

Deep Learning

A course about theory & practice



Learning as Optimization

Marco Piastra

*About why they did not use
Deep Networks
from the beginning*

Problem: vanishing or exploding Gradients

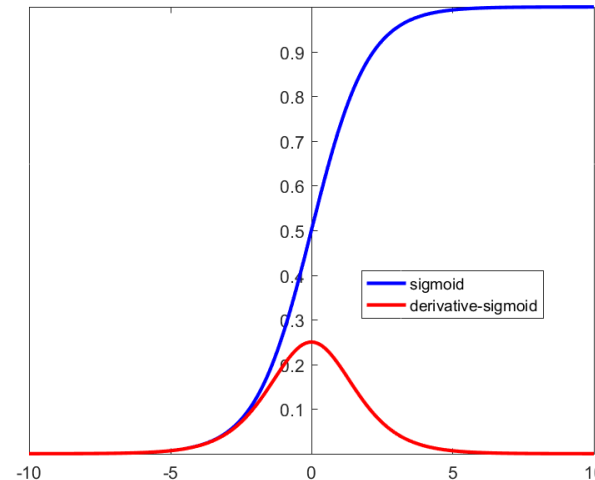
The gradient descent method implies updating the parameters at each step: making sure that the gradient does not either *vanish* or *explode* is not easy

For instance, in

$$\Delta \mathbf{W} = -\eta \frac{\partial L}{\partial \mathbf{W}}(\tilde{y}^{(i)}, y^{(i)})$$

the gradient contains a multiplicative term $\frac{\partial}{\partial x} g(x)$
which can be $\ll 1.0$

e.g. for the sigmoid function:



Problem: vanishing or exploding Gradients

The gradient descent method implies updating the parameters at each step: making sure that the gradient does not *vanish* or *explode* is not easy ...

Consider a special deep network

$$\tilde{y} = \mathbf{w} \cdot g(\mathbf{W}^{[k]} \dots g(\mathbf{W}^{[1]} \mathbf{x} + \mathbf{b}^{[1]}) \dots + \mathbf{b}^{[k]}) + b$$

in which

- g is the identity function and all $\mathbf{b}^{[i]}$ and b are zero;
- all hidden layers have the same size d of the input (i.e., all matrices are square);
- all $\mathbf{W}^{[i]}$ are identical and diagonalizable, with eigenbasis $(\mathbf{e}_1, \dots, \mathbf{e}_d)$.

This means that

$$\begin{aligned} \mathbf{W}^{[k]} \dots \mathbf{W}^{[1]} \mathbf{x} &= \mathbf{W}^k \mathbf{x} = \lambda_1^k \overset{\text{eigenvalue raised to the } k\text{-th power}}{(\mathbf{e}_1 \cdot \mathbf{x})} \mathbf{e}_1 + \dots + \lambda_d^k (\mathbf{e}_d \cdot \mathbf{x}) \mathbf{e}_d \\ &= \lambda_1^k x_1 \mathbf{e}_1 + \dots + \lambda_d^k x_d \mathbf{e}_d \end{aligned}$$

Moral: any $\lambda_i > 1$ leads to explosion while any $\lambda_i < 1$ leads to vanishing

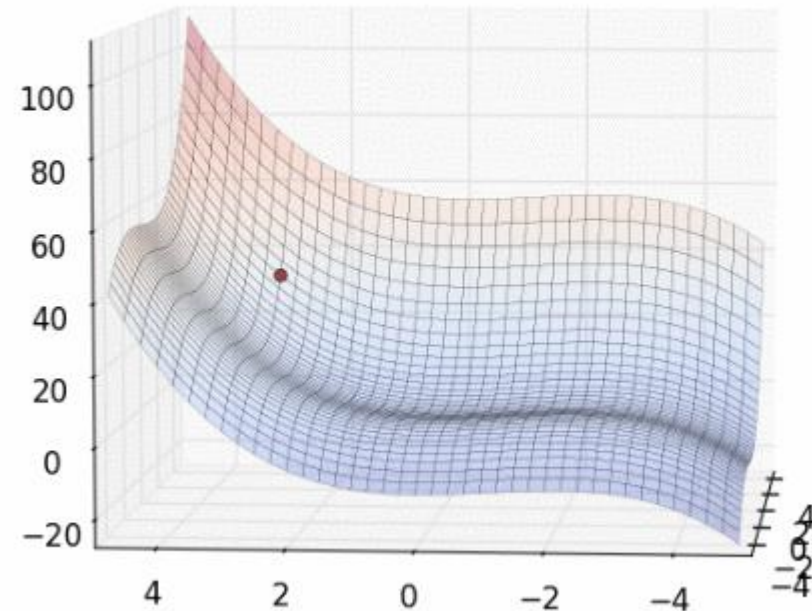
Problem: initial values of the parameters

However, the main problem of training is that of *initial values*...

Gradient Descent can only discover minima that are close to the initial values

*Using deep networks
can only make this problem worse:
intuitively, with deeper networks,
the 'surface' can be even rougher...*

$x=3.00000$, $y=3.00000$, $f(x,y)=34.20000$



[Image from <http://cpmarkchang.logdown.com/posts/434534-optimization-method-momentum>]

Improving optimization

Improving optimization

■ **SGD (or MBGD)**

Standard, decaying learning rate

Update step:

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \underset{\substack{\text{decaying} \\ \text{learning rate}}}{\eta} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(\underset{\substack{\text{mini-batch,} \\ \text{possibly a singleton}}}{B}, \boldsymbol{\vartheta}^{(t-1)})$$

Improving optimization

■ **SGD (or MBGD)**

Standard, decaying learning rate

Update step:

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \underset{\substack{\text{decaying} \\ \text{learning rate}}}{\eta} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(\underset{\substack{\text{mini-batch,} \\ \text{possibly a singleton}}}{B}, \boldsymbol{\vartheta}^{(t-1)})$$

Many different ways to improve performance and speed rate:

- add some *momentum*
- take in account 2^{nd} order derivatives
- make the *learning rate adaptive*

Improving optimization

■ SGD (or MBGD)

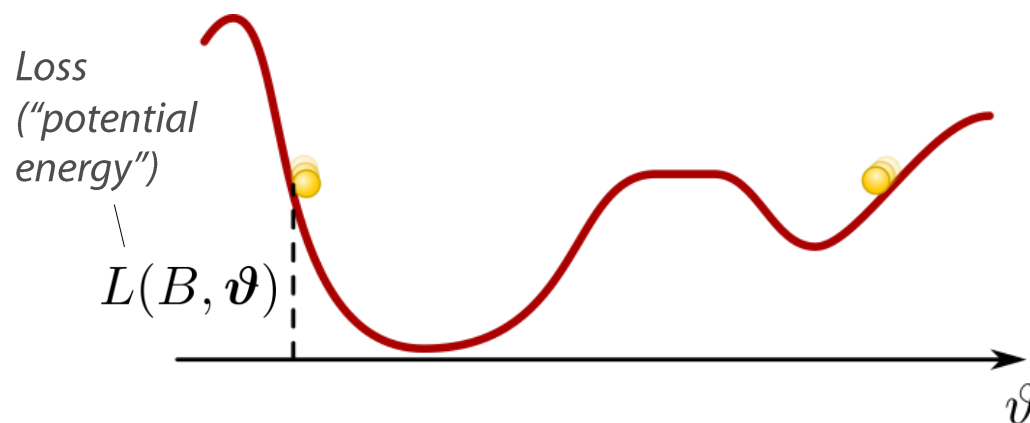
Standard, decaying learning rate


Update step:

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)})$$

$$\boldsymbol{\vartheta}^{(t)} - \boldsymbol{\vartheta}^{(t-1)} = -\eta \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)})$$

“velocity”




$$-\eta \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)})$$

“force felt by the ball”

$$\boldsymbol{f} = -\frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta})$$

“acceleration”

$$\boldsymbol{f} = m\boldsymbol{a}$$

$$\boldsymbol{a} \propto -\frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta})$$

... the gradient directly affects the velocity
(not the position)

Momentum

■ Momentum

"Let the ball run"

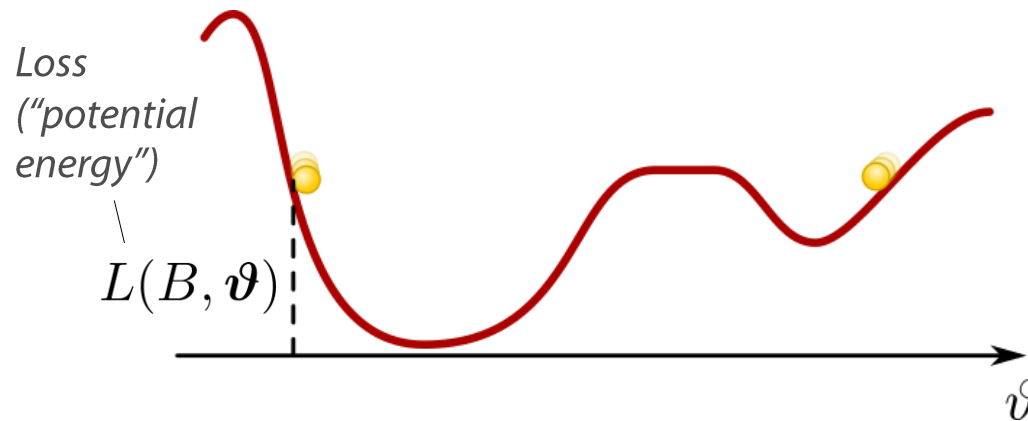
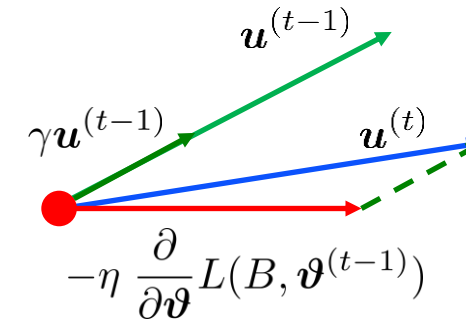
momentum term:

tendency to keep running at the same velocity and direction

$$\mathbf{u}^{(t)} = \gamma \mathbf{u}^{(t-1)} - \eta \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)}), \quad \mathbf{u}^{(0)} = \mathbf{0}$$

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} + \mathbf{u}^{(t)} \quad 0 < \gamma < 1$$

"coefficient of friction"



