

Optimization

Florin Gogianu,
Bitdefender ML Team

Deep Learning Course Structure



Bibliography

Textbooks:

- Deisenroth, Mathematics for Machine Learning, [Ch. 5 - Vector Calculus]
- Goodfellow, Deep Learning
- Prince, Understanding Deep Learning

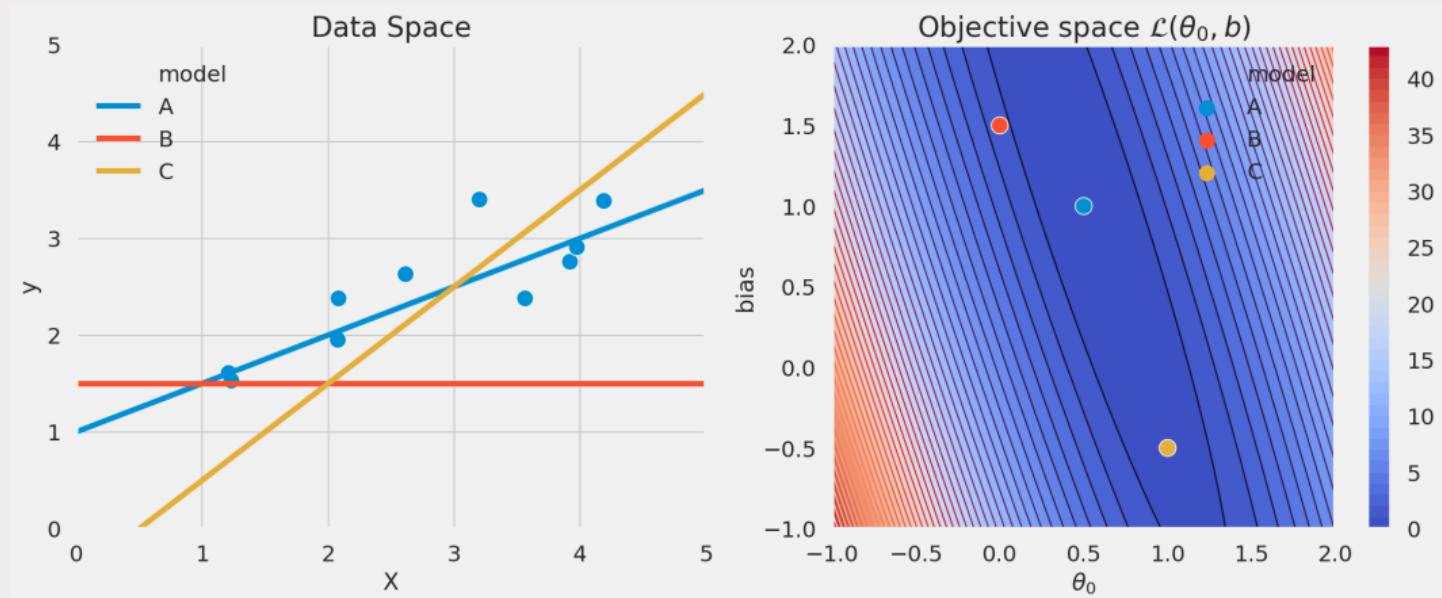
Courses:

- Stanford, CS231n
- NYU Deep Learning

Recap++

Recap. Linear Models

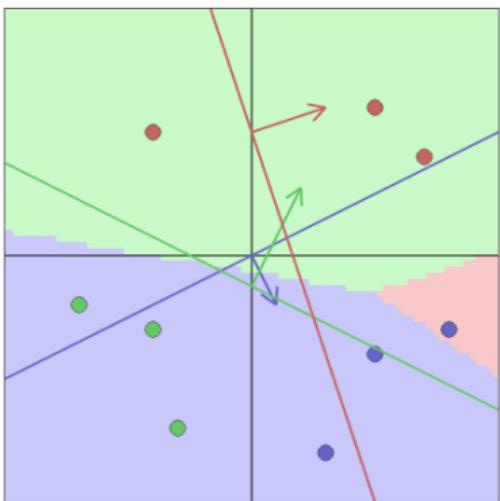
$$\mathbf{y} = \mathbf{Wx} + \mathbf{b}$$



B

Recap. Linear Models

blue line shows the set of points (x_0, x_1) that give score of zero. The blue arrow draws the vector $(W_{0,0}, W_{0,1})$, which shows the direction of score increase and its length is proportional to how steep the increase is.
Note: you can drag the datapoints.



parameters.

$W[0,0]$	$W[0,1]$	$b[0]$
1.00 -0.38	2.00 0.07	0.00 0.11
$W[1,0]$	$W[1,1]$	$b[1]$
2.00 0.51	-4.00 -0.58	0.50 -0.11
$W[2,0]$	$W[2,1]$	$b[2]$
3.00 0.17	-1.00 0.36	-0.50 0.00

Step size: 0.10000

Single parameter update

Start repeated update

x[0]	x[1]	y	s[0]	s[1]	s[2]	L
0.50	0.40	0	1.30	-0.10	0.60	0.30
0.80	0.30	0	1.40	0.90	1.60	1.70
0.30	0.80	0	1.90	-2.10	-0.40	0.00
-0.40	0.30	1	0.20	-1.50	-2.00	3.20
-0.30	0.70	1	1.10	-2.90	-2.10	6.80
-0.70	0.20	1	-0.30	-1.70	-2.80	2.40
0.70	-0.40	2	-0.10	3.50	2.00	2.50
0.50	-0.60	2	-0.70	3.90	1.60	3.30
-0.40	-0.50	2	-1.40	1.70	-1.20	4.70
mean:						2.77

Total data loss: 2.77
Regularization loss: 3.50
Total loss: 6.27

L2 Regularization strength: 0.10000

Multiclass SVM loss formulation:

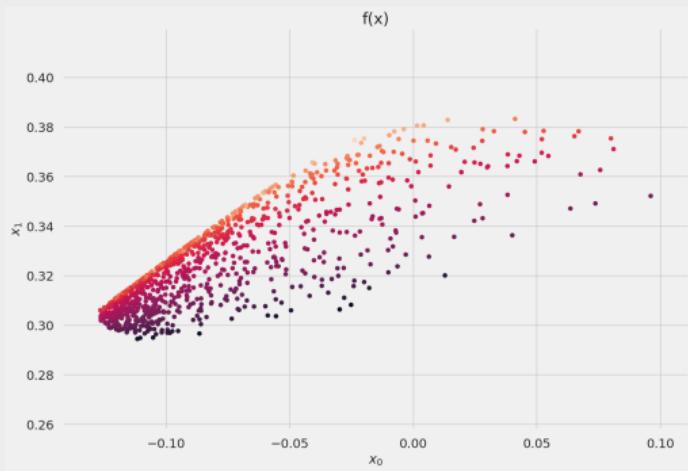
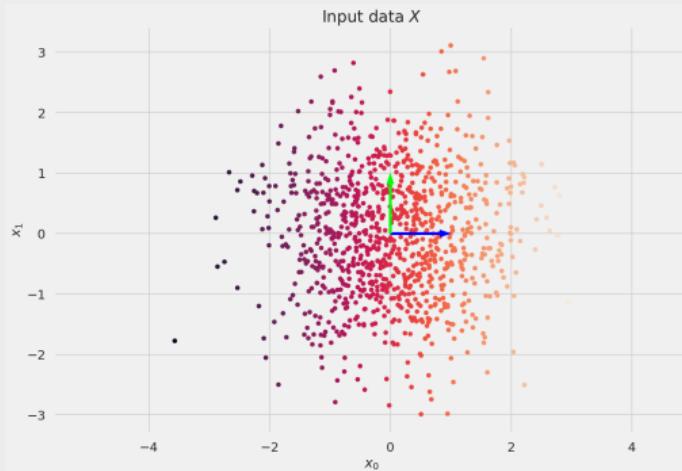
© Weston Watkins 1999

