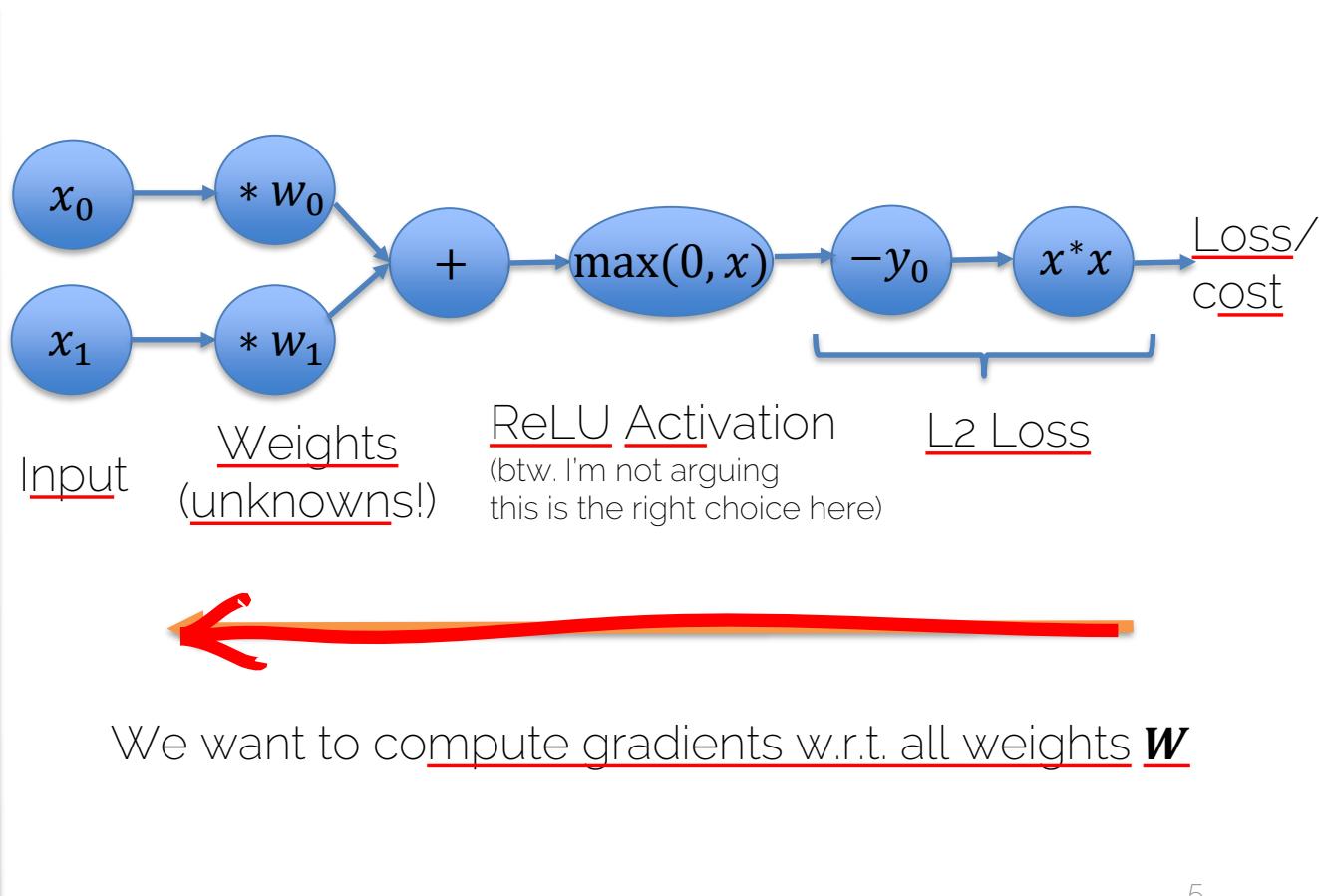
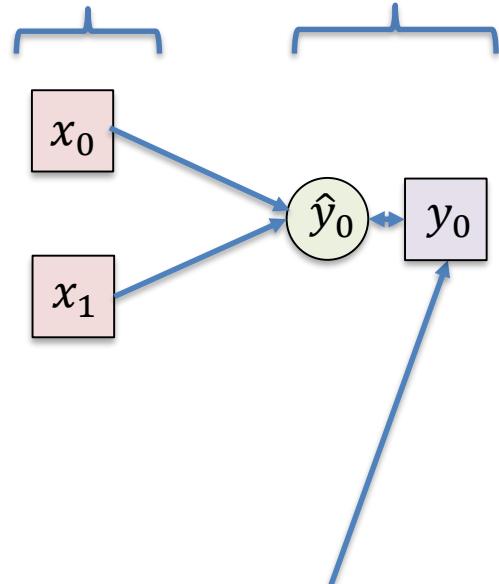
Source: <http://cs231n.github.io/neural-networks-1/>

Neural Network

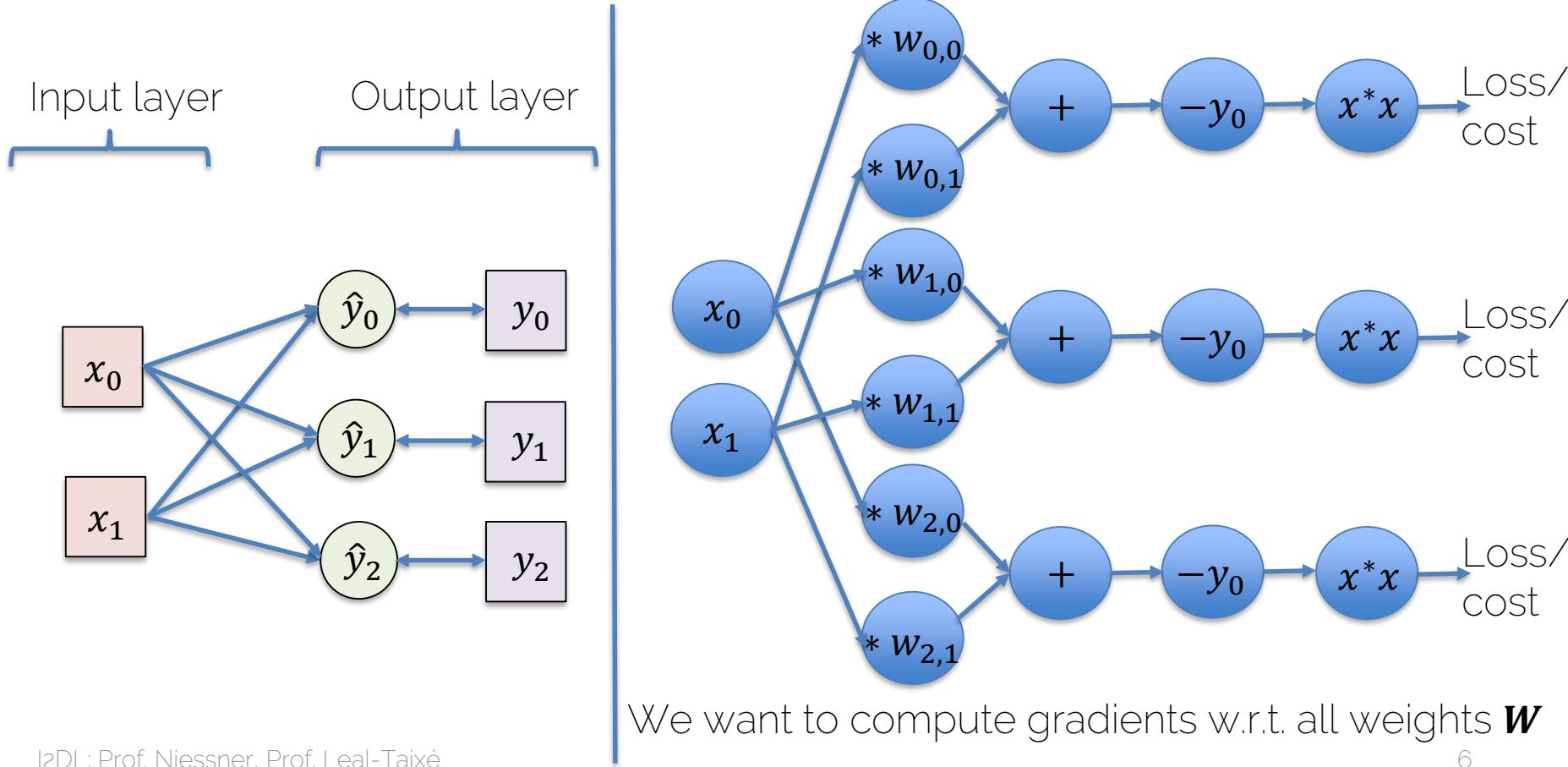


Compute Graphs → Neural Networks

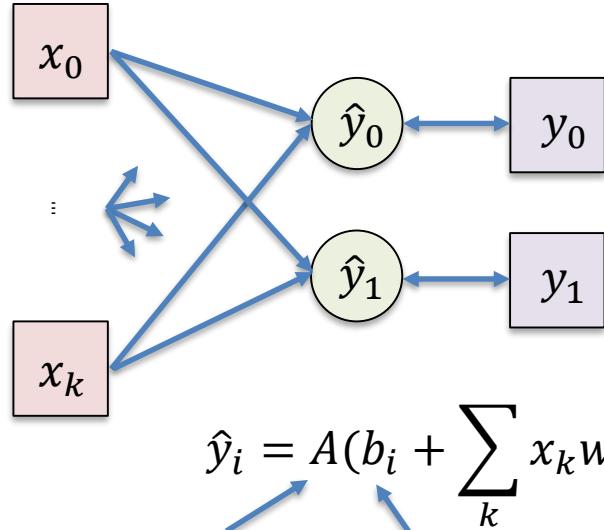
Input layer Output layer



Compute Graphs → Neural Networks



Compute Graphs → Neural Networks



Activation function
bias

NB

We want to compute gradients w.r.t.
all weights **w** **AND** all biases **b**

Goal: We want to compute gradients of the **loss function** **L** w.r.t. all weights **w**

$$L = \sum_i L_i$$

L: sum over **loss per sample**, e.g.
L2 loss → simply sum up squares:

$$L_i = (\hat{y}_i - y_i)^2$$

→ use **chain rule** to compute partials

$$\frac{\partial L}{\partial w_{i,k}} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial w_{i,k}}$$