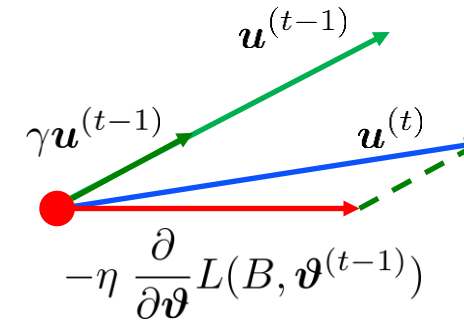
Momentum

■ Momentum

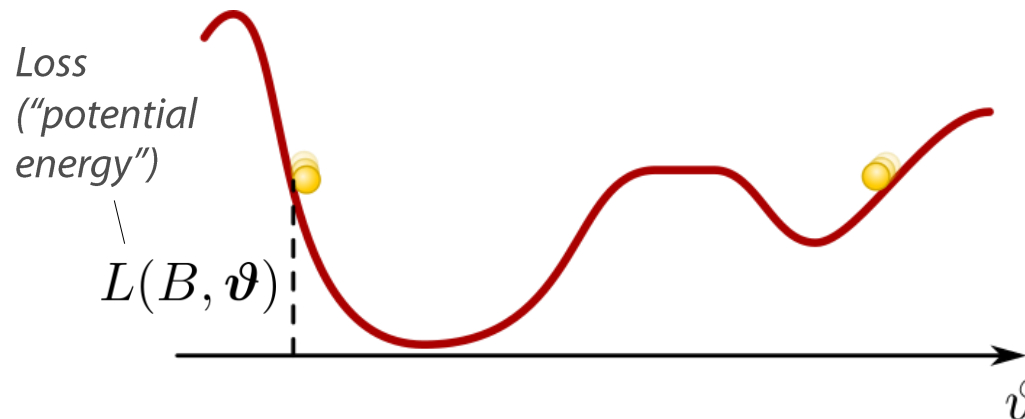
"Let the ball run"

$$\mathbf{u}^{(t)} = \gamma \mathbf{u}^{(t-1)} - \eta \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)}), \quad \mathbf{u}^{(0)} = \mathbf{0}$$

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} + \mathbf{u}^{(t)}$$



Consider $\boldsymbol{\vartheta}$ as a position ...



"velocity"

$$\mathbf{u} := \frac{\partial}{\partial t} \boldsymbol{\vartheta} \approx \boldsymbol{\vartheta}^{(t)} - \boldsymbol{\vartheta}^{(t-1)}$$

"acceleration"

$$\mathbf{a} \approx \mathbf{u}^{(t)} - \mathbf{u}^{(t-1)} \propto -\frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta})$$

... the gradient directly affects the velocity
(not the position)

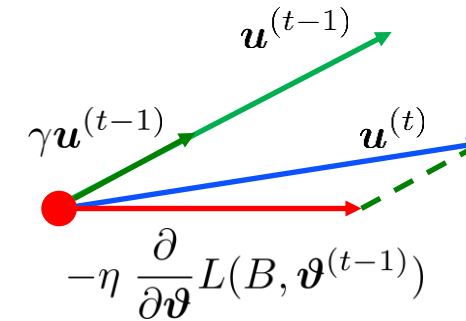
NAG

■ Momentum

"Let the ball run"

$$\mathbf{u}^{(t)} = \gamma \mathbf{u}^{(t-1)} - \eta \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)}), \quad \mathbf{u}^{(0)} = \mathbf{0}$$

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} + \mathbf{u}^{(t)}$$

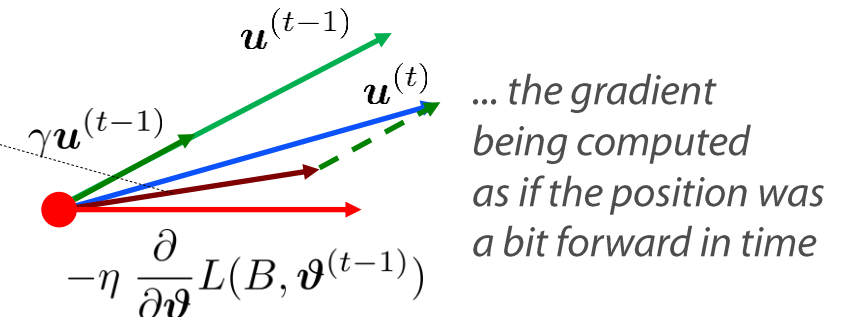


■ Nesterov Accelerated Gradient (NAG)

"Let the ball run but be predictive"

$$\mathbf{u}^{(t)} = \gamma \mathbf{u}^{(t-1)} - \eta \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)} + \gamma \mathbf{u}^{(t-1)})$$

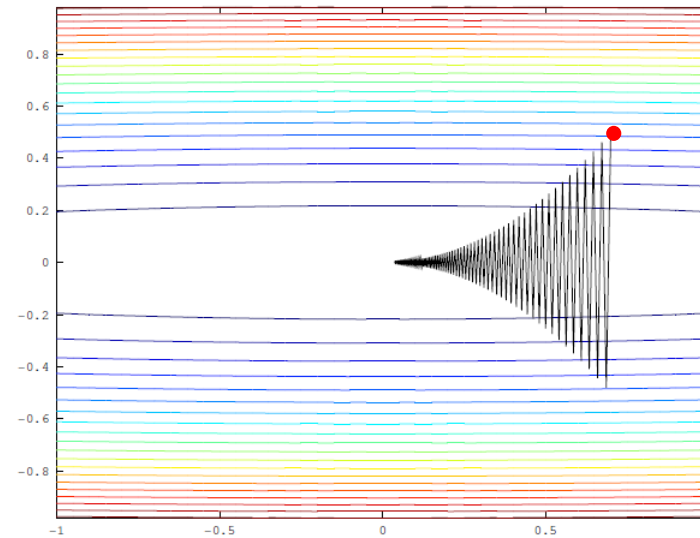
$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} + \mathbf{u}^{(t)}$$



2nd order methods

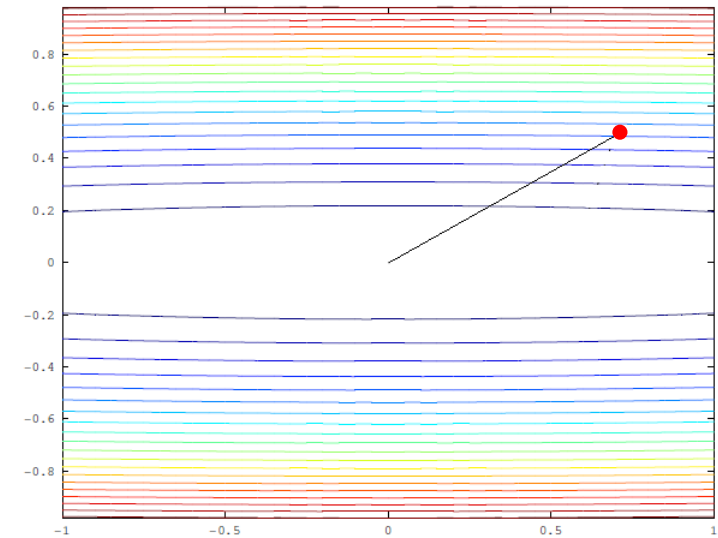
In this example (geometric view)

Gradient Descent



The level curves of
a quadratic form in 2D
are ellipses centered
in the origin

Newton-Raphson



2nd order methods

■ Taylor's expansion

$$L(B, \boldsymbol{\vartheta}) = L(B, \boldsymbol{\vartheta}^{(t-1)}) + \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)}) \cdot (\boldsymbol{\vartheta} - \boldsymbol{\vartheta}^{(t-1)}) + \frac{1}{2} (\boldsymbol{\vartheta} - \boldsymbol{\vartheta}^{(t-1)}) \cdot \mathbf{H} (\boldsymbol{\vartheta} - \boldsymbol{\vartheta}^{(t-1)}) + \dots$$

The Hessian Matrix

All terms in blue are constant

where:

$$\mathbf{H} := \frac{\partial}{\partial \boldsymbol{\vartheta}} \left(\frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)}) \right) \text{ — The Hessian Matrix}$$

■ Differentiate both sides and take $\boldsymbol{\vartheta} = \boldsymbol{\vartheta}^*$ — The argmin

$$\frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^*) = \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)}) + \mathbf{H} (\boldsymbol{\vartheta}^* - \boldsymbol{\vartheta}^{(t-1)})$$

set this equal to 0

then:

$$\boldsymbol{\vartheta}^* - \boldsymbol{\vartheta}^{(t-1)} = -\mathbf{H}^{-1} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)})$$

2nd order methods

- **Gradient Descent**

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)})$$

- **Newton-Raphson's optimization method**

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta \mathbf{H}^{-1} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)})$$

where:

$$\mathbf{H} := \frac{\partial}{\partial \boldsymbol{\vartheta}} \left(\frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)}) \right)$$

Why is the Newton-Raphson's method better than GD?

2nd order methods

■ Newton-Raphson's optimization method

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta \mathbf{H}^{-1} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)}) \quad \mathbf{H} := \frac{\partial}{\partial \boldsymbol{\vartheta}} \left(\frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)}) \right)$$

Example

$$L(B, \boldsymbol{\vartheta}) = \boldsymbol{\vartheta} \cdot \mathbf{A} \boldsymbol{\vartheta}$$

a quadratic form, centered in the origin

where;

$$\mathbf{A} := \begin{bmatrix} a_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & a_d \end{bmatrix}, \quad a_i > 0 \quad \forall i = 1, \dots, d$$

*a diagonal, positive definite matrix
(therefore, L is convex)*

$$\frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}) = 2\mathbf{A} \boldsymbol{\vartheta}$$

$$\mathbf{H} = \frac{\partial}{\partial \boldsymbol{\vartheta}} \left(\frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}) \right) = 2\mathbf{A} \quad \mathbf{H}^{-1} = \frac{1}{2} \mathbf{A}^{-1} = \frac{1}{2} \begin{bmatrix} 1/a_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 1/a_d \end{bmatrix}$$

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta \frac{1}{2} \mathbf{A}^{-1} 2\mathbf{A} \boldsymbol{\vartheta}^{(t-1)} = \boldsymbol{\vartheta}^{(t-1)} - \eta \boldsymbol{\vartheta}^{(t-1)} = (1 - \eta) \boldsymbol{\vartheta}^{(t-1)}$$

What??