B

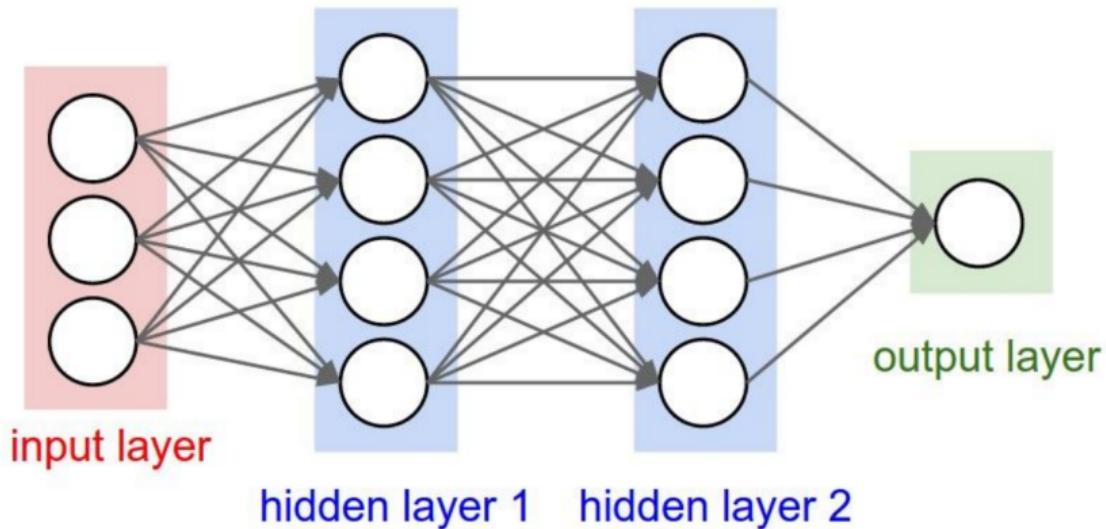
Figure: Linear classification visualization

Recap. NN computations

- $\mathbf{y} = \mathbf{Wx}$ computes: rotation, scale reflection transformations
- $\sigma(\mathbf{Wx})$ adds non-linear transformations

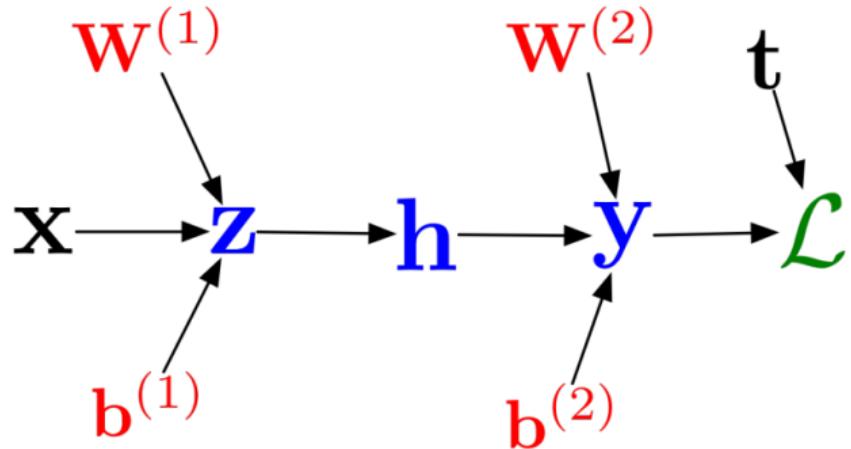


Recap. Fully Connected NN



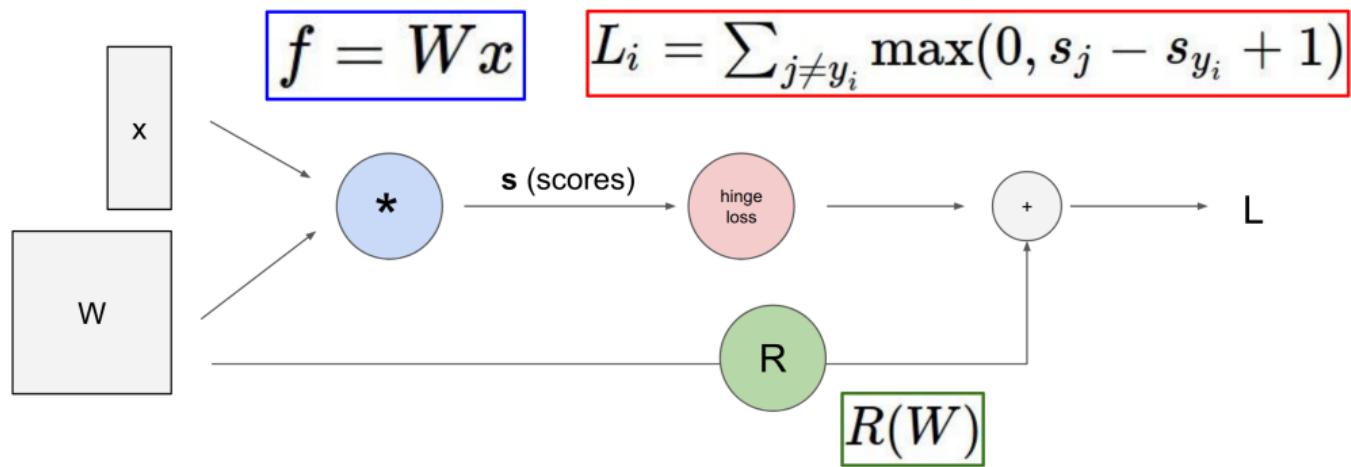
- Linear case: $f(\mathbf{x}) = \mathbf{Wx}$
- Neural net: $f(\mathbf{x}) = \mathbf{W}_3\sigma(\mathbf{W}_2\sigma(\mathbf{W}_1\mathbf{x}))$