Summary

- We have
 - (Directional) compute graph
 - Structure graph into layers
 - Compute partial derivatives w.r.t. weights (unknowns)
- Next
 - Find weights based on gradients

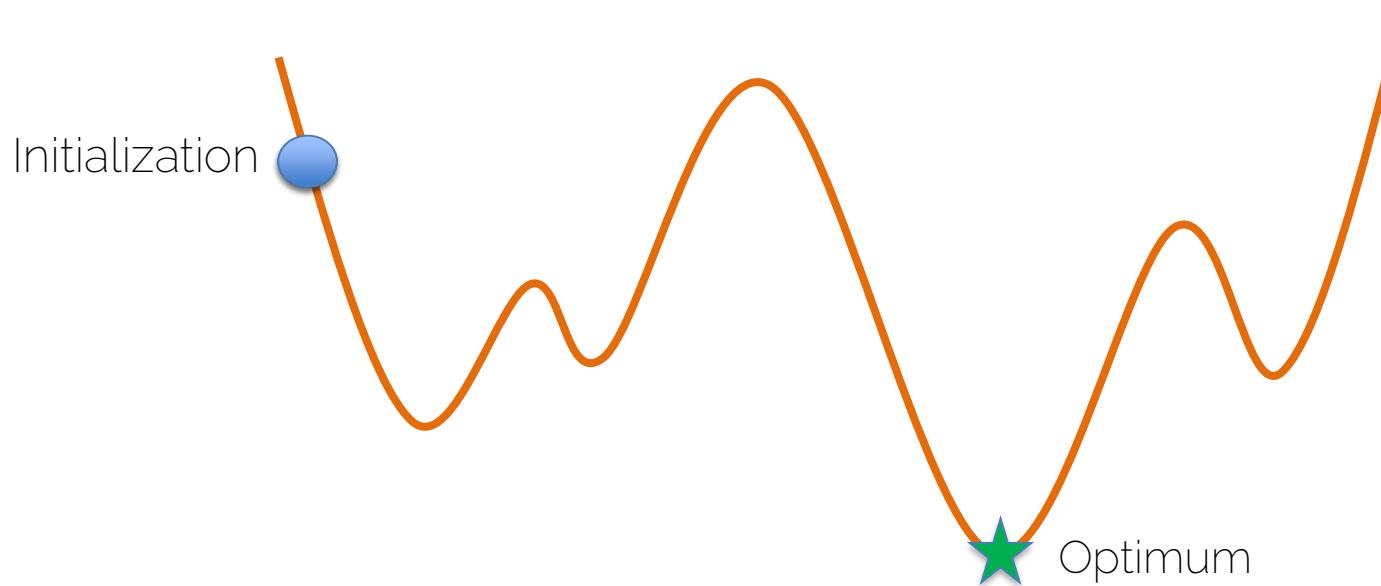
$$\nabla_{\mathbf{W}} f_{\{x,y\}}(\mathbf{W}) = \begin{bmatrix} \frac{\partial f}{\partial w_{0,0,0}} \\ \dots \\ \frac{\partial f}{\partial w_{l,m,n}} \\ \dots \\ \frac{\partial f}{\partial b_{l,m}} \end{bmatrix}$$

Gradient step:
 $\mathbf{W}' = \mathbf{W} - \alpha \nabla_{\mathbf{W}} f_{\{x,y\}}(\mathbf{W})$