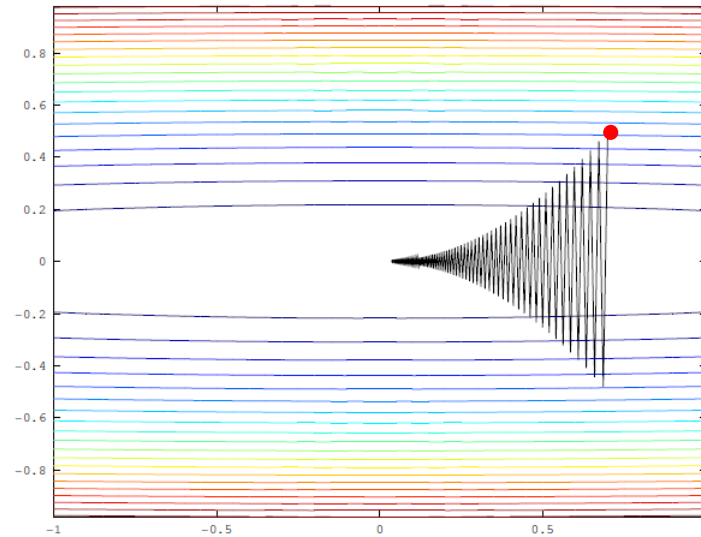
2nd order methods

In this example (geometric view)

$$L(B, \boldsymbol{\vartheta}) = \boldsymbol{\vartheta} \cdot \mathbf{A} \boldsymbol{\vartheta} \quad \text{with} \quad \mathbf{A} := \begin{bmatrix} a_1 & 0 \\ 0 & a_2 \end{bmatrix}, \quad a_1 \ll a_2$$

Gradient Descent

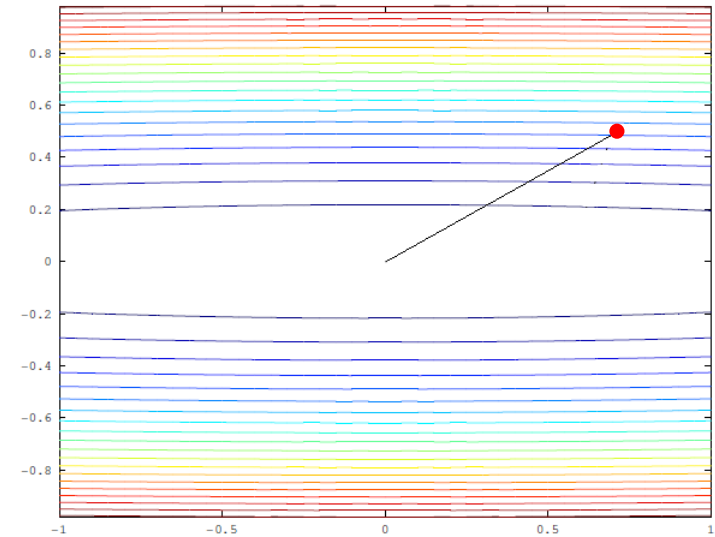
$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta 2\mathbf{A} \boldsymbol{\vartheta}^{(t-1)}$$



The level curves of
a quadratic form in 2D
are ellipses centered
in the origin

Newton-Raphson

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta \boldsymbol{\vartheta}^{(t-1)}$$

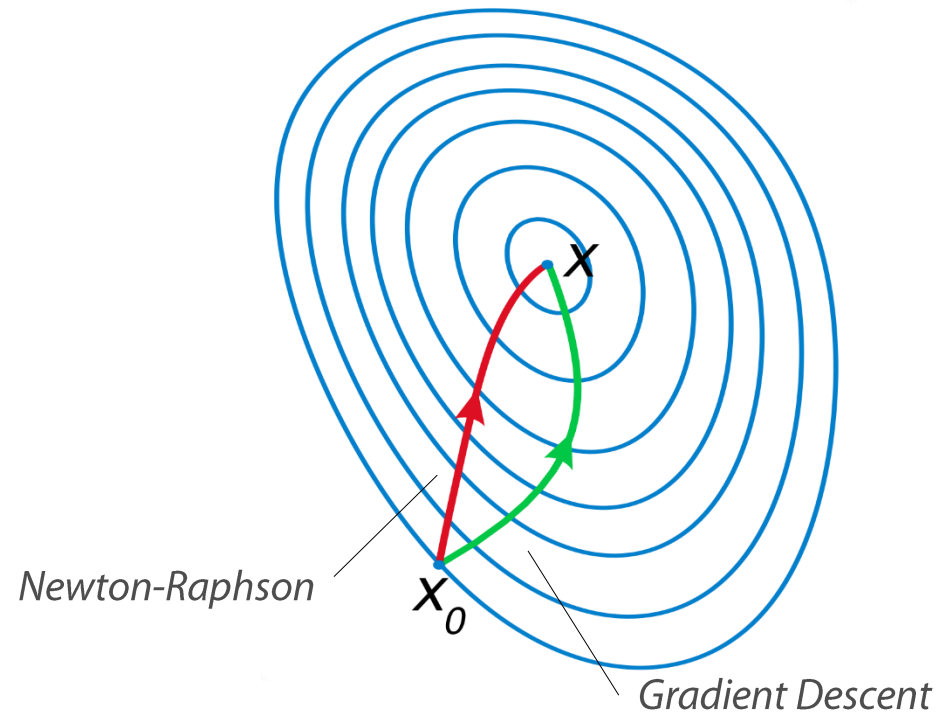


2nd order methods

- **Newton-Raphson's optimization method**

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta \mathbf{H}^{-1} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)}) \quad \mathbf{H} := \frac{\partial}{\partial \boldsymbol{\vartheta}} \left(\frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)}) \right)$$

The (inverse of the) Hessian Matrix considers the curvature



AdaGrad

■ **Newton-Raphson's optimization method**

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta \boldsymbol{H}^{-1} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)}) \quad \boldsymbol{H} := \frac{\partial}{\partial \boldsymbol{\vartheta}} \left(\frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)}) \right)$$

However

- Computing the inverse Hessian matrix is not easy, in general
- It requires $\mathcal{O}(d^3)$ time versus $\mathcal{O}(d)$ of the gradient — d is the number of parameters

AdaGrad

■ Newton-Raphson's optimization method

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta \mathbf{H}^{-1} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)}) \quad \mathbf{H} := \frac{\partial}{\partial \boldsymbol{\vartheta}} \left(\frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)}) \right)$$

However

- Computing the inverse Hessian matrix is not easy, in general
- It requires $\mathcal{O}(d^3)$ time versus $\mathcal{O}(d)$ of the gradient — d is the number of parameters

■ AdaGrad approximation

$$G_i^{(t)} := \sqrt{\sum_{j=1}^t \left(\frac{\partial}{\partial \vartheta_i} L(B, \boldsymbol{\vartheta}^{(j)}) \right)^2} \quad \mathbf{G}^{(t)} := \begin{bmatrix} G_1^{(t)} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & G_d^{(t)} \end{bmatrix}$$

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta (\mathbf{G}^{(t-1)})^{-1} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)})$$

AdaGrad

Gradient Descent

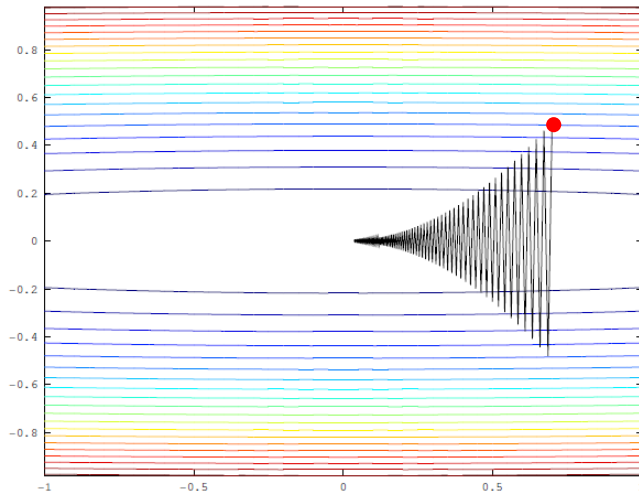
$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)})$$

Newton-Raphson

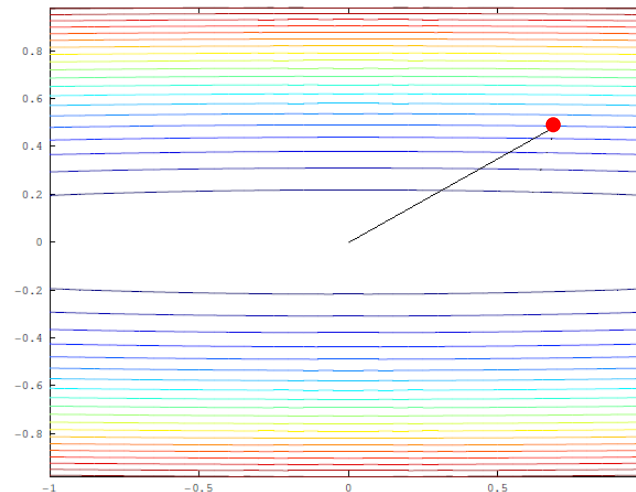
$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta \mathbf{H}^{-1} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)})$$

AdaGrad

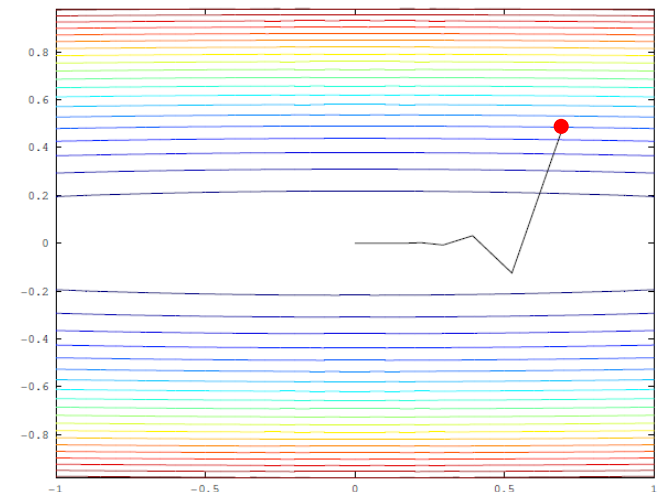
$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta (\mathbf{G}^{(t-1)})^{-1} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)})$$



Gradient Descent



Newton-Raphson



AdaGrad

RMSprop

■ AdaGrad approximation

$$G_i^{(t)} := \sqrt{\sum_{j=1}^t \left(\frac{\partial}{\partial \vartheta_i} L(B, \boldsymbol{\vartheta}^{(j)}) \right)^2}$$

■ RMSprop approximation

The overall sum is replaced by the exponential moving average (EMA)

$$g_i^{(t)} := \frac{\partial}{\partial \vartheta_i} L(B, \boldsymbol{\vartheta}^{(t)})$$

$$\text{EMA}(g_i^2)^{(t)} := \gamma (g_i^{(t)})^2 + (1 - \gamma) \text{EMA}(g_i^2)^{(t-1)}$$

$$G_i^{(t)} := \sqrt{\text{EMA}(g_i^2)^{(t)}}$$

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta (\mathbf{G}^{(t-1)})^{-1} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)})$$

$$\mathbf{G}^{(t)} := \begin{bmatrix} G_1^{(t)} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & G_d^{(t)} \end{bmatrix}$$