Recap. Fully Connected NN



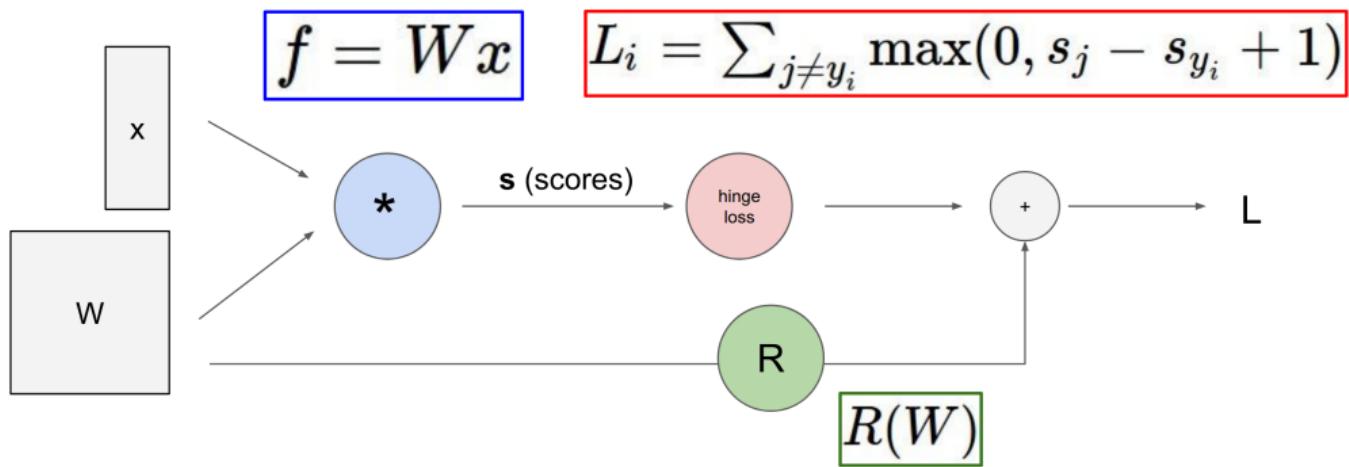
- Neural net: $f(\mathbf{x}) = \mathbf{y} = \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2$
- Loss function: $L = \frac{\mathbf{y} \cdot \mathbf{t}}{\|\mathbf{y}\| \|\mathbf{t}\|}$

Computation graphs + automatic differentiation



¹Grosse, Ba, CS421

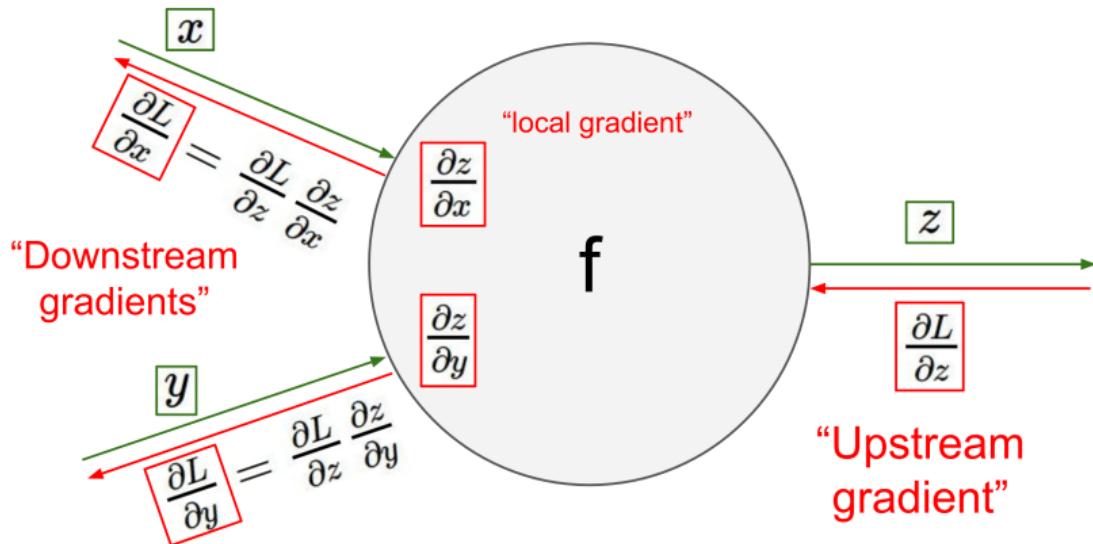
Computation graphs + automatic differentiation



Automatic differentiation: takes a program which computes a value, and automatically constructs a procedure for computing derivatives of that value¹.

¹Grosse, Ba, CS421

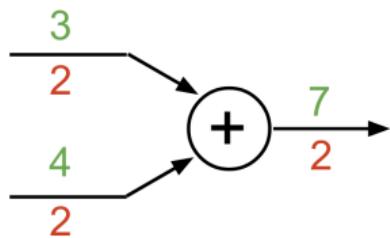
Backprop



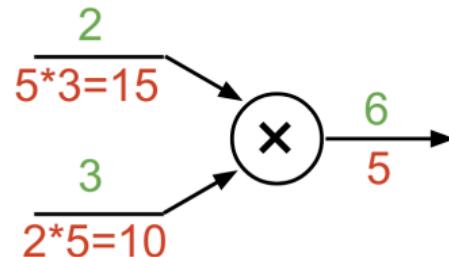
¹from CS231n, lecture 4, 2019

Recap. Backprop patterns

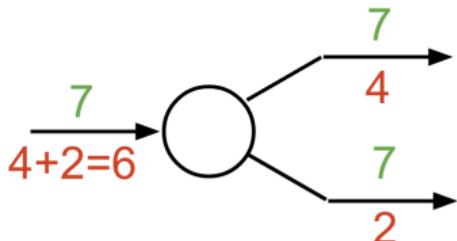
add gate: gradient distributor



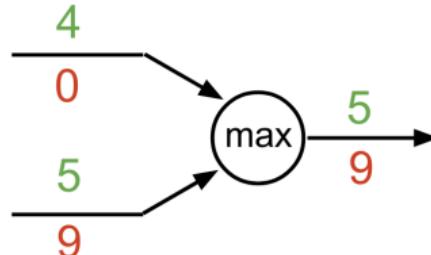
mul gate: “swap multiplier”



copy gate: gradient adder



max gate: gradient router

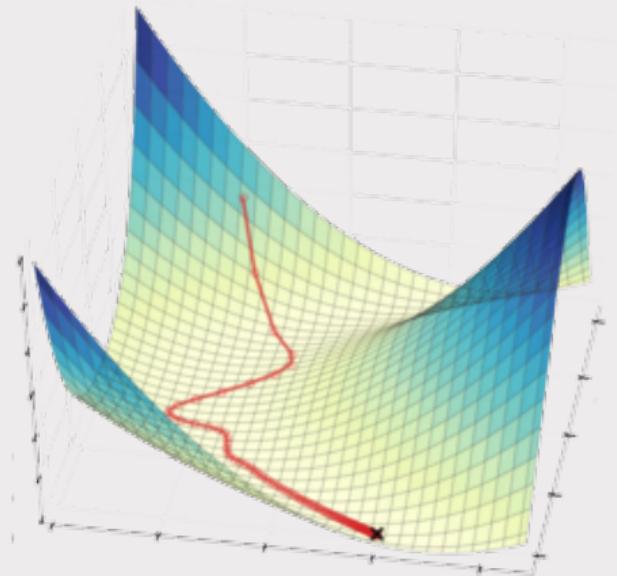


Optimization concepts

Gradient Descent

1. Measure the error of our model using $L(D, \theta)$,
2. Compute the gradients w.r.t. θ
3. Find incremental solutions that better explain the data using:

$$\theta_{j+1} \leftarrow \theta_j - \eta \nabla_{\theta_j} L$$



B

Gradient Descent

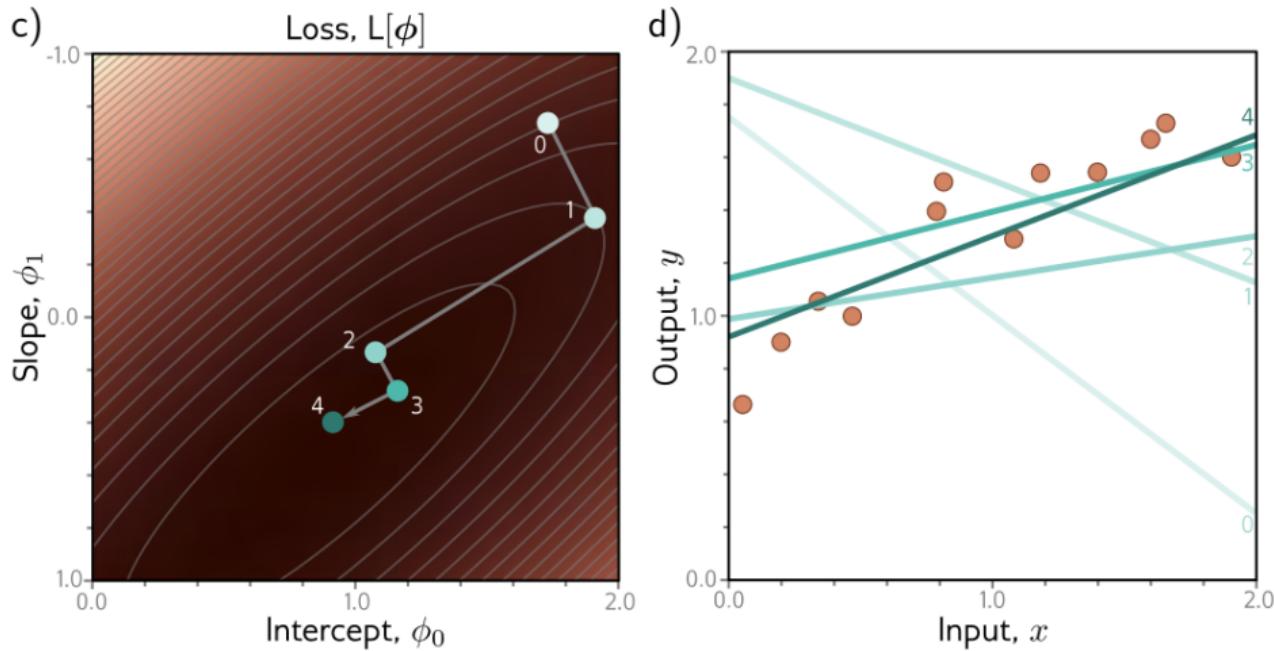


Figure: Gradient descent for the linear model.

B

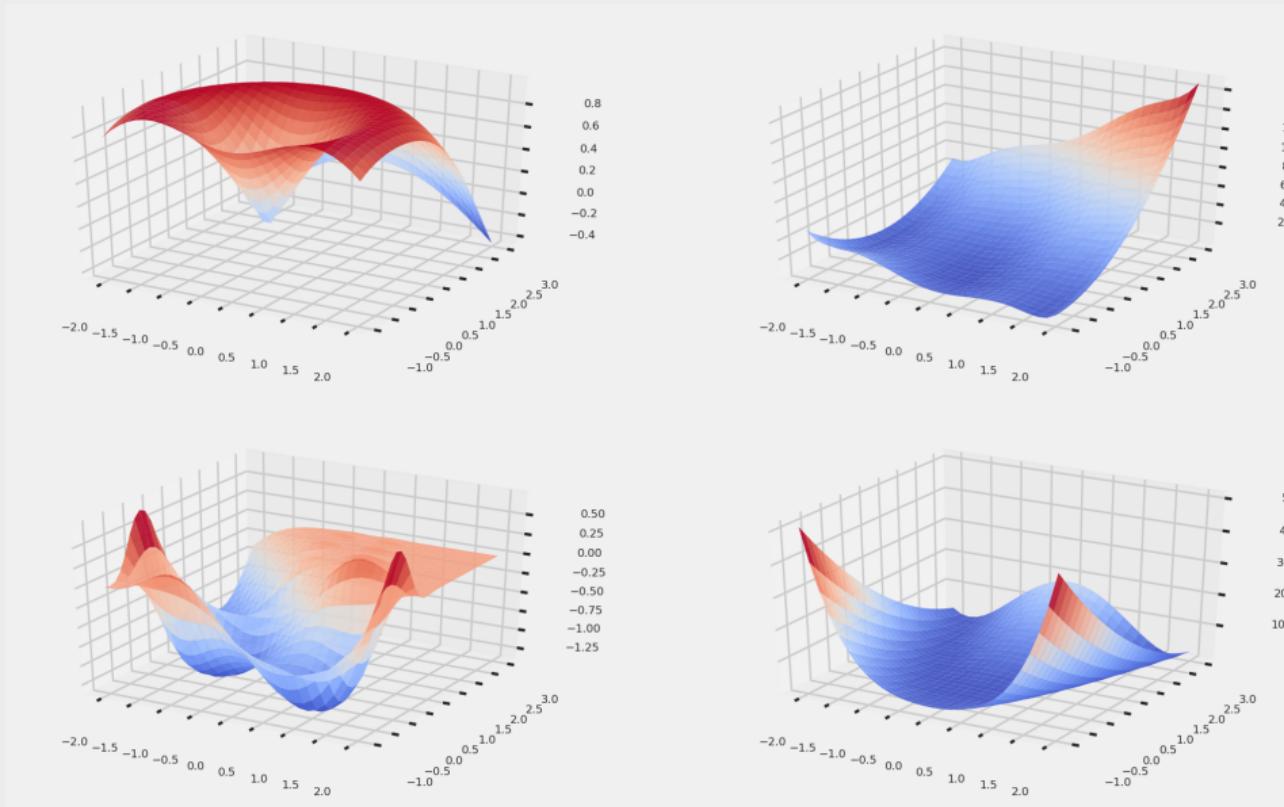
Stochastic Gradient Descent

Consider instead objective functions that are the sum of the losses – $L(\boldsymbol{\theta}) = \sum_{n=1}^N L_n(\boldsymbol{\theta})$, resulting in the update:

$$\boldsymbol{\theta}_{j+1} \leftarrow \boldsymbol{\theta}_j - \eta \sum_{n=1}^N \nabla_{\boldsymbol{\theta}_j} L_n(\boldsymbol{\theta})$$

- evaluating the sum of gradients can be expensive
- instead we can randomly choose a subset of L_n
- for gradient descent to converge we only require the samples to be unbiased

Optimization landscapes



B

Taylor Approximation

- Arbitrary non-linear functions are hard to globally analyze.
- Ideally, we would like to work with simple functions: [polynomials](#).
- [Solution](#): get a polynomial that approximates the function in a *neighbourhood* well enough

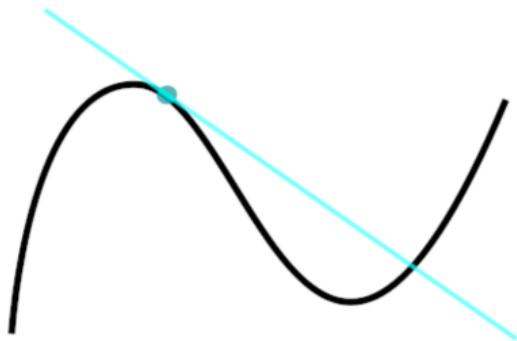
$$L(\boldsymbol{\theta}) \approx L(\boldsymbol{\theta}_0) + \nabla L(\boldsymbol{\theta}_0)^T (\boldsymbol{\theta} - \boldsymbol{\theta}_0) + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \mathbf{H}(\boldsymbol{\theta}_0) (\boldsymbol{\theta} - \boldsymbol{\theta}_0)$$

Taylor Approximation

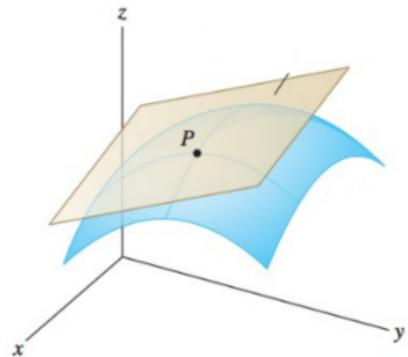


Figure: $\sin(x)$ and its Taylor approximations, polynomials of degree 1, 3, 5, 7, 9, 11, 13
By [IkamusumeFan](#), CC BY-SA 3.0.