Optimization

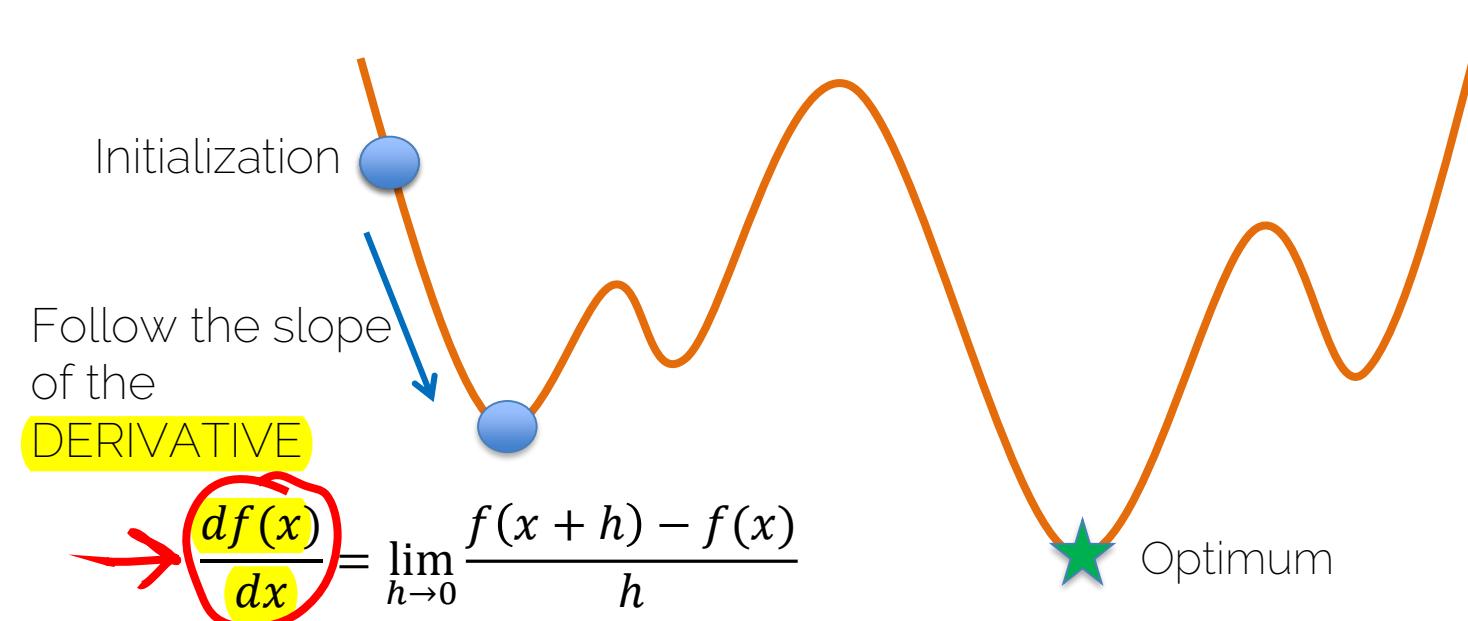
Gradient Descent

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



Gradient Descent

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



Understand

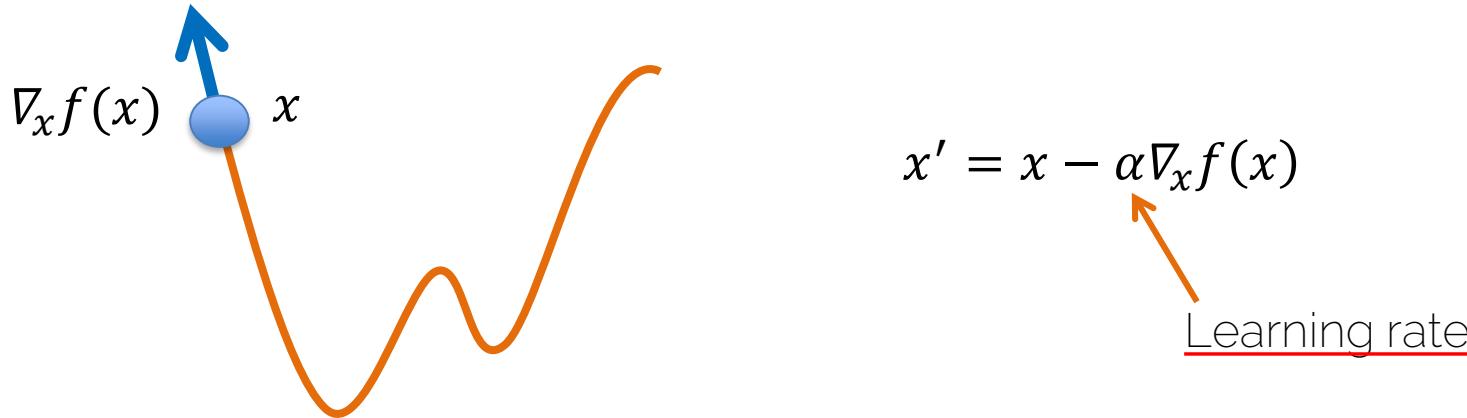
Gradient Descent

- 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



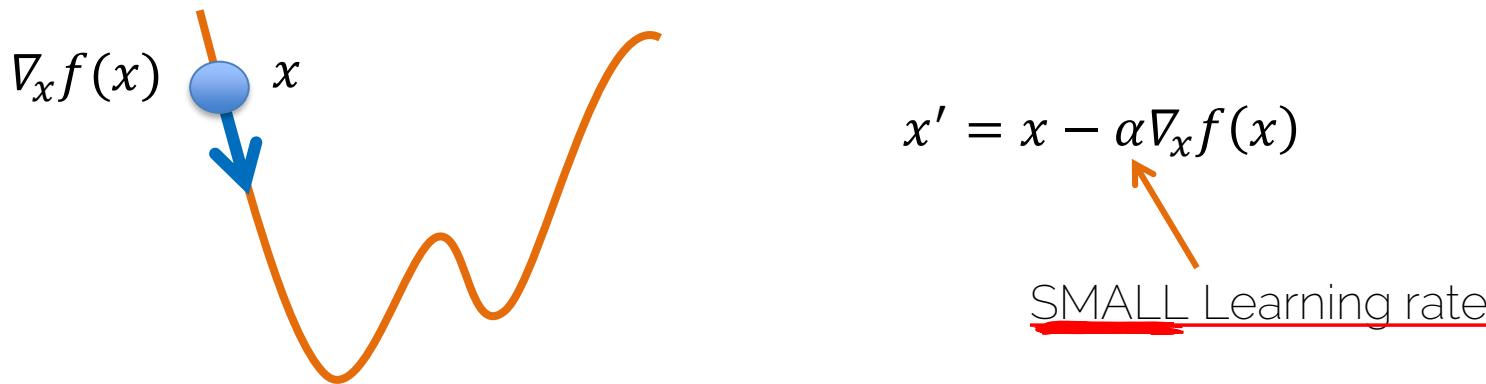
Gradient Descent

- 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



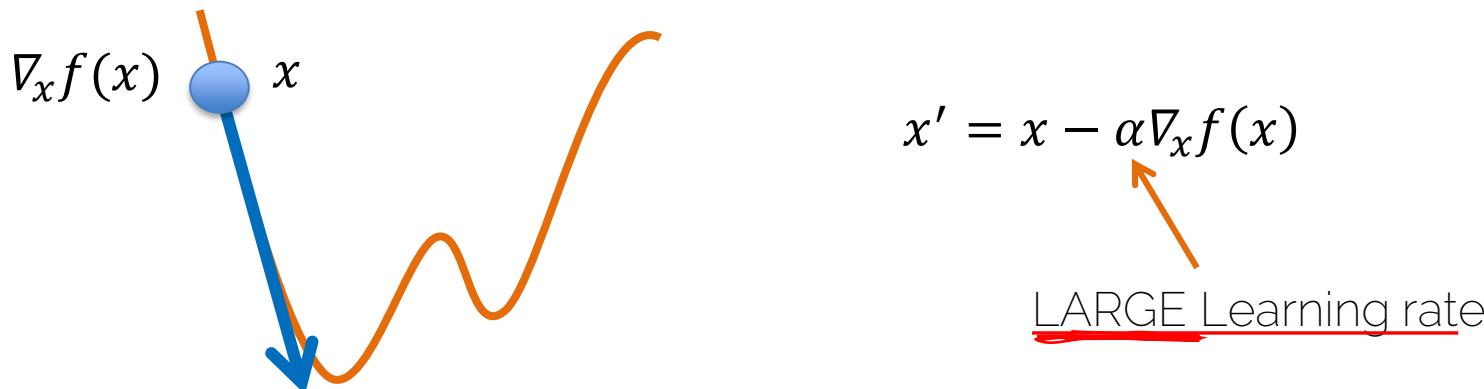
Gradient Descent

- 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

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

Initialization

What is the gradient when we reach this point?

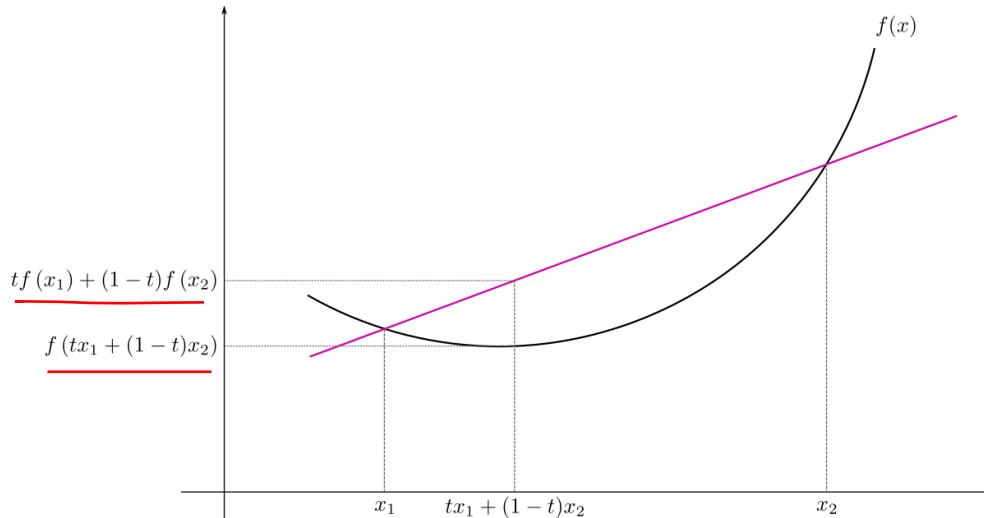
Optimum



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

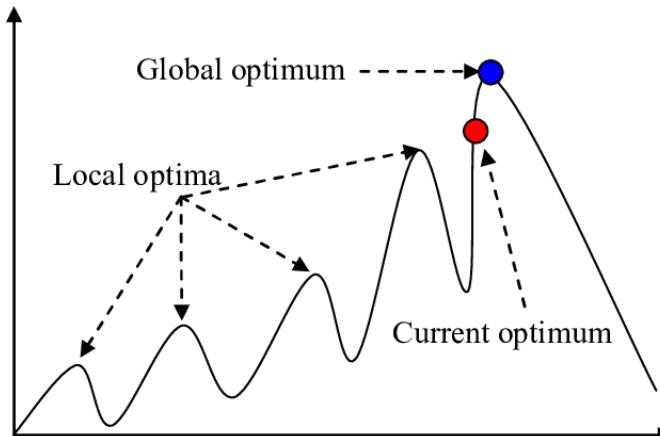
If line/plane segment between any two points lies above or on the graph

Understand

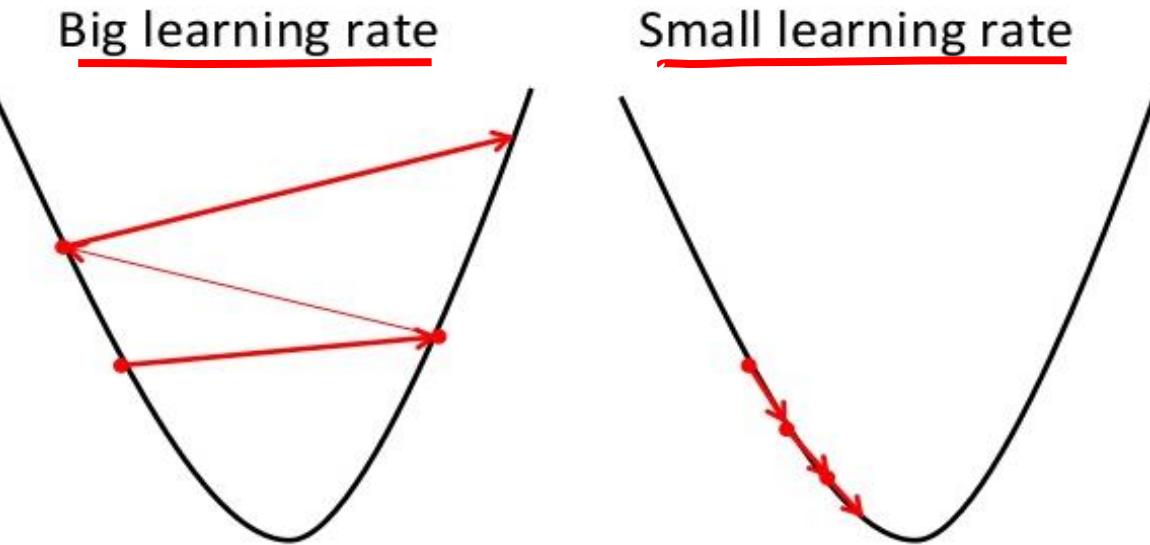


Neural networks are non-convex

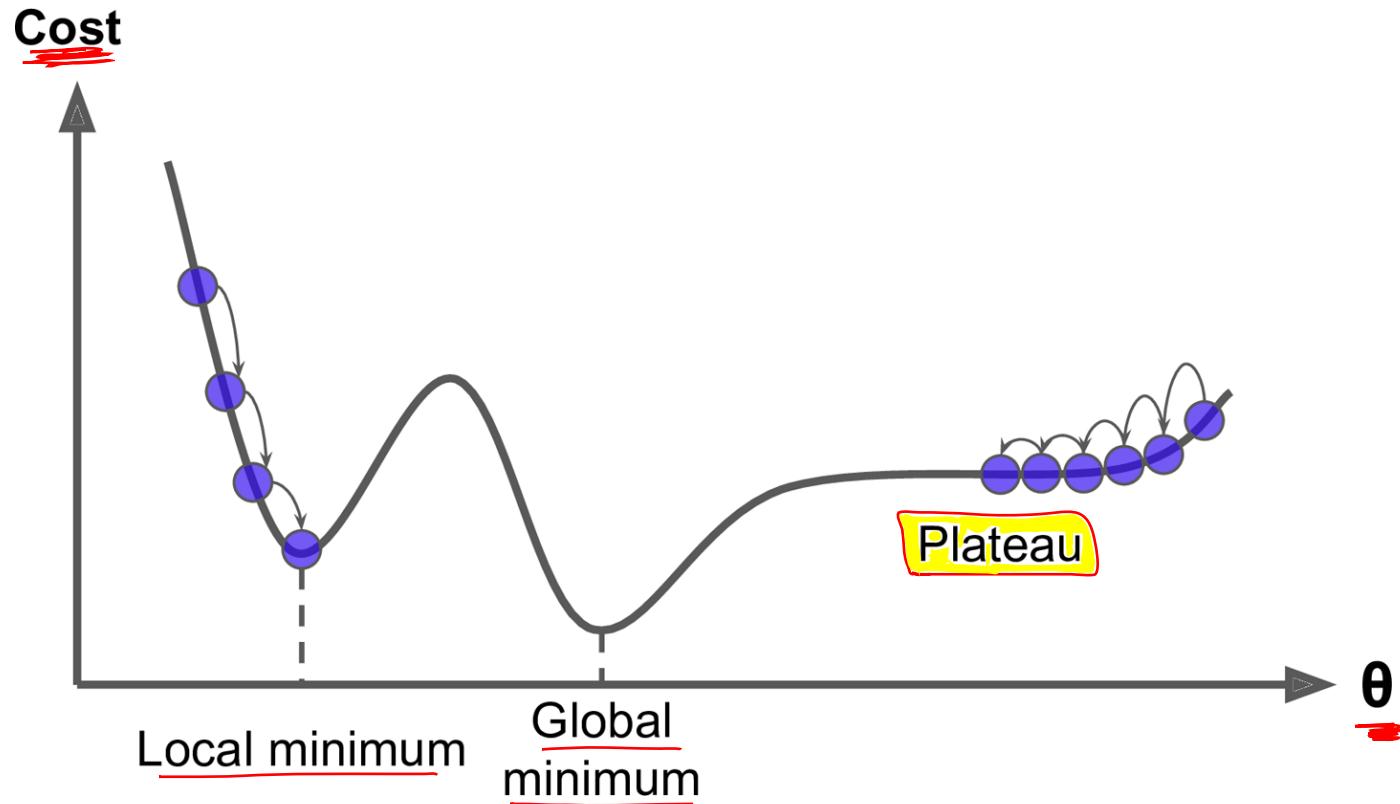
- many (different) local minima
- no (practical) way to say which is globally optimal