AdaDelta

■ RMSprop approximation

$$g_i^{(t)} := \frac{\partial}{\partial \vartheta_i} L(B, \boldsymbol{\vartheta}^{(t)})$$

$$\text{EMA}(g_i^2)^{(t)} := \gamma (g_i^{(t)})^2 + (1 - \gamma) \text{EMA}(g_i^2)^{(t-1)}$$

$$G_i^{(t)} := \sqrt{\text{EMA}(g_i^2)^{(t)}}$$

— Hessian approximation

$$\mathbf{G}^{(t)} := \begin{bmatrix} G_1^{(t)} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & G_d^{(t)} \end{bmatrix}$$

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta (\mathbf{G}^{(t-1)})^{-1} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)})$$

■ AdaDelta approximation

$$D_i^{(t)} := \sqrt{\text{EMA}(\Delta \vartheta_i^2)^{(t)}}$$

— 'momentum' factor

$$\mathbf{D}^{(t)} := \begin{bmatrix} D_1^{(t)} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & D_d^{(t)} \end{bmatrix}$$

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta \mathbf{D}^{(t-1)} (\mathbf{G}^{(t-1)})^{-1} \frac{\partial}{\partial \boldsymbol{\vartheta}} L(B, \boldsymbol{\vartheta}^{(t-1)})$$

Improving optimization

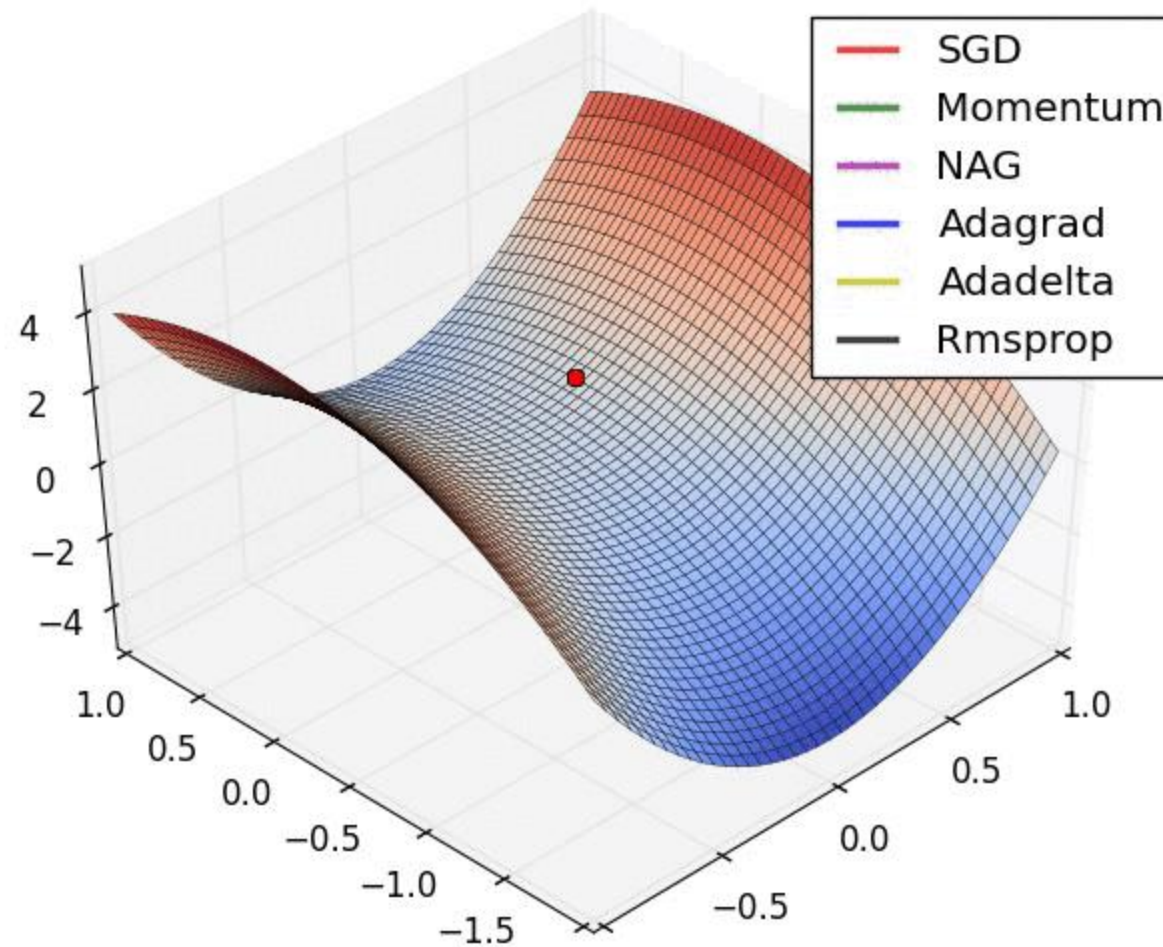


Image from <https://imgur.com/a/Hqolp>

Improving optimization

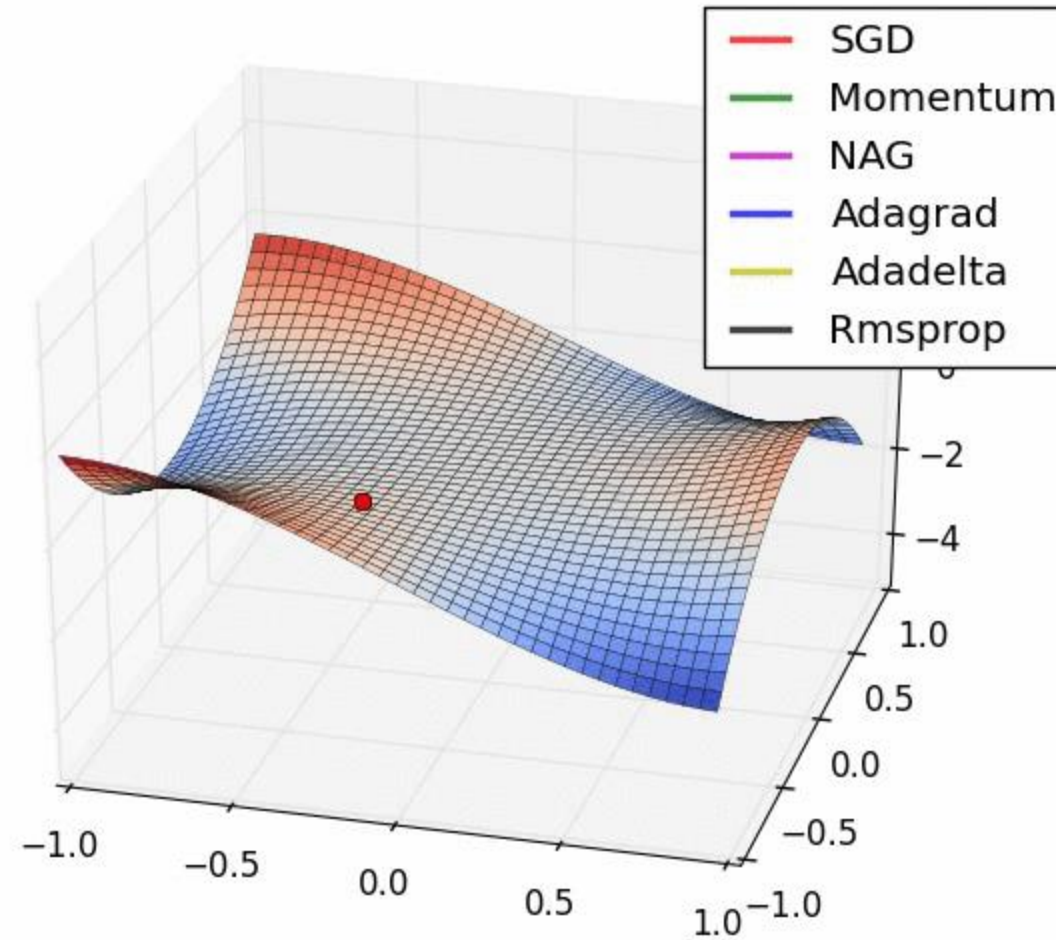


Image from <https://imgur.com/a/Hqolp>

Improving optimization

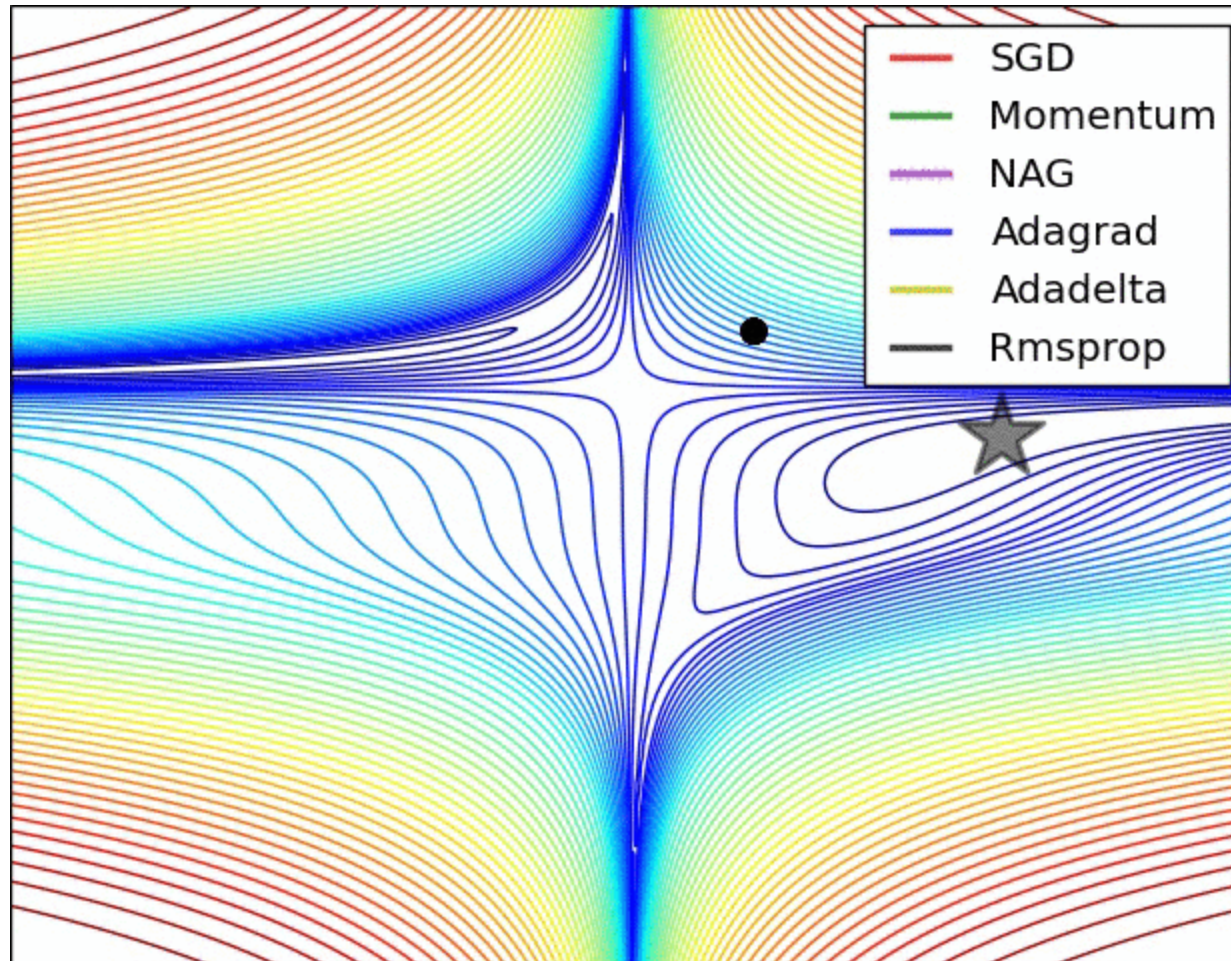


Image from <https://imgur.com/a/Hqolp>

Adam

- **Replace components with their EMAs ...**

$$m_i^{(t)} := \beta_1(g_i^{(t)}) + (1 - \beta_1)m_i^{(t-1)}$$

$$\mathbf{m}^{(t)} := \begin{bmatrix} m_1^{(t)} \\ \vdots \\ m_d^{(t)} \end{bmatrix} \quad \text{EMA of the gradient}$$