First-order information



Tangent line ¹



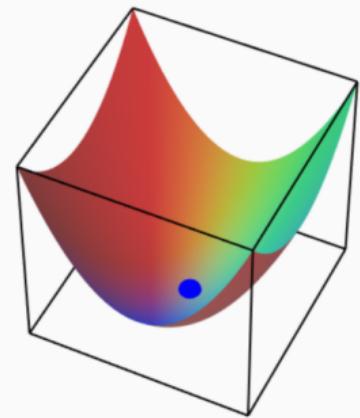
Tangent space ²

Growth of a function when a) changing θ and b) moving in the direction of $\theta - \theta_0$.

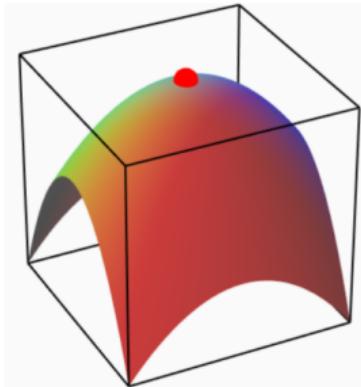
$$L(\theta) \approx L(\theta_0) + \nabla L(\theta_0)^T (\theta - \theta_0)$$

B

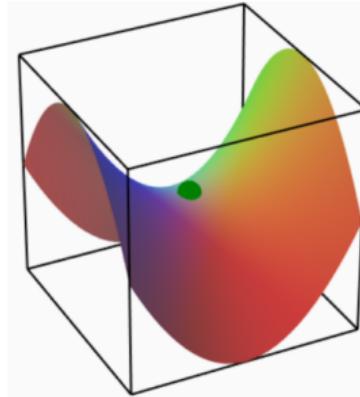
Second-order information



(a) Minima



(b) Maxima

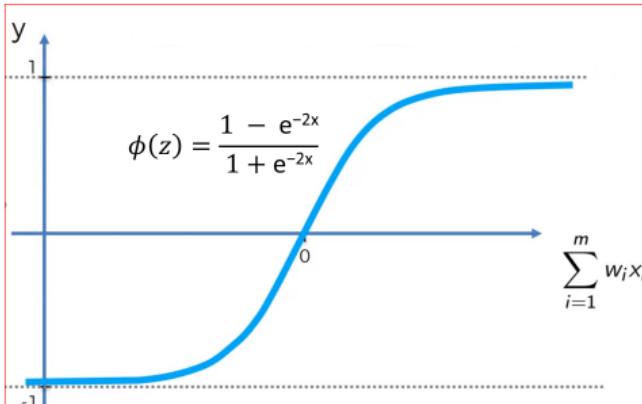
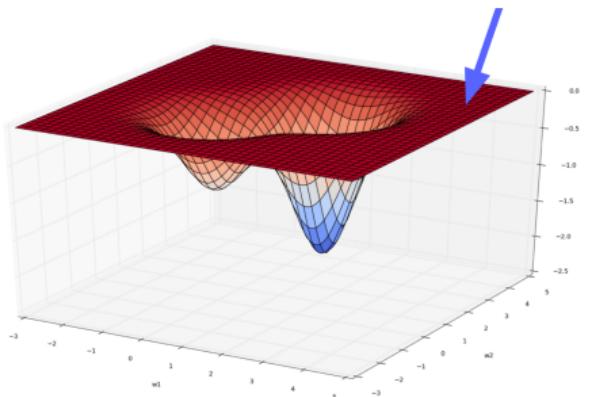


(c) Saddle

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

How does moving in a direction affects the rate of change.

Pathologies. Plateaux

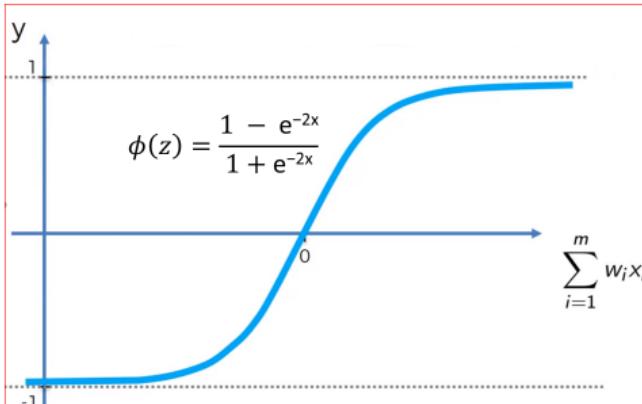
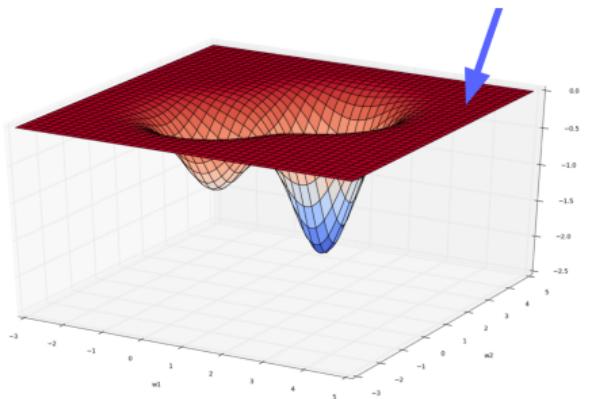


In cases where a neuron is [saturated](#):

$$\delta_\theta = \delta_z \mathbf{x}$$

$$\delta_z = \delta_h \frac{\partial h}{\partial z} = \delta_h \phi'(z)$$

Pathologies. Plateaux



In cases where a neuron is [saturated](#):

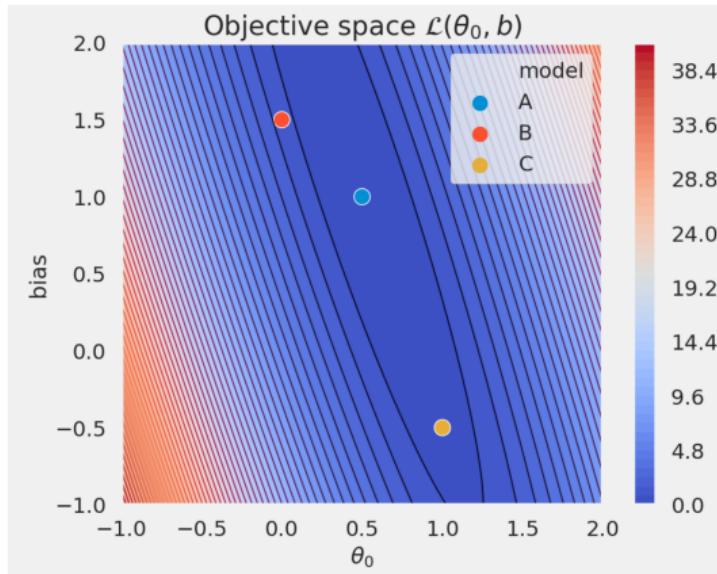
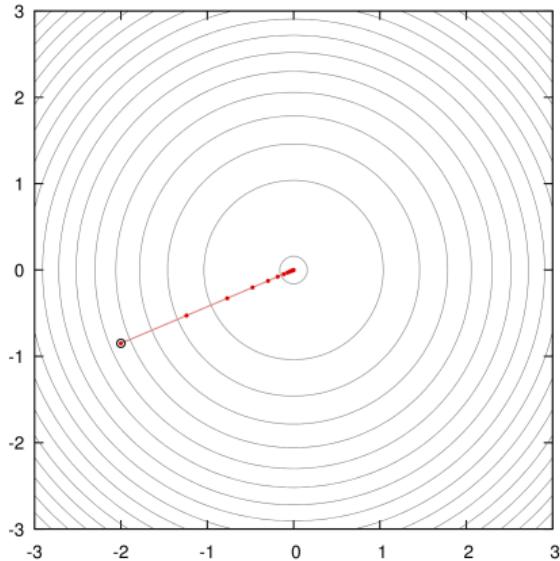
$$\delta_\theta = \delta_z \mathbf{x}$$