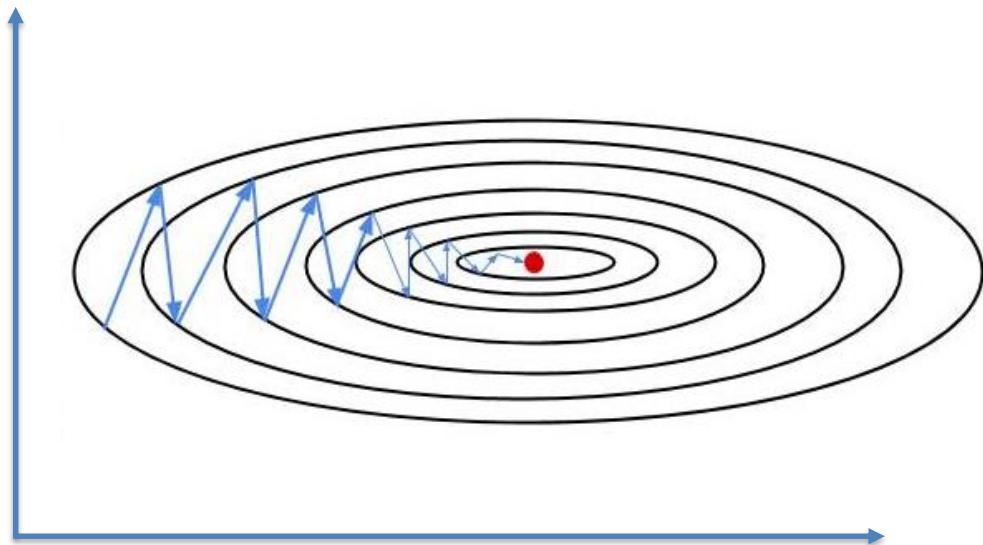
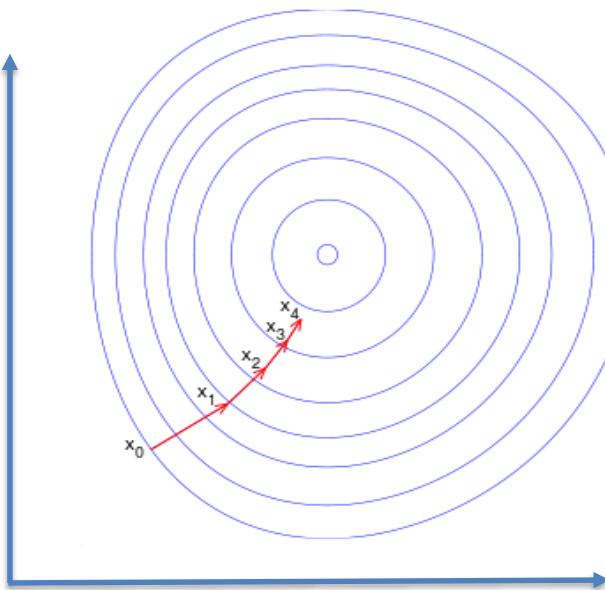
Source: Li, Qi. (2006). Challenging Registration of Geologic Image Data