$$r_i^{(t)} := \beta_2(g_i^{(t)})^2 + (1 - \beta_2)r_i^{(t-1)}$$

$$\mathbf{r}^{(t)} := \begin{bmatrix} r_1^{(t)} \\ \vdots \\ r_d^{(t)} \end{bmatrix} \quad \text{EMA of the Hessian approximation (vector form)}$$

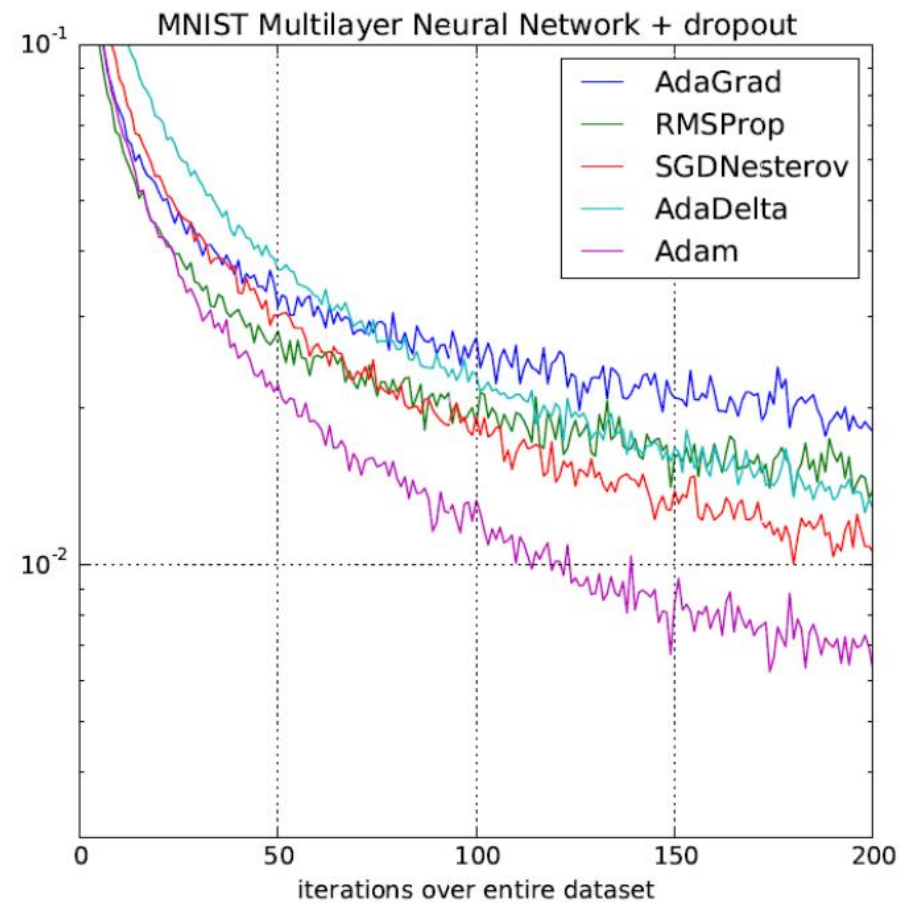
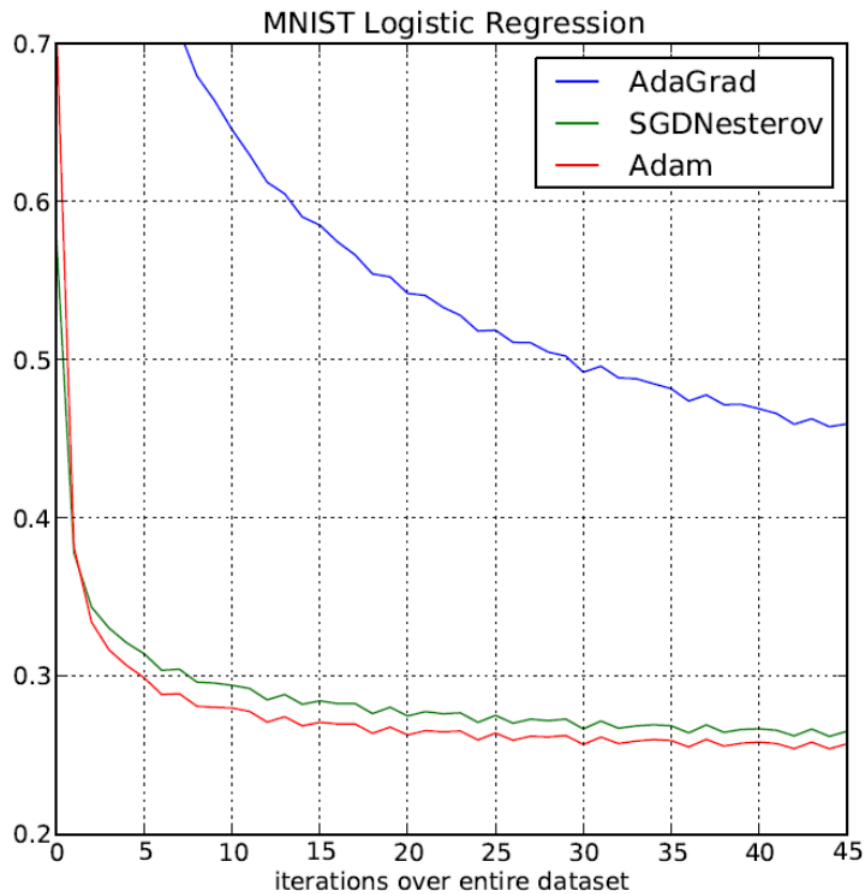
$$\hat{\mathbf{m}}^{(t)} := \frac{\mathbf{m}^{(t)}}{1 - (1 - \beta_1)^t} \quad \text{bias corrections (decay with time)}$$

$$\hat{\mathbf{r}}^{(t)} := \frac{\mathbf{r}^{(t)}}{1 - (1 - \beta_2)^t}$$

$$\boldsymbol{\vartheta}^{(t)} = \boldsymbol{\vartheta}^{(t-1)} - \eta \frac{\hat{\mathbf{m}}^{(t-1)}}{\sqrt{\hat{\mathbf{r}}^{(t-1)}}} \quad \text{(elementwise)}$$

Adam

■ Experimentally



Improving optimization

■ *Messages to take home*

- Improved optimizers adopt a combination of intuition and mathematical modeling
- In particular, some of them are approximators to 2nd order optimization methods
- As such, there is no formal guarantee that they will be effective in all cases

Moral: in general, their effectiveness will depend on the optimization problem and the representation being used

A bag of wonderful tricks

Why ReLU is better (sometimes)

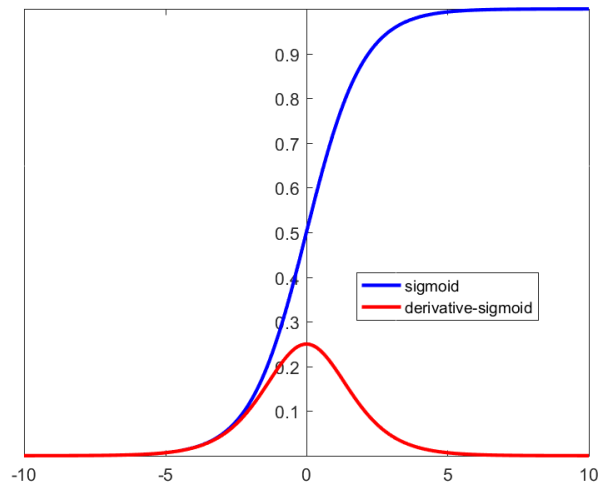
The gradient descent method implies updating the parameters at each step:
making sure that the gradient does not either *vanish* or *explode* is not easy

For instance, in

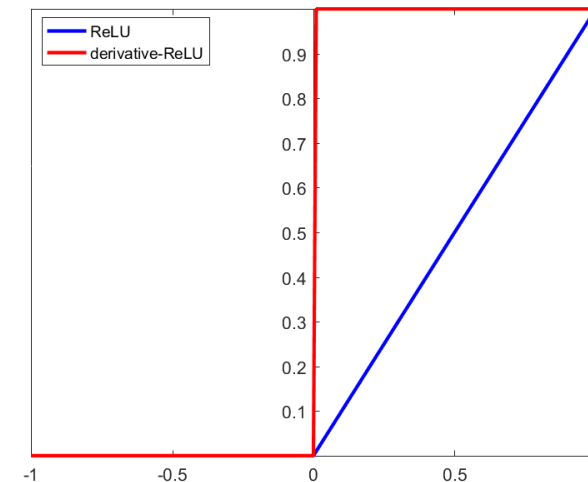
$$\Delta \mathbf{W} = -\eta \frac{\partial L}{\partial \mathbf{W}}(\tilde{y}^{(i)}, y^{(i)})$$

the gradient contains a multiplicative term
which can be $\ll 1.0$

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



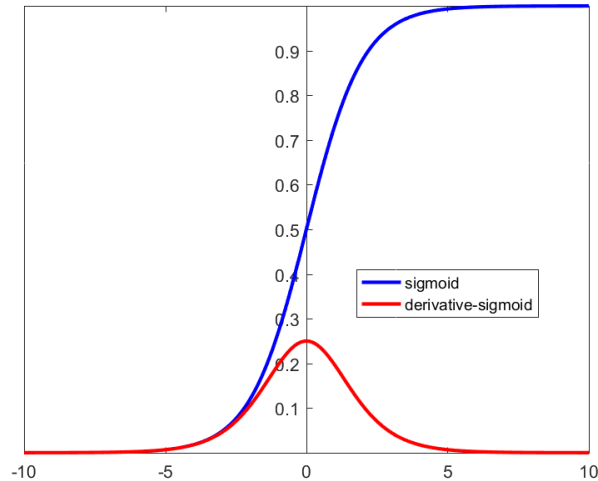
In general,
the derivative of ReLU
does not suffer
from the same problem