$$\delta_z = \delta_h \frac{\partial h}{\partial z} = \delta_h \phi'(z)$$

Use [ReLU](#) activations and good [initialization](#) schemes.

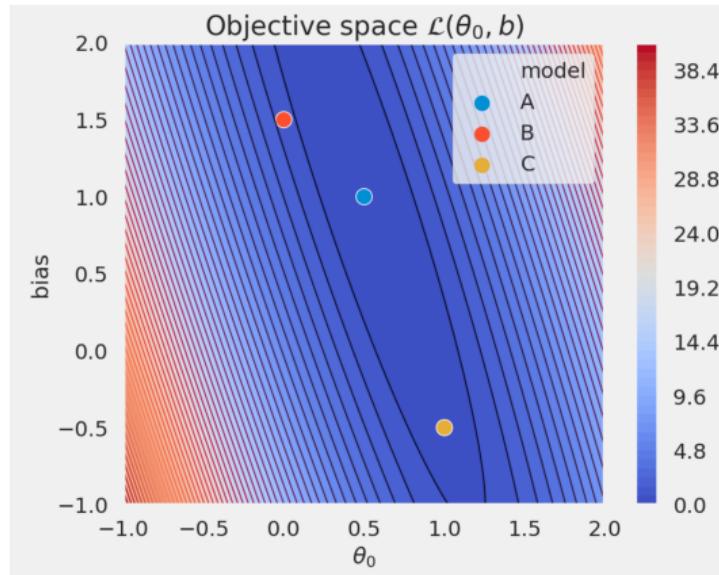
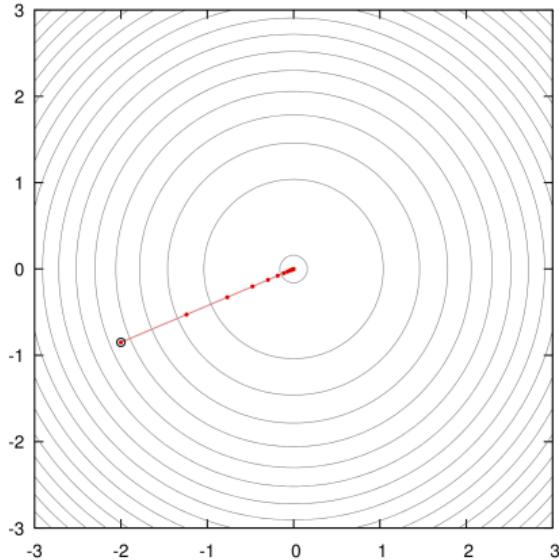
B

Pathologies. High curvature



Remember: $\theta_{j+1} \leftarrow \theta_j - \alpha \nabla_{\theta_j} L$. Badly conditioned curvatures arise even in simple cases, when the inputs have slightly different scales or are off-centered.

Pathologies. High curvature



Remember: $\theta_{j+1} \leftarrow \theta_j - \alpha \nabla_{\theta_j} L$. Badly conditioned curvatures arise even in simple cases, when the inputs have slightly different scales or are off-centered.

Solutions: i) standardize your input data, ii) normalize every pre-activation.

Demo time!

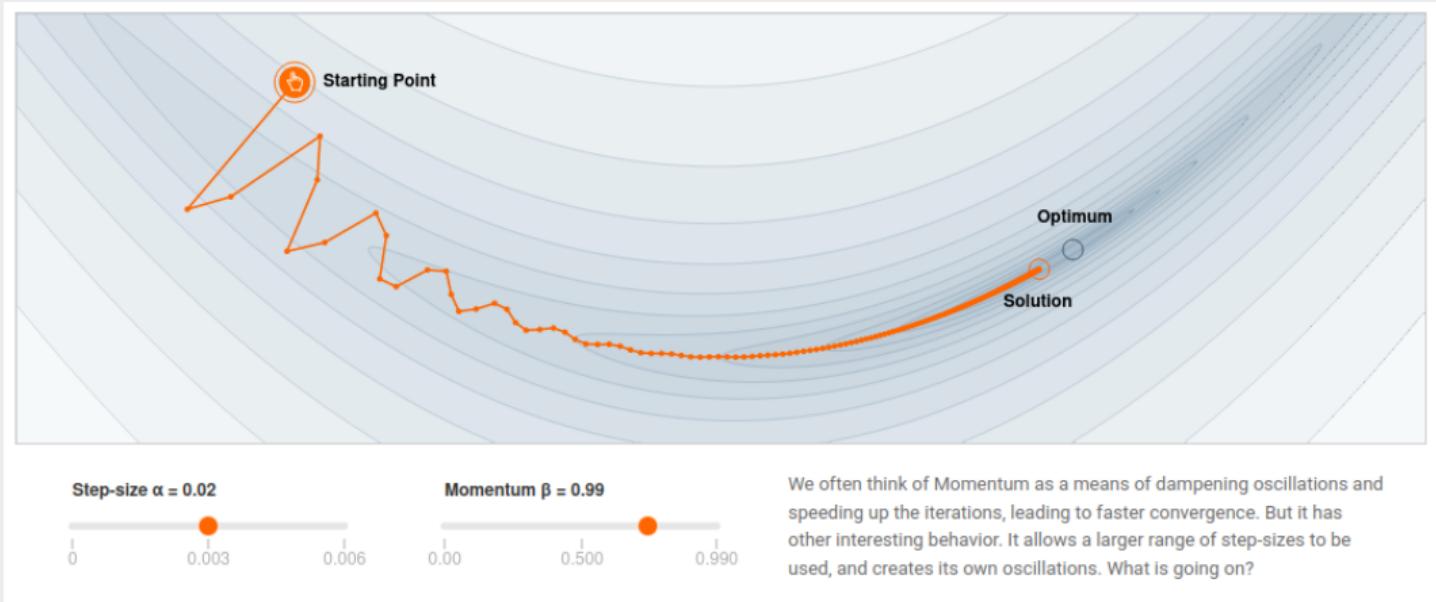


Figure: Example of high curvature inducing optimization problems.