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



Gradient Descent: Multiple Dimensions



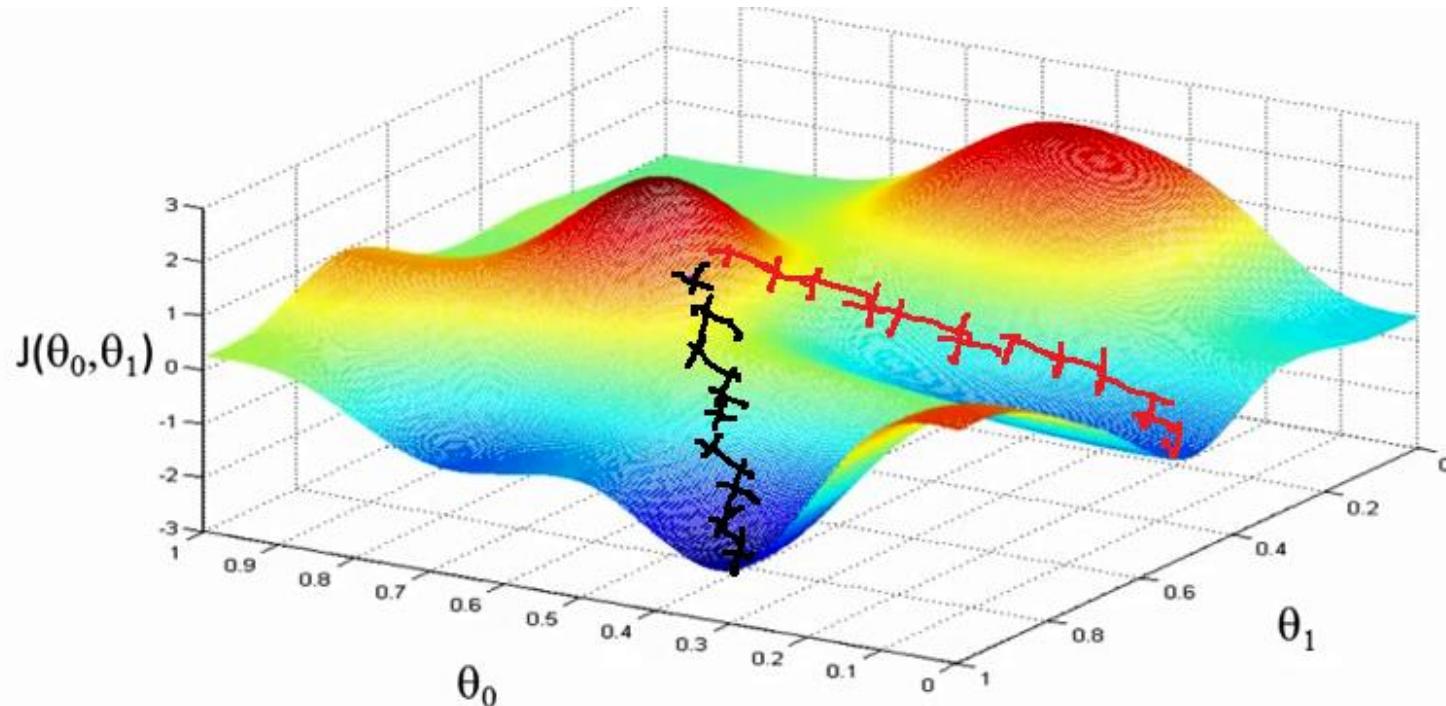
NB

grad shall be
orthogonal to isolines
along the way

Source: builtin.com/data-science/gradient-descent

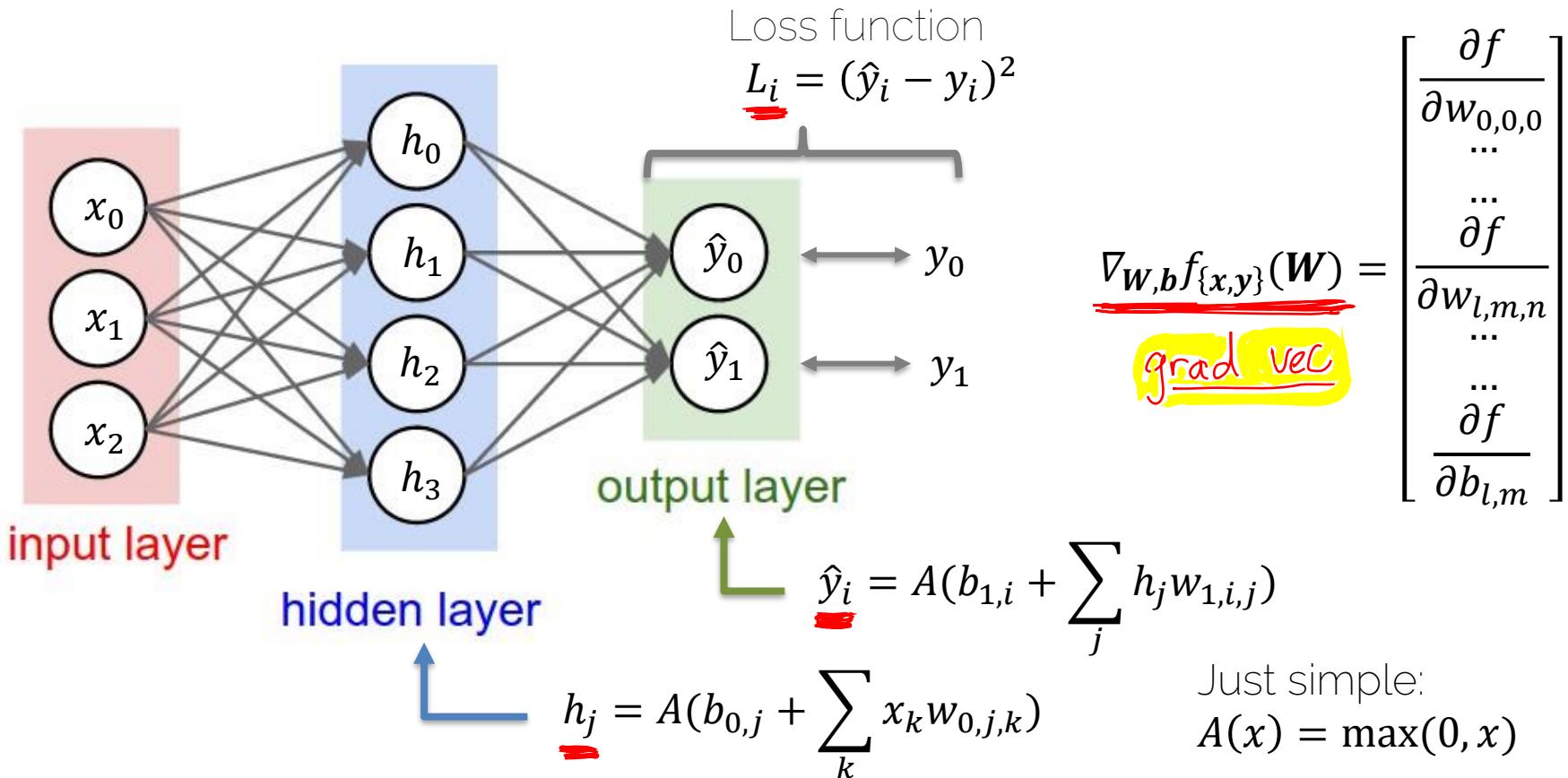
Various ways to visualize...

Gradient Descent: Multiple Dimensions



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

Gradient Descent for Neural Networks



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: $|\boldsymbol{\theta}^{k+1} - \boldsymbol{\theta}^k| < \epsilon$

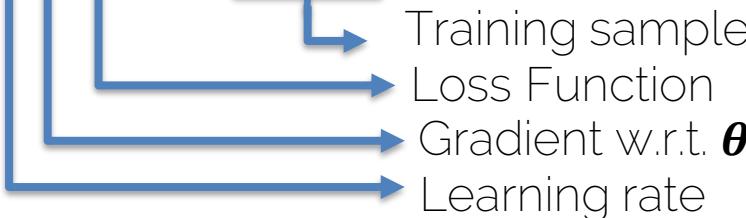
Update step

* Of course, it's less ideal to train
Simple 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)



NB

- $\nabla_{\theta} L_i(\theta^k, x_i, y_i)$ computed via backpropagation

- Typically: $\dim(\nabla_{\theta} L_i(\theta^k, x_i, y_i)) = \dim(\theta) \gg 1 \text{ million}$

looping over all layers
over all parameters

Gradient Descent: Multiple Training Samples

- Given a loss function L and multiple (n) training samples $\{x_i, y_i\}$
- Find best model parameters $\theta = \{W, b\}$

- Cost $L = \frac{1}{n} \sum_{i=1}^n L_i(\theta, x_i, y_i)$

↙ - $\theta = \arg \min L$

[Global Loss]

- Update step for multiple samples

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

- Gradient is average / sum over residuals

Sum is over all training Samples

$$\nabla_{\theta} L(\theta^k, x_{\{1..n\}}, y_{\{1..n\}}) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} L_i(\theta^k, x_i, y_i)$$

Reminder: this comes from backprop.

NB

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

* An update Step of the grad shall be repeated till convergence

Remember

Side Note: Optimal Learning Rate

Can compute optimal learning rate α using Line Search
(optimal for a given set)

1. Compute gradient: $\nabla_{\theta} L = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} L_i$
2. Optimize for optimal step α :

$$\arg \min_{\alpha} L(\underbrace{\theta^k - \alpha \nabla_{\theta} L}_{\theta^{k+1}})$$

3. $\theta^{k+1} = \theta^k - \alpha \nabla_{\theta} L$

Not that practical for DL since we need to solve huge system every step...

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

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

- 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 $m \ll n$

$$B_i = \{\{\mathbf{x}_1, \mathbf{y}_1\}, \{\mathbf{x}_2, \mathbf{y}_2\}, \dots, \{\mathbf{x}_m, \mathbf{y}_m\}\}$$

$$\{B_1, B_2, \dots, B_{n/m}\}$$



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 $n = 2^{20}$ (about 1 million) images
 - With batch size $m = 64$: $B_1 \dots n/m = B_1 \dots 16,384$ minibatches

(Epoch = complete pass through training set)

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

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

k now refers to k -th iteration

one forward
&
one backward
for each batch size

m training samples in the current minibatch

Gradient for the k -th minibatch

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.

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

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

Mathematically, one can tell
when it will converge

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

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

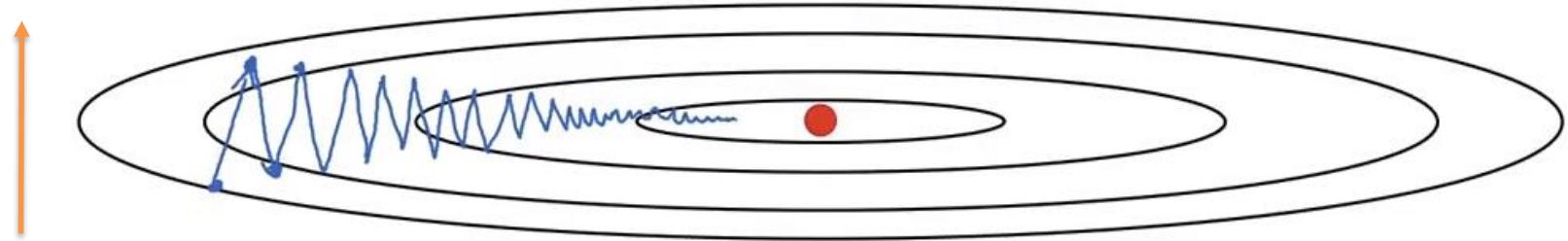
- 1) $\alpha_n \geq 0, \forall n \geq 0$
- 2) $\sum_{n=1}^{\infty} \alpha_n = \infty$
- 3) $\sum_{n=1}^{\infty} \alpha_n^2 < \infty$
- 4) $F(\theta)$ is strictly convex

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



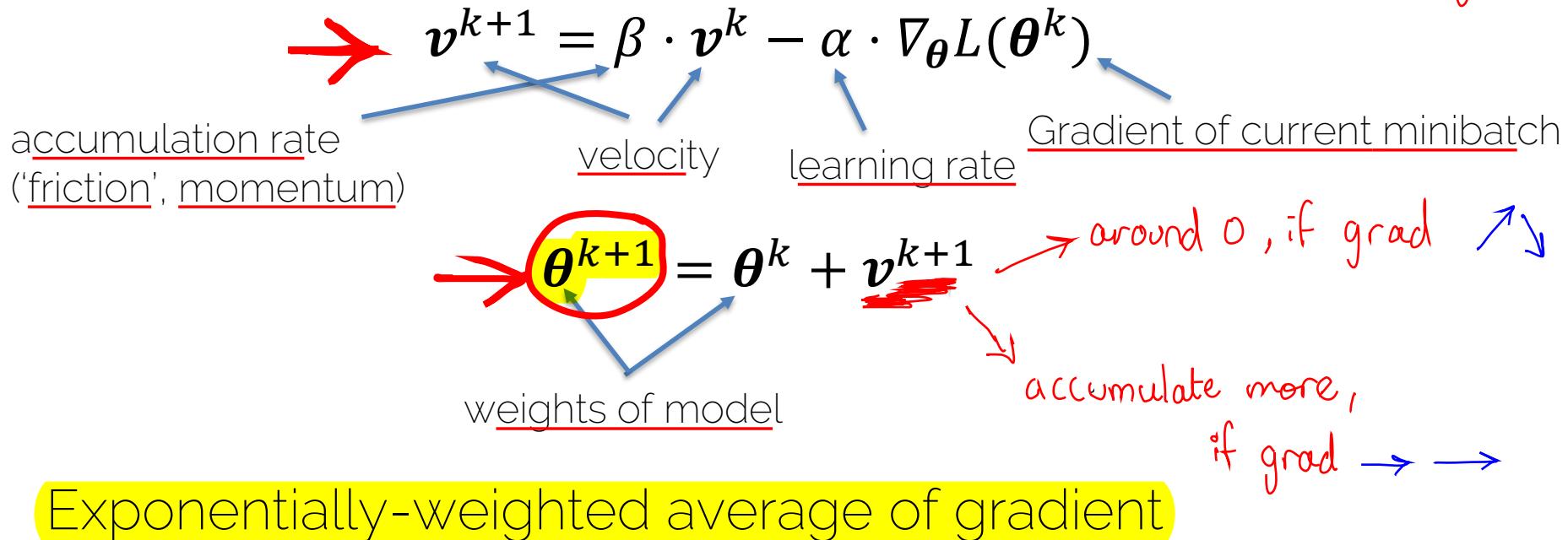
Source: A. Ng

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

i.e history of all grads

First Momentum
mean of grads

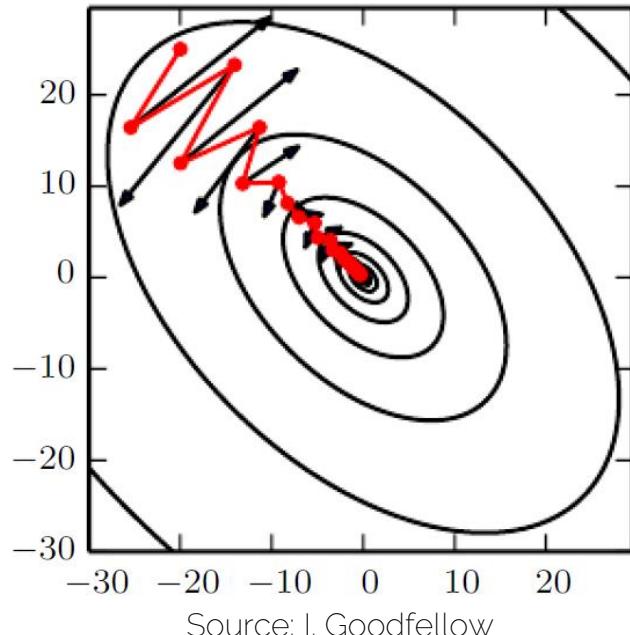


Exponentially-weighted average of gradient

Important: velocity v^k is vector-valued! → on partial deriv basis;

in order to accumulate
grad in each dim

* Needless to say , setting β to Zero is just naive GD .