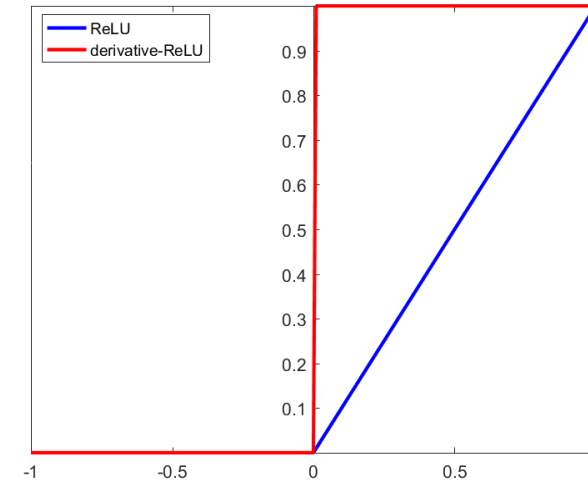
Why ReLU is better (sometimes)

In experimental practice (*sometimes*):

- ReLU alleviates the problem of initial values (i.e. when initial values are too far away and cause sigmoid or tanh to saturate)



*In general,
the derivative of ReLU
does not suffer
from the same problem*



Why ReLU is better (sometimes)

In experimental practice (*sometimes*):

- ReLU alleviates the problem of initial values (i.e. when initial values are too far away and cause sigmoid or tanh to saturate)
- ReLU may accelerate the training process

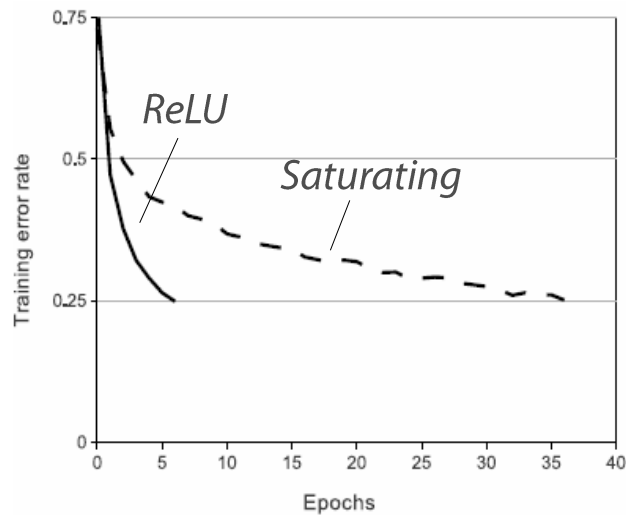
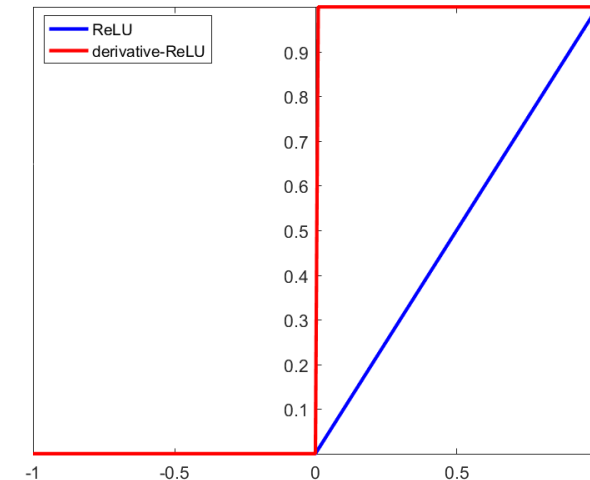


Image from [Krizhevsky, Sutskever & Hinton, 2012]



Xavier (Glorot) Initialization

■ The problem of initialization

Have a look to the demo at: <https://www.deeplearning.ai/ai-notes/initialization/index.html>

See what happens when the initial weights are either too small or too large

■ Objectives

1. *The mean of the activations should be zero*
2. *The variance of the activations should stay the same across layers*

■ Strategy

$$\mathbf{W}^{[l]} \sim \mathcal{N}\left(0, \frac{1}{h^{[l-1]}}\right) \quad \text{or} \quad \mathbf{W}^{[l]} \sim \mathcal{N}\left(0, \frac{2}{h^{[l-1]} + h^{[l]}}\right)$$

$$\mathbf{b}^{[l]} = 0$$

where: $\mathbf{W}^{[l]} \in \mathbb{R}^{h^{[l]} \times h^{[l-1]}}$

Under some simplifying assumptions, this makes all layers have the same variance
(see the link above for a complete mathematical justification)

Input Normalization

■ Intuition

Consider the (very simple) layer

$$h(\mathbf{x}) := g(\mathbf{w}\mathbf{x} + b) = g(w_1x_1 + w_2x_2 + b)$$

and suppose $x_1 \in [1000, 2000]$, $x_2 \in [0.1, 0.2]$ —

- w_1 influences h a lot more than w_2
- training w_2 is challenging and slow

x_1 and x_2 are in completely different scales

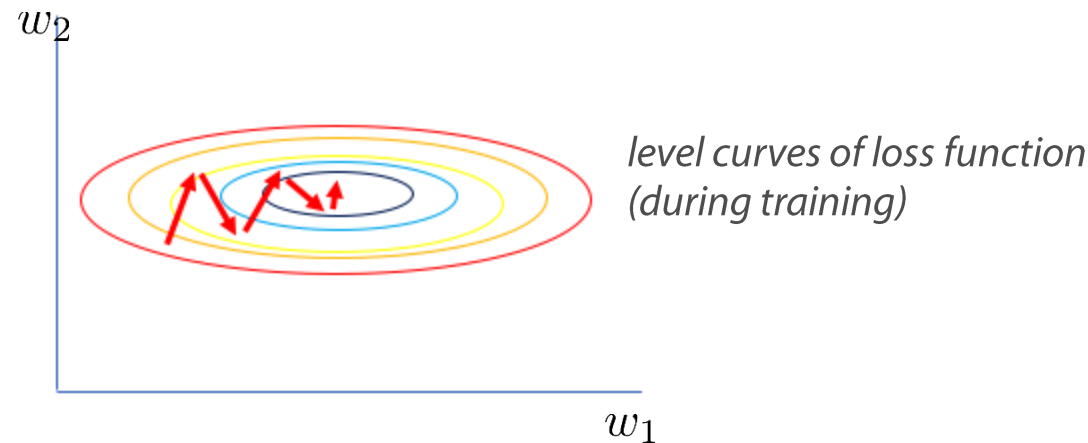
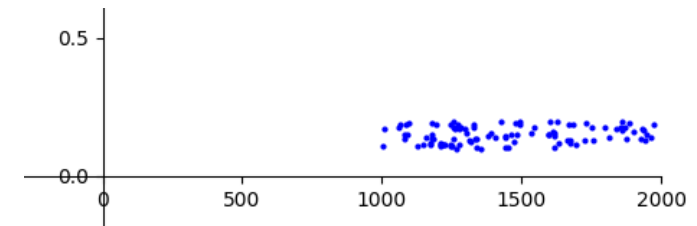


Image from <https://www.jeremyjordan.me/batch-normalization/>

Input Normalization

■ Input normalization

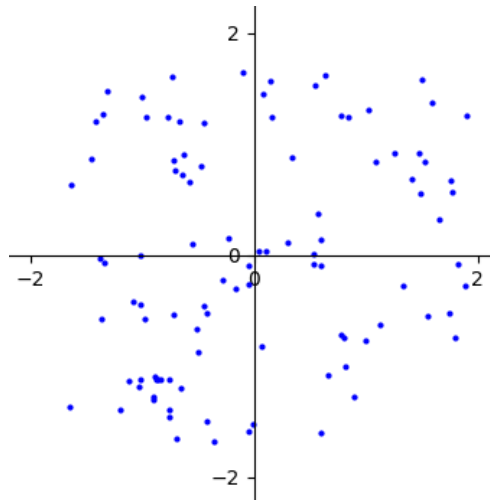
- 1) compute **mean** μ and (component-wise) **variance** σ^2 of inputs over dataset D

$$\mu := \frac{1}{|D|} \sum_{x \in D} x \quad \sigma^2 := (\sigma_1^2, \dots, \sigma_d^2) \quad \text{with } \sigma_i^2 := \frac{1}{|D|} \sum_{x \in D} (x_i - \mu_i)^2$$

- 2) normalize all inputs, component-wise

$$\hat{x} := (\hat{x}_1, \dots, \hat{x}_d), \quad \text{with } \hat{x}_i := \frac{x_i - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$

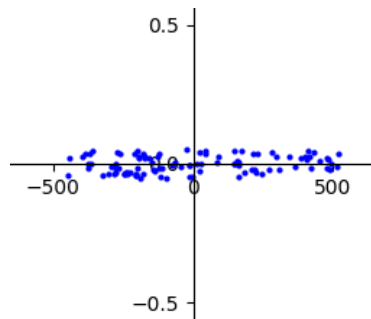
to avoid division by zero



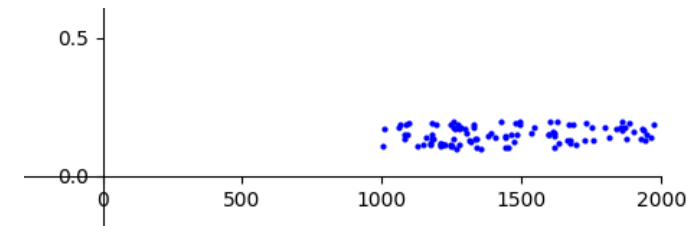
rescale
each component

by

$$\frac{1}{\sqrt{\sigma_i^2 + \epsilon}}$$



shift by μ



Input Normalization

■ Input normalization

1) compute **mean** μ and (component-wise) **variance** σ^2 of inputs over dataset D

$$\mu := \frac{1}{|D|} \sum_{\mathbf{x} \in D} \mathbf{x} \quad \sigma^2 := (\sigma_1^2, \dots, \sigma_d^2) \quad \text{with } \sigma_i^2 := \frac{1}{|D|} \sum_{\mathbf{x} \in D} (x_i - \mu_i)^2$$