Solution: Normalize pre-activations

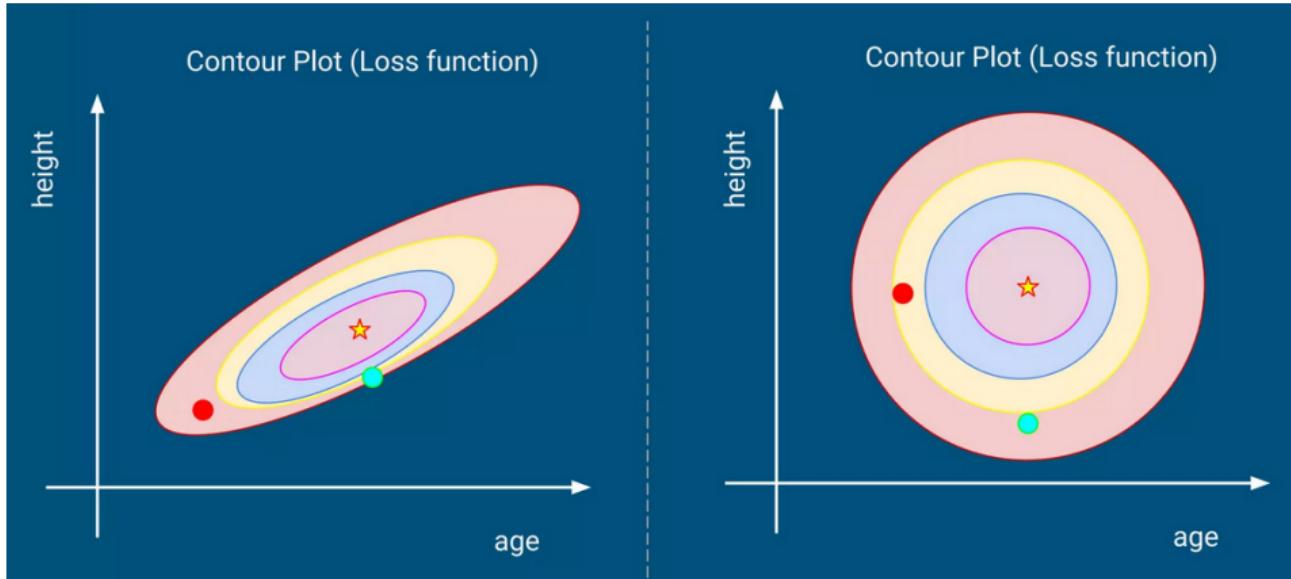


Figure: Feature normalization using BatchNormalization

$$y = \frac{x - \mathbb{E}[x]}{\sqrt{\text{Var}[x] + \varepsilon}} * \gamma + \beta$$

B

Optimization dynamics

Closed-form Gradient Descent update

Consider the objective:

$$L(\boldsymbol{\theta}) = \frac{1}{2} \boldsymbol{\theta}^\top \underbrace{(\mathbf{X}^\top \mathbf{X})}_{\mathbf{A}} \boldsymbol{\theta},$$

Write the gradient descent update as:

$$\begin{aligned}\boldsymbol{\theta}_{t+1} &\leftarrow \boldsymbol{\theta}_t - \alpha \nabla L(\boldsymbol{\theta}_t) \\ &= \boldsymbol{\theta}_t - \alpha \mathbf{A} \boldsymbol{\theta}_t \\ &= (\mathbf{I} - \alpha \mathbf{A}) \boldsymbol{\theta}_t\end{aligned}$$

Breaking the recursion:

$$\boldsymbol{\theta}_t = (\mathbf{I} - \alpha \mathbf{A})^t \boldsymbol{\theta}_0.$$

Gradient Descent, component-wise

Spectral decomposition: $\mathbf{A} = \mathbf{Q}\Lambda\mathbf{Q}^T$ so that we can do a change of basis:

$$\begin{aligned}(\mathbf{I} - \alpha\mathbf{A})^t \boldsymbol{\theta}_0 &= (\mathbf{I} - \alpha\mathbf{Q}\Lambda\mathbf{Q}^T)^t \boldsymbol{\theta}_0 \\&= [\mathbf{Q}(\mathbf{I} - \alpha\Lambda)\mathbf{Q}^T]^t \boldsymbol{\theta}_0 \\&= \mathbf{Q}(\mathbf{I} - \alpha\Lambda)^t \mathbf{Q}^T \boldsymbol{\theta}_0\end{aligned}$$

In the eigenspace of \mathbf{Q} , each coordinate is multiplied by $(1 - \alpha\lambda_i)^t$

Behaviours:

- $0 < \alpha\lambda_i \leq 1$: decays to 0 at a rate depending on λ_i
- $1 < \alpha\lambda_i \leq 2$: oscillations
- $\alpha\lambda_i > 2$: diverges

How does it look in practice?

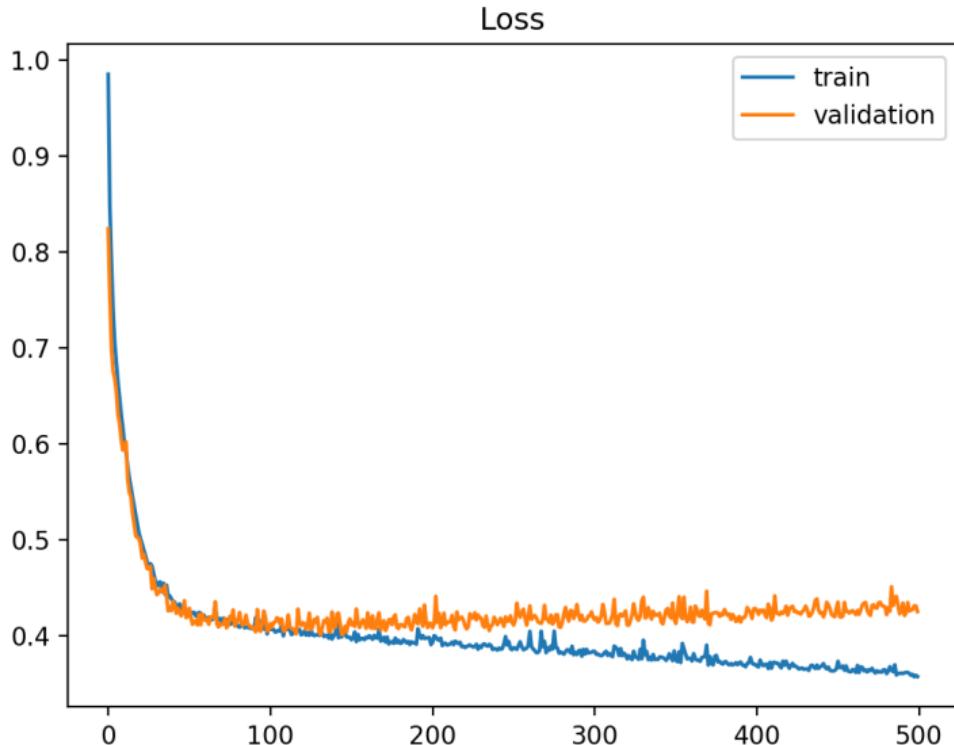


Figure: The optimizer is quickly descending along the direction of the largest eigenvalues.

Algorithms

Fix #1. Momentum

Plain gradient descent:

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t - \alpha \nabla_{\boldsymbol{\theta}_t} L$$

With [momentum](#):

$$\mathbf{v}_{t+1} \leftarrow \mu \mathbf{v}_t + \nabla_{\boldsymbol{\theta}_t} L$$

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t - \alpha \mathbf{v}_{t+1}.$$

- μ : [damping factor](#) and controls *how much momentum* we use.
 - Rolling ball on a slope analogy.
 - Keeps the *historic* direction in low curvature regions and *dampens* oscillations in high curvature directions.