Very Important

idea

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

Speeding things up

Hyperparameters are α, β
 β is often set to 0.9

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

- Can it overcome local minima? YES.



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

Nesterov Momentum

- Look-ahead momentum

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

directly looking ahead
from prev velo v^k

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

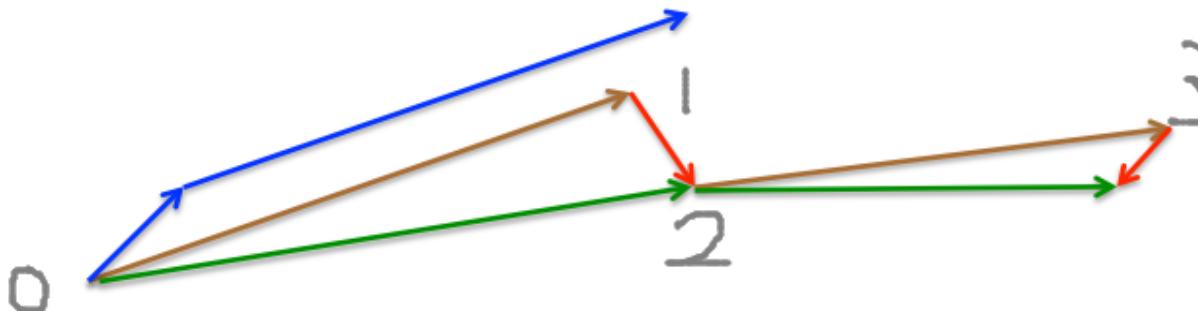
$$\rightarrow \theta^{k+1} = \theta^k + v^{k+1}$$

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.

* All those tricks are computationally ~~in~~expensive



- 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.



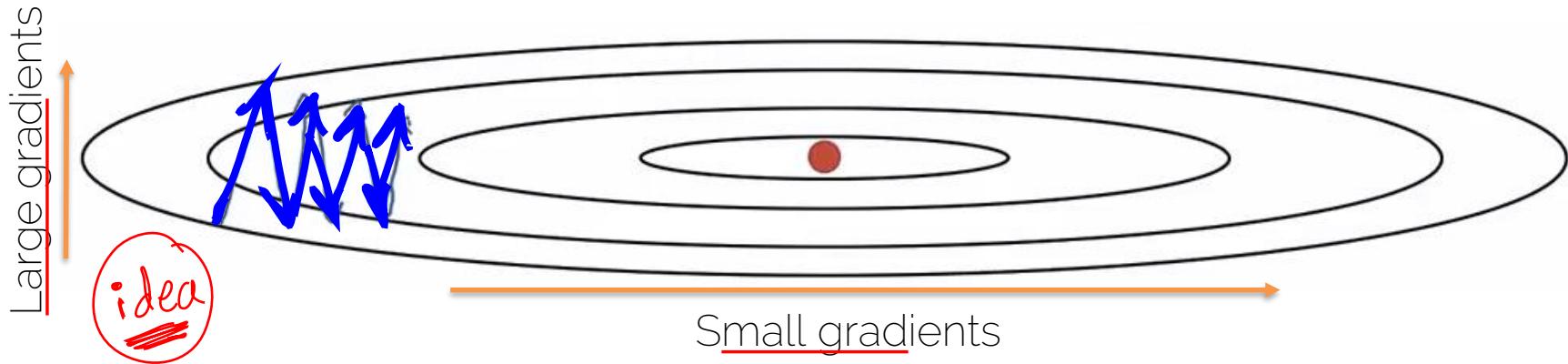
brown vector = jump, red vector = correction, green vector = accumulated gradient

blue vectors = standard momentum

Source: G. Hinton

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

Root Mean Squared Prop (RMSProp)



to dampen this jittering

Source: Andrew. Ng

- RMSProp divides the learning rate by an exponentially-decaying average of squared gradients.

Hinton et al. "Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude." COURSERA: Neural networks for machine learning 4.2 (2012): 26-31.



$$\mathbf{s}^{k+1} = \beta \cdot \mathbf{s}^k + (1 - \beta) [\nabla_{\theta} L \circ \nabla_{\theta} L]$$



~~$$\theta^{k+1} = \theta^k - \alpha \cdot \frac{\nabla_{\theta} L}{\sqrt{\mathbf{s}^{k+1}} + \epsilon}$$~~

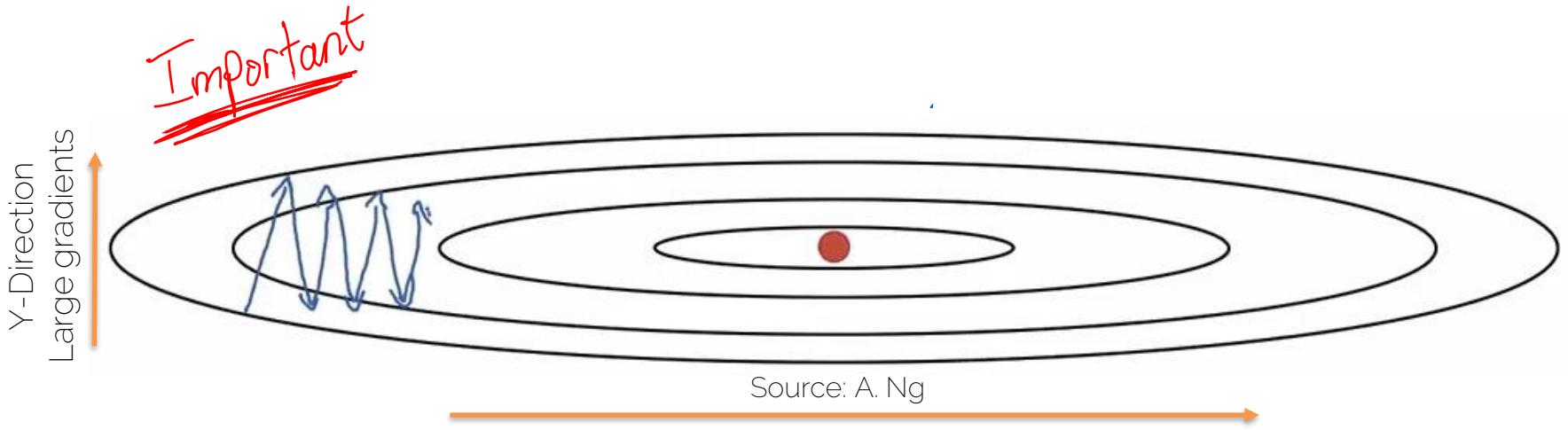
Element-wise multiplication

Hyperparameters: α , β , ϵ

Needs tuning!

Often 0.9

Typically 10^{-8}



(Uncentered) variance of gradients

→ second momentum

$$\boxed{s^{k+1} = \beta \cdot s^k + (1 - \beta)[\nabla_{\theta} L \circ \nabla_{\theta} L]}$$

We're dividing by square gradients:

- Division in Y-Direction will be large
- Division in X-Direction will be small

$$\theta^{k+1} = \theta^k - \alpha \cdot \frac{\nabla_{\theta} L}{\sqrt{s^{k+1}} + \epsilon}$$

Can increase learning rate!