2) normalize all inputs, component-wise

$$\hat{\mathbf{x}} := (\hat{x}_1, \dots, \hat{x}_d), \quad \text{with } \hat{x}_i := \frac{x_i - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$

3) apply $h(\hat{\mathbf{x}}) := g(\mathbf{w}\hat{\mathbf{x}} + b) = g(w_1\hat{x}_1 + w_2\hat{x}_2 + b)$

Input Normalization

■ Input normalization

1) compute **mean** μ and (component-wise) **variance** σ^2 of inputs over dataset D

$$\mu := \frac{1}{|D|} \sum_{x \in D} x \quad \sigma^2 := (\sigma_1^2, \dots, \sigma_d^2) \quad \text{with } \sigma_i^2 := \frac{1}{|D|} \sum_{x \in D} (x_i - \mu_i)^2$$

2) normalize all inputs, component-wise

$$\hat{x} := (\hat{x}_1, \dots, \hat{x}_d), \quad \text{with } \hat{x}_i := \frac{x_i - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$

3) apply $h(\hat{x}) := g(\mathbf{w}\hat{x} + b) = g(w_1\hat{x}_1 + w_2\hat{x}_2 + b)$

- training becomes faster and more stable
(also allowing higher learning rates)

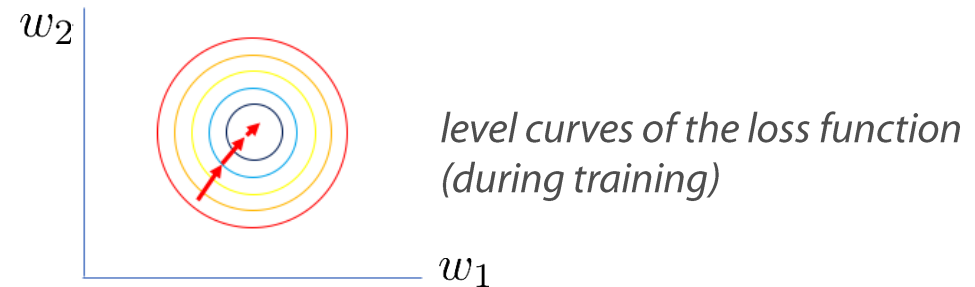


Image from <https://www.jeremyjordan.me/batch-normalization/>

Batch Normalization

- **Normalizing in between layers**

In a DNN

$$\tilde{y} = h^{[n]}(h^{[n-1]}(\dots(h^{[2]}(h^{[1]}(x))\dots))$$

each layer $h^{[i]}$ has an input of its own, which should be *normalized*

How?

Batch Normalization

■ Normalizing in between layers

In a DNN

$$\tilde{y} = h^{[n]}(h^{[n-1]}(\dots(h^{[2]}(h^{[1]}(x))\dots))$$

each layer $h^{[i]}$ has an input of its own, which should be *normalized*

Normalizing in between layers during training would require:

- pre-computing the input to each layer, for *each data item* in D
- applying normalization before proceeding further upwards
- doing it again after *each* updating the DNN parameters

Moral: *it's impossible*

Batch Normalization

- For each mini-batch:

$$B = \left\{ \mathbf{x}^{(i)} \right\}_{i=1}^m$$

(all operations are performed element-wise)

$$\text{BN}_{\beta, \gamma}(\mathbf{x}^{(i)}) := \gamma \hat{\mathbf{x}}^{(i)} + \beta$$

trainable parameters

$$\hat{\mathbf{x}}^{(i)} = \frac{\mathbf{x}^{(i)} - \boldsymbol{\mu}_B}{\sqrt{\boldsymbol{\sigma}_B^2 + \epsilon}}$$

avoid division by zero

$$\boldsymbol{\sigma}_B^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)} - \boldsymbol{\mu}_B)$$

$$\boldsymbol{\mu}_B = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)}$$

Batch Normalization

■ Training

- at step t : $\mu_{B^{(t)}}$ and $\sigma_{B^{(t)}}^2$ are computed over the current mini-batch $B^{(t)}$
- parameters γ and β (for each BN-layer) are trained *in the same way as the other parameters in the DNN*
- *exponential moving averages* of mean and variance of the mini-batches $B^{(t)}$ are collected

$$\begin{aligned} \text{MA}(\mu)^{(t)} &:= \delta \cdot \mu_{B^{(t)}} + (1 - \delta) \cdot \text{MA}(\mu)^{(t-1)}, & \text{MA}(\mu)^{(1)} &:= \mu_{B^{(1)}} \\ \text{MA}(\sigma^2)^{(t)} &:= \delta \cdot \sigma_{B^{(t)}}^2 + (1 - \delta) \cdot \text{MA}(\sigma^2)^{(t-1)}, & \text{MA}(\sigma^2)^{(1)} &:= \sigma_{B^{(1)}}^2 \end{aligned}$$

■ Inference

Inference is typically performed for fewer inputs, possibly just one ...

Batch Normalization

■ Training

- at step t : $\mu_{B^{(t)}}$ and $\sigma_{B^{(t)}}^2$ are computed over the current mini-batch $B^{(t)}$
- parameters γ and β (for each BN-layer) are trained *in the same way as the other parameters in the DNN*
- *exponential moving averages* of mean and variance of the mini-batches $B^{(t)}$ are collected

$$\begin{aligned} \text{MA}(\mu)^{(t)} &:= \delta \cdot \mu_{B^{(t)}} + (1 - \delta) \cdot \text{MA}(\mu)^{(t-1)}, & \text{MA}(\mu)^{(1)} &:= \mu_{B^{(1)}} \\ \text{MA}(\sigma^2)^{(t)} &:= \delta \cdot \sigma_{B^{(t)}}^2 + (1 - \delta) \cdot \text{MA}(\sigma^2)^{(t-1)}, & \text{MA}(\sigma^2)^{(1)} &:= \sigma_{B^{(1)}}^2 \end{aligned}$$

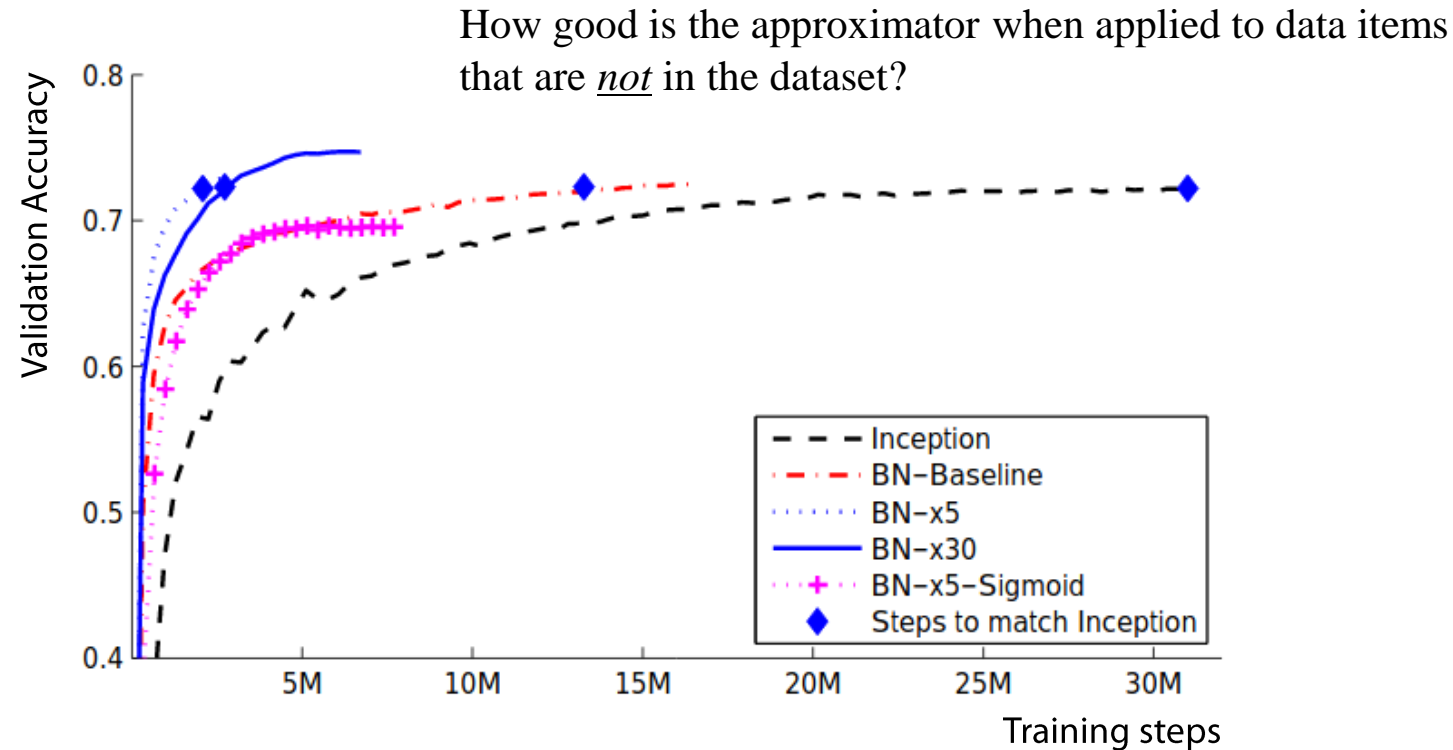
■ Inference

Normalize using the moving averages collected during training

- $\mu := \text{MA}(\mu)^{(T)}$
 - $\sigma^2 := \text{MA}(\sigma^2)^{(T)}$
- as collected during the training process*

Batch Normalization

■ Does it work?



- Batch normalization acts as a *reparametrization* of the optimization process that
 1. makes the loss function smoother
 2. allows higher learning rates
 3. reduces chances to getting stuck into local minima

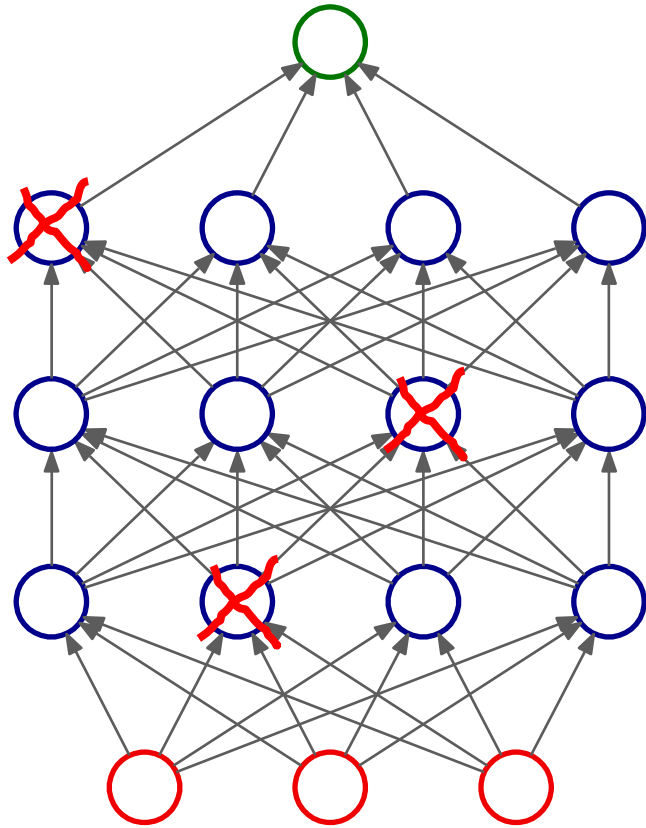
Image from [Ioffe and Szegedy 2015]

Fighting Overfitting

Dropout

- **Knocking-out at random**

For each mini-batch, a small percentage of 'units' is de-activated

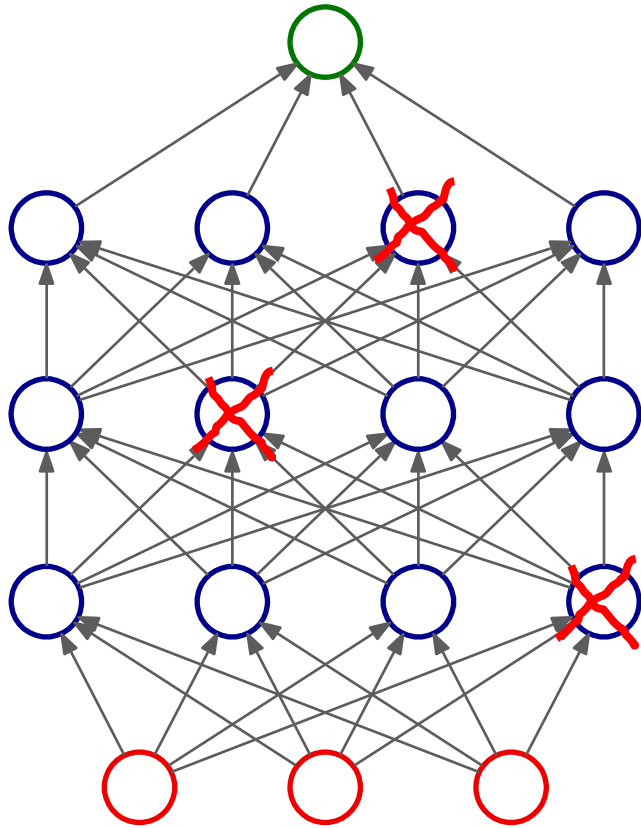


Training: mini-batch 1

Dropout

- **Knocking-out at random**

For each mini-batch, a small percentage of 'units' is de-activated

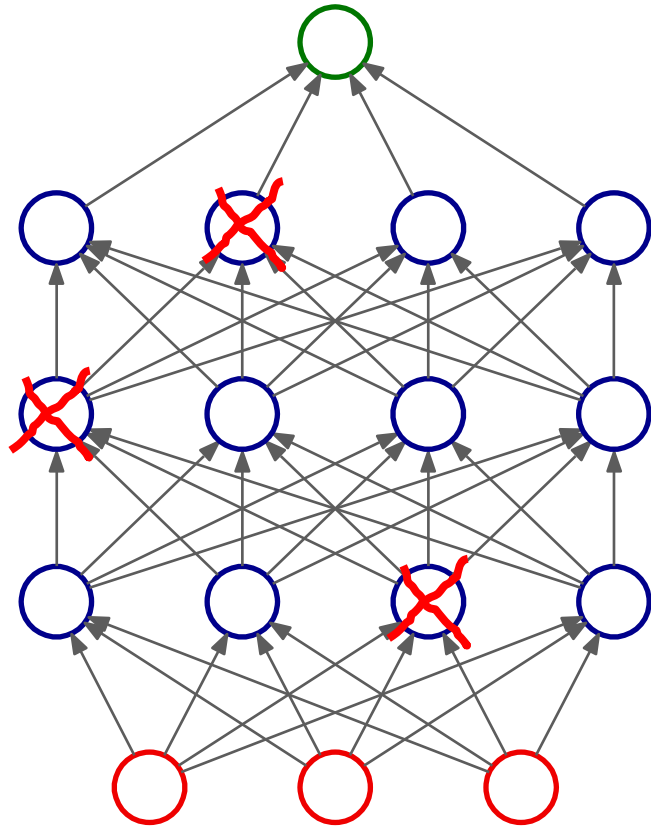


Training: mini-batch 2

Dropout

- **Knocking-out at random**

For each mini-batch, a small percentage of 'units' is de-activated

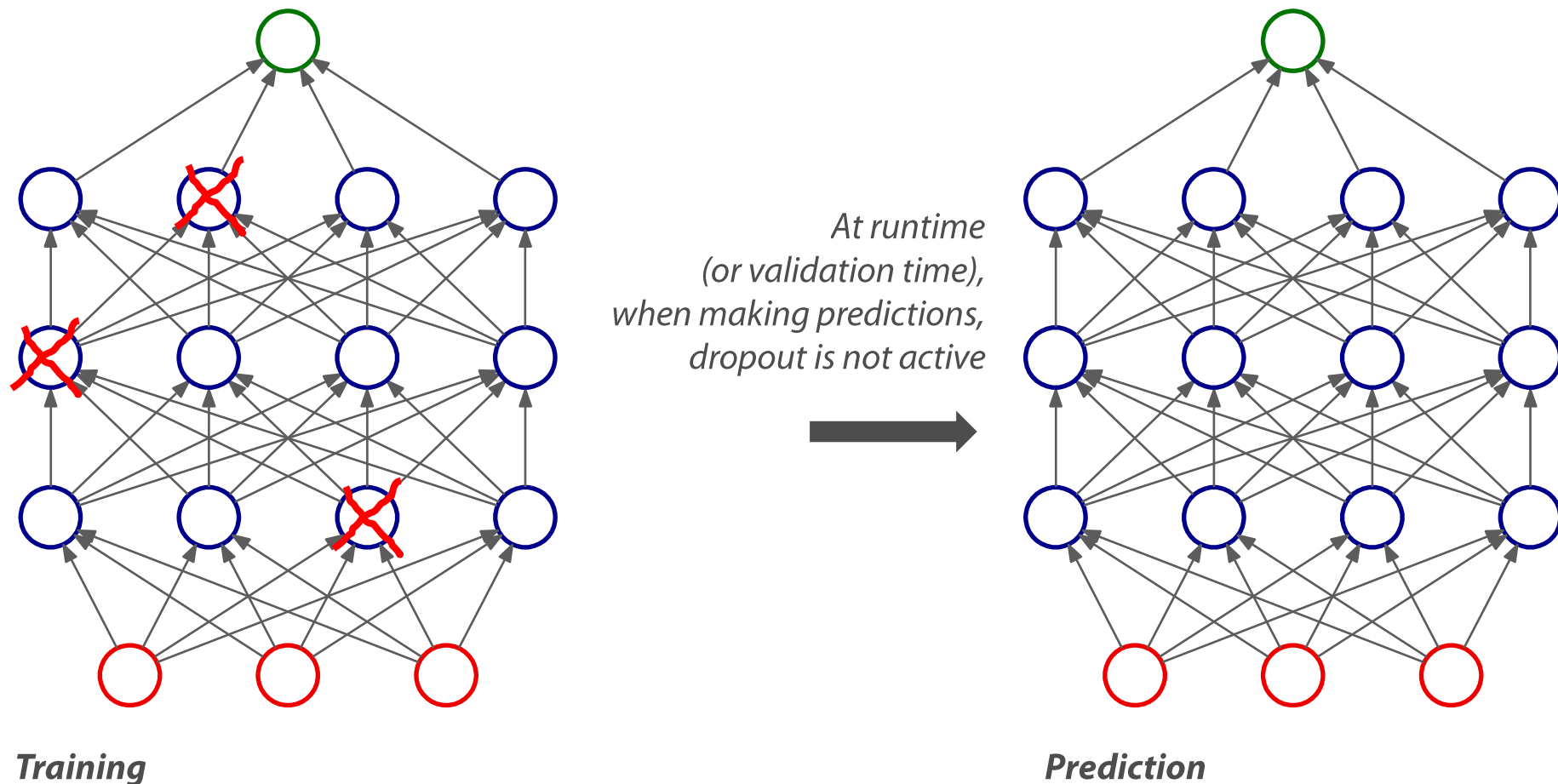


Training: mini-batch 3

Dropout

- **Knocking-out at random**

For each mini-batch, a small percentage of 'units' is de-activated



Contrasting Overfitting

■ Applying Dropout

In a typical experiment

- initially, the performance on D_{val} improves slowly
- then it becomes better and more resilient to *overfitting* (to be explained next)

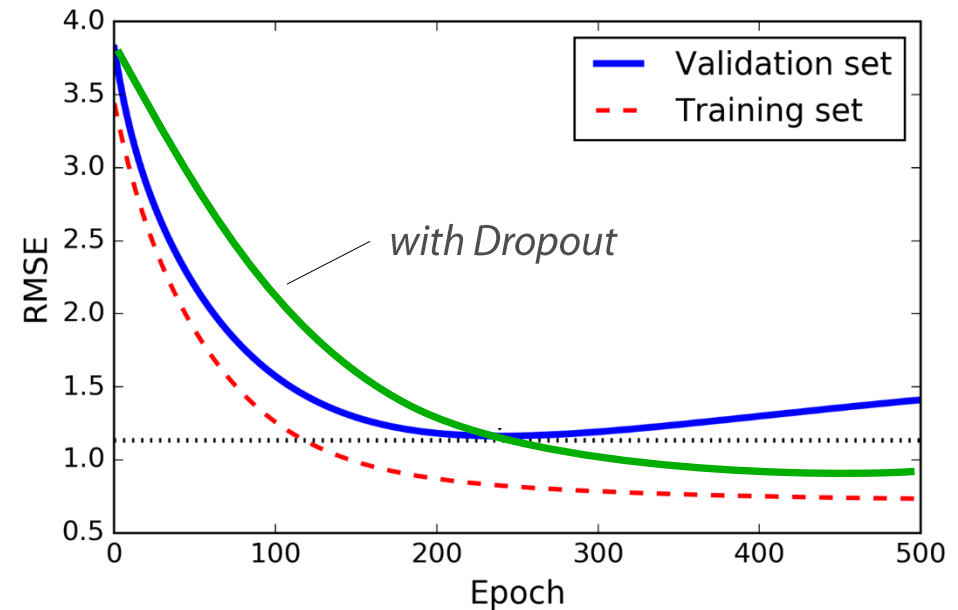


Image from <https://www.safaribooksonline.com/library/view/hands-on-machine-learning/9781491962282/ch04.html>