Dampening the oscillations for high-variance directions

- Can use faster learning rate because it is less likely to diverge
 - Speed up learning speed
 - Second moment

Adaptive Moment Estimation (Adam)

Idea : Combine Momentum and RMSProp

$$\mathbf{m}^{k+1} = \beta_1 \cdot \mathbf{m}^k + (1 - \beta_1) \nabla_{\theta} L(\theta^k)$$

First momentum:
mean of gradients

$$\mathbf{v}^{k+1} = \beta_2 \cdot \mathbf{v}^k + (1 - \beta_2) [\nabla_{\theta} L(\theta^k) \circ \nabla_{\theta} L(\theta^k)]$$

$$\theta^{k+1} = \theta^k - \alpha \cdot \frac{\mathbf{m}^{k+1}}{\sqrt{\mathbf{v}^{k+1}} + \epsilon}$$

Note : This is not the
update rule of Adam

Second momentum:
variance of gradients

Q. What happens at $k = 0$?

A. We need bias correction as $\mathbf{m}^0 = \mathbf{0}$ and $\mathbf{v}^0 = \mathbf{0}$

Adam : Bias Corrected

- Combines Momentum and RMSProp

$$\mathbf{m}^{k+1} = \beta_1 \cdot \mathbf{m}^k + (1 - \beta_1) \nabla_{\theta} L(\theta^k) \quad \mathbf{v}^{k+1} = \beta_2 \cdot \mathbf{v}^k + (1 - \beta_2) [\nabla_{\theta} L(\theta^k) \circ \nabla_{\theta} L(\theta^k)]$$

* \mathbf{m}^k and \mathbf{v}^k are initialized with zero
→ bias towards zero
→ Need bias-corrected moment updates

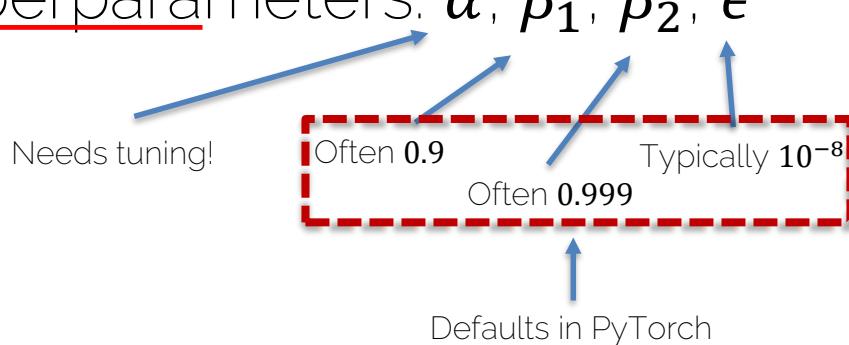
Update rule of Adam

$$\hat{\mathbf{m}}^{k+1} = \frac{\mathbf{m}^{k+1}}{1 - \beta_1^{k+1}} \quad \hat{\mathbf{v}}^{k+1} = \frac{\mathbf{v}^{k+1}}{1 - \beta_2^{k+1}} \quad \rightarrow \theta^{k+1} = \theta^k - \alpha \cdot \frac{\hat{\mathbf{m}}^{k+1}}{\sqrt{\hat{\mathbf{v}}^{k+1}} + \epsilon}$$



Exponentially-decaying mean and variance of gradients (combines first and second order momentum)

- Hyperparameters: $\alpha, \beta_1, \beta_2, \epsilon$



$$\mathbf{m}^{k+1} = \beta_1 \cdot \mathbf{m}^k + (1 - \beta_1) \nabla_{\theta} L(\theta^k)$$

$$\mathbf{v}^{k+1} = \beta_2 \cdot \mathbf{v}^k + (1 - \beta_2) [\nabla_{\theta} L(\theta^k) \circ \nabla_{\theta} L(\theta^k)]$$

$$\hat{\mathbf{m}}^{k+1} = \frac{\mathbf{m}^{k+1}}{1 - \beta_1^{k+1}} \quad \hat{\mathbf{v}}^{k+1} = \frac{\mathbf{v}^{k+1}}{1 - \beta_2^{k+1}}$$

$$\theta^{k+1} = \theta^k - \alpha \cdot \frac{\hat{\mathbf{m}}^{k+1}}{\sqrt{\hat{\mathbf{v}}^{k+1}} + \epsilon}$$

There are a few others...

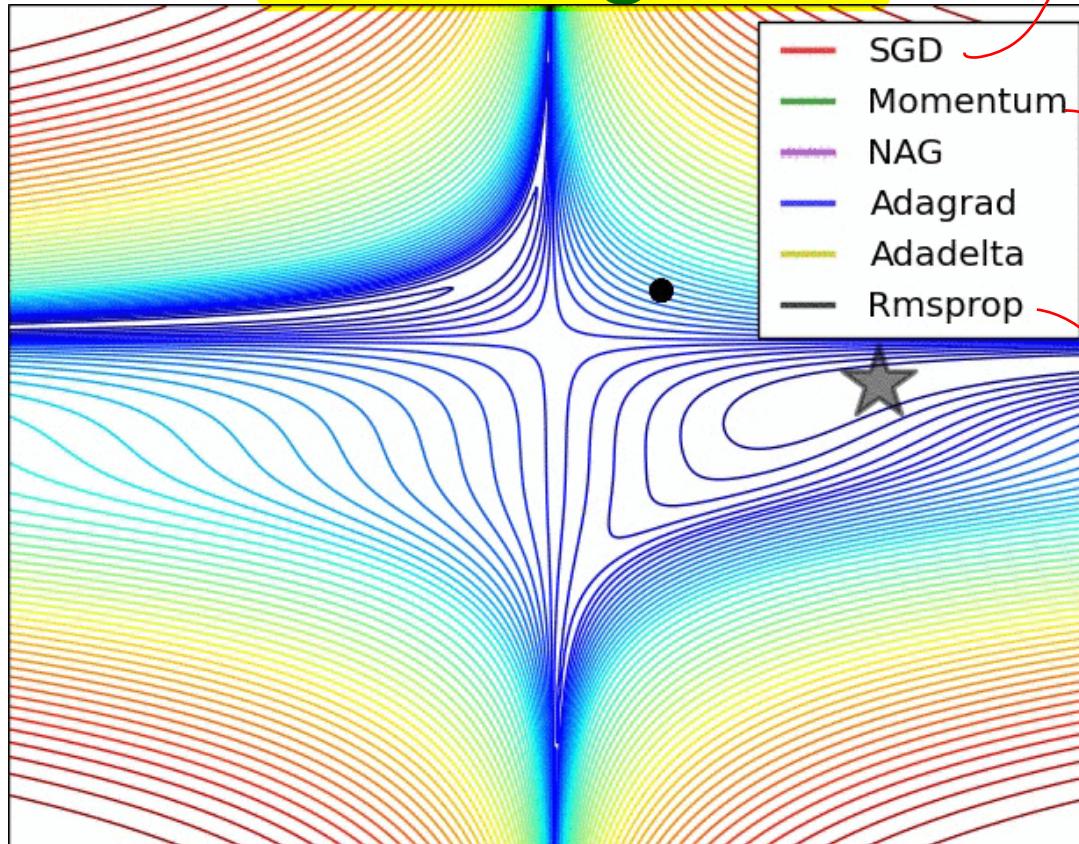
- 'Vanilla' SGD
- Momentum
- RMSProp
- Adagrad
- Adadelta
- AdaMax
- Nada
- AMSGrad



Adam is mostly method
of choice for neural networks!

It's actually fun to play around with SGD updates.
It's easy and you get pretty immediate feedback ☺

Convergence



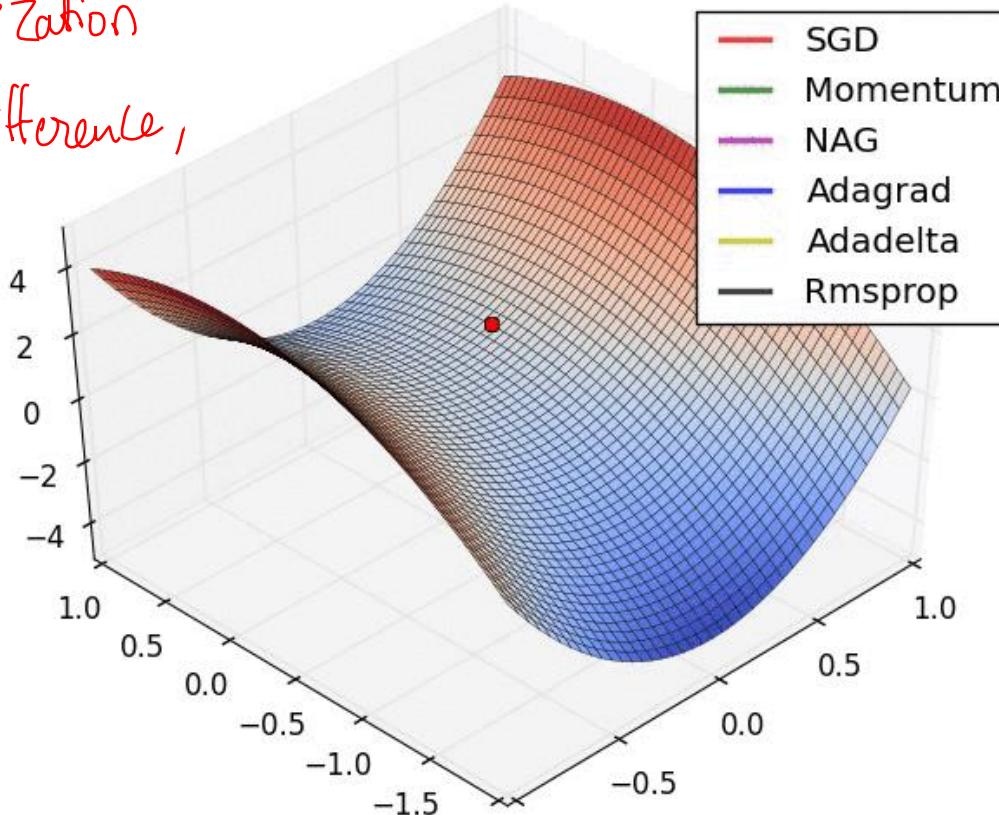
Slow, since there
is no accumulation
of grads

→ usually overshoots

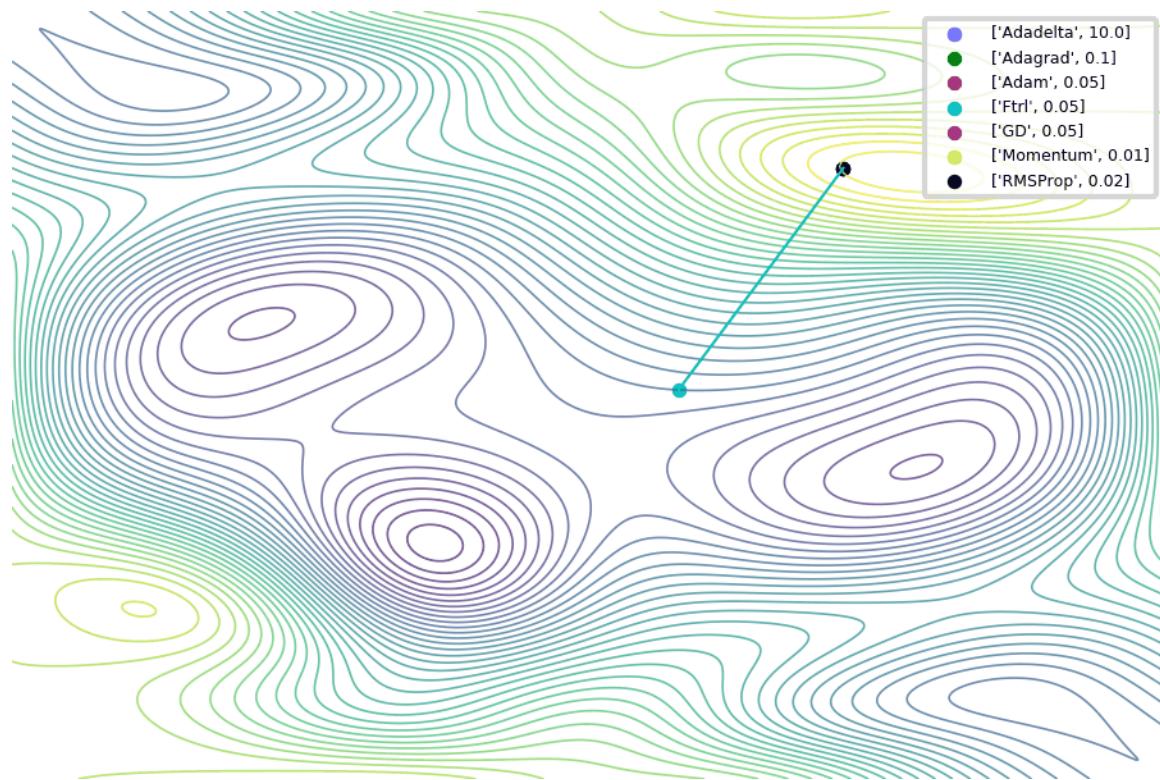
→ pretty good, yet
slower to make
step in forward dir

Convergence

* Choose of optimization makes huge difference, depending on the problem



Convergence



Source: <https://github.com/Jaewan-Yun/optimizer-visualization>

Jacobian and Hessian

- Derivative

$$f: \mathbb{R} \rightarrow \mathbb{R} \quad \frac{df(x)}{dx}$$

- Gradient

$$f: \mathbb{R}^m \rightarrow \mathbb{R} \quad \nabla_x f(x) \quad \left(\frac{df(x)}{dx_1}, \frac{df(x)}{dx_2} \right)$$

- Jacobian

$$f: \mathbb{R}^m \rightarrow \mathbb{R}^n \quad J \in \mathbb{R}^{n \times m}$$

- Hessian

$$f: \mathbb{R}^m \rightarrow \mathbb{R} \quad H \in \mathbb{R}^{m \times m}$$

SECOND
DERIVATIVE

Important

Newton's Method

- Approximate our function by a second-order Taylor series expansion

$$L(\boldsymbol{\theta}) \approx L(\boldsymbol{\theta}_0) + (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_0) + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \mathbf{H}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)$$

First derivative

Second derivative (curvature)

More info: https://en.wikipedia.org/wiki/Taylor_series

- Differentiate and equate to zero

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} L(\theta)$$

Update step

We got rid of the learning rate!

* H is somehow replacing learning rate, scaling grad across every parameter dim

Remember

SGD $\theta_{k+1} = \theta_k - \alpha \nabla_{\theta} L(\theta_k, x_i, y_i)$

(NB)

2nd-order optimizer will do much better than a grad optimizer.

- Differentiate and equate to zero

* No need to invert \mathbf{H} ,
just solve the linear sys

$$\boldsymbol{\theta}^* = \boldsymbol{\theta}_0 - \mathbf{H}^{-1} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$

Update step

NB

Parameters of a
network (millions)

k

Number of
elements in the
Hessian

k^2

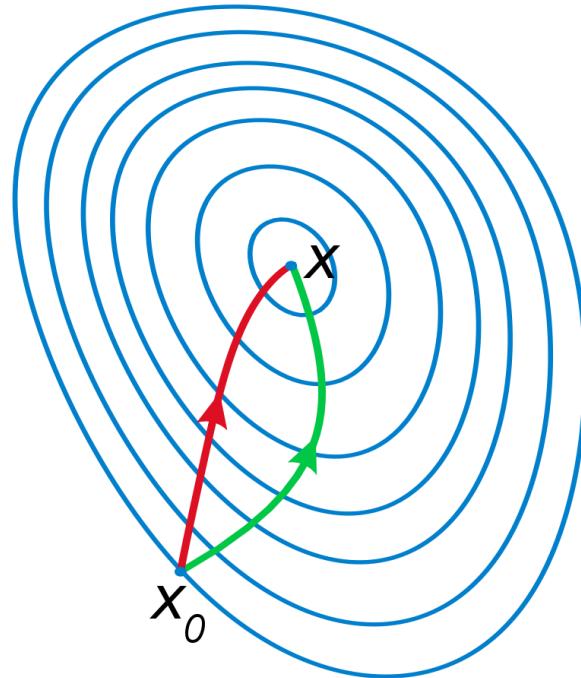
Computational
complexity of 'inversion'
per iteration

$\mathcal{O}(k^3)$

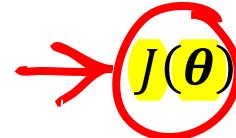
- Gradient Descent (green)

* Newton's method exploits
the curvature to take a
more direct route

fewer, yet more expensive
steps needed



Source: https://en.wikipedia.org/wiki/Newton%27s_method_in_optimization


$$J(\theta) = (\mathbf{y} - \mathbf{X}\theta)^T(\mathbf{y} - \mathbf{X}\theta)$$

Can you apply Newton's
method for linear regression?
What do you get as a result?

YES.

it will solve it
in one step!

Variations of
Newton's

BFGS and L-BFGS

- Broyden-Fletcher-Goldfarb-Shanno algorithm
- Belongs to the family of quasi-Newton methods
- Have an approximation of the inverse of the Hessian

$$\theta^* = \theta_0 - \mathbf{H}^{-1} \nabla_{\theta} L(\theta)$$

Computational Complexity

- BFGS $\mathcal{O}(n^2)$
- Limited memory: L-BFGS $\mathcal{O}(n)$

Gauss-Newton

- $x_{k+1} = x_k - \underline{H_f(x_k)^{-1} \nabla f(x_k)}$
 - 'true' 2nd derivatives are often hard to obtain (e.g., numerics)
 - $H_f \approx 2J_F^T J_F$
- Gauss-Newton (GN):

$$x_{k+1} = x_k - \underline{[2J_F(x_k)^T J_F(x_k)]^{-1} \nabla f(x_k)}$$
- Solve linear system (again, inverting a matrix is unstable):

$$2(J_F(x_k)^T J_F(x_k)) \underbrace{(x_k - x_{k+1})}_{\text{Solve for delta vector}} = \nabla f(x_k)$$

Solve for delta vector

Levenberg

- Levenberg
 - “damped” version of Gauss-Newton:
 - $(J_F(x_k)^T J_F(x_k) + \underline{\lambda} \cdot I) \cdot (x_k - x_{k+1}) = \nabla f(x_k)$
 - The damping factor λ is adjusted in each iteration ensuring:
 - $f(x_k) > f(x_{k+1})$
 - if the equation is not fulfilled increase λ
 - → Trust region

Tikhonov
regularization

* → “Interpolation” between Gauss-Newton (small λ) and Gradient Descent (large λ)

Levenberg-Marquardt

- Levenberg-Marquardt (LM)

$$(J_F(x_k)^T J_F(x_k) + \underline{\lambda} \cdot \underline{\text{diag}(J_F(x_k)^T J_F(x_k))}) \cdot (x_k - x_{k+1}) \\ = \boxed{\nabla f(x_k)}$$

* Instead of a plain Gradient Descent for large λ , scale each component of the gradient according to the curvature.

- Avoids slow convergence in components with a small gradient

Which, What, and When?

- Standard: Adam
- Fallback option: SGD with momentum

* Newton, L-BFGS, GN, LM only if you can do full batch updates (doesn't work well for minibatches!!)

This practically never happens for DL
Theoretically, it would be nice though due to fast convergence

General Optimization

- Linear Systems ($Ax = b$)
 - LU, QR, Cholesky, Jacobi, Gauss-Seidel, CG, PCG, etc.
- Non-linear (gradient-based)
 - Newton, Gauss-Newton, LM, (L)BFGS \leftarrow second order
 - Gradient Descent, SGD \leftarrow first order
- Others
 - Genetic algorithms, MCMC, Metropolis-Hastings, etc.
 - Constrained and convex solvers (Langrange, ADMM, Primal-Dual, etc.)

Please Remember!

- Think about your problem and optimization at hand
- * SGD is specifically designed for minibatch
- When you can, use 2nd order method → it's just faster

* GD or SGD is not a way to solve a linear system!

Next Lecture

- This week:
 - Check exercises
 - Check office hours ☺
- Next lecture
 - Training Neural networks

See you next week 😊

Some References to SGD Updates

- Goodfellow et al. "Deep Learning" (2016),
 - Chapter 8: Optimization
- Bishop "Pattern Recognition and Machine Learning" (2006),
 - Chapter 5.2: Network training (gradient descent)
 - Chapter 5.4: The Hessian Matrix (second order methods)
- <https://ruder.io/optimizing-gradient-descent/index.html>
- PyTorch Documentation (with further readings)
 - <https://pytorch.org/docs/stable/optim.html>