Fix #1. Momentum

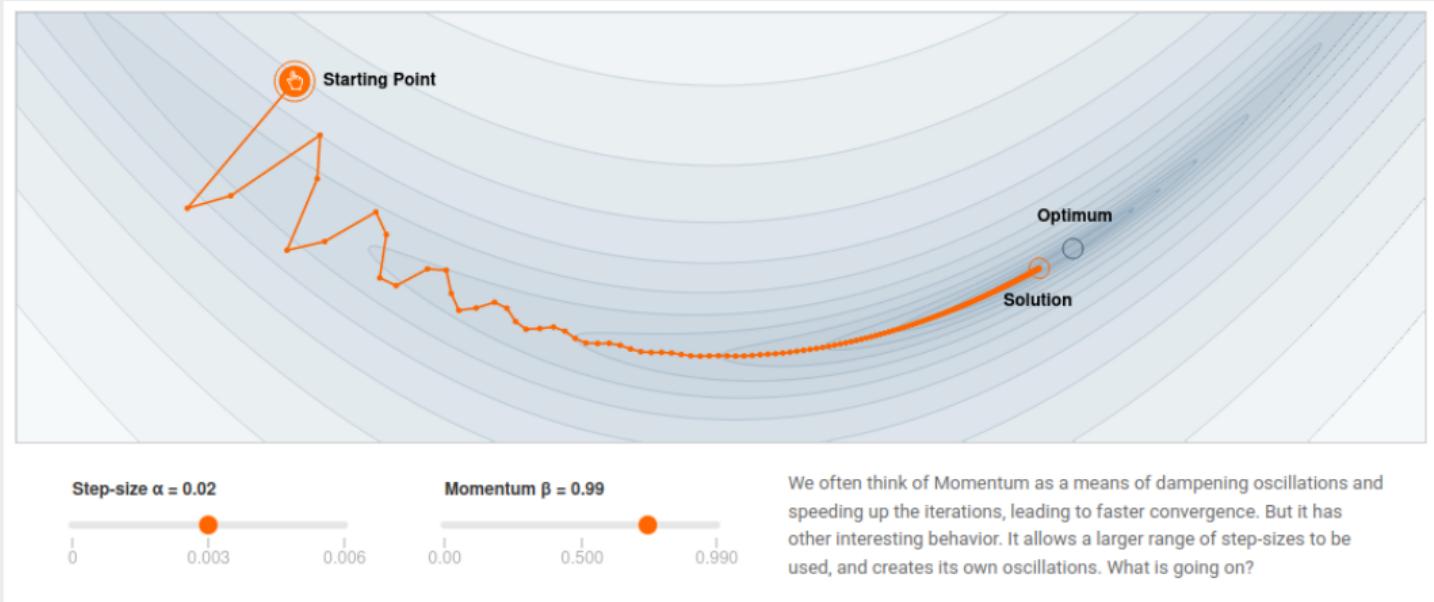


Figure: Why Momentum Really Works

Fix #1. Momentum

- no reason not to use it, it almost always helps...
- ...except when modelling distributions and small changes in θ means a large change in the distribution,
- low memory complexity,
- also check Nesterov momentum.

Fix #2. RMSprop

Plain gradient descent:

$$\theta_{t+1} \leftarrow \theta_t - \alpha \nabla_{\theta_t} L$$

With [RMSprop](#):

$$\begin{aligned}s_t &\leftarrow (1 - \gamma)s_t + \gamma[\nabla_{\theta_t} L]^2, \\ \theta_{t+1} &\leftarrow \theta_t - \frac{\alpha}{\sqrt{s_t + \epsilon}} \nabla_{\theta_t} L.\end{aligned}$$

- adaptive learning rate!
- [small components](#) will use a higher step size
- alternate view: the gradient is rescaled to have norm 1

Fix #3. Adam

Bring back **momentum**!

$$\begin{aligned}\mathbf{v}_t &\leftarrow \mu \mathbf{v}_{t-1} + (1 - \alpha) \nabla_{\theta_t} L \\ \mathbf{s}_t &\leftarrow \gamma \mathbf{s}_{t-1} + (1 - \gamma) [\nabla_{\theta_t} L]^2,\end{aligned}$$

The parameter update is then:

$$\theta_{t+1} \leftarrow \theta_t - \frac{\alpha}{\sqrt{\mathbf{s}_t + \epsilon}} \mathbf{v}_t.$$

Comparison

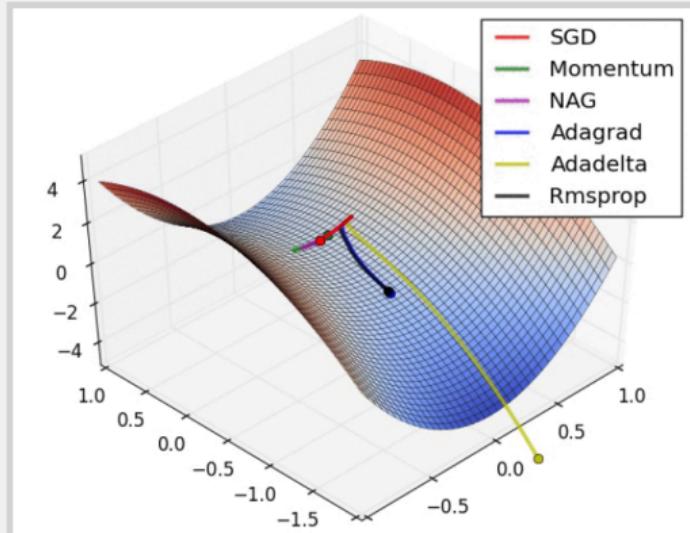
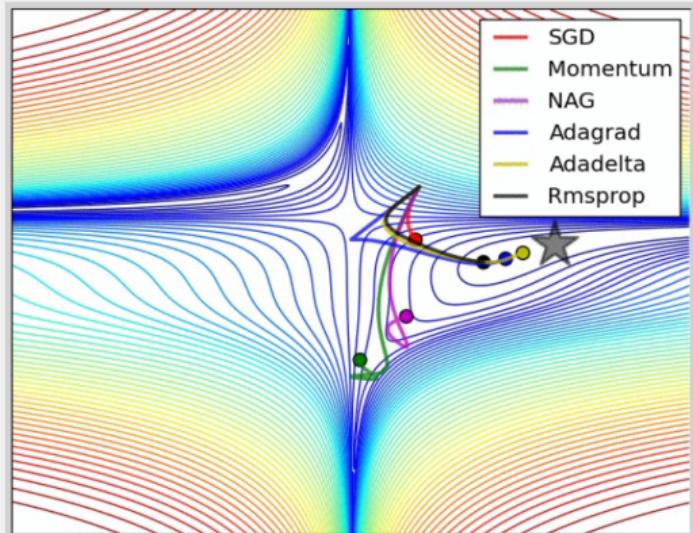


Figure: [check animation](#)

Regularization

Regularization. The modified objective

Remember the our purpose is to find Θ that minimizez the objective:

$$\begin{aligned}\hat{\Theta} &= \arg \min_{\Theta} [L(\Theta)] \\ &= \arg \min_{\Theta} \left[\sum_i l_i [\mathbf{x}_i, \mathbf{y}_i] \right]\end{aligned}$$

We can bias the optimization towards certain solutions:

$$\hat{\Theta} = \arg \min_{\Theta} \left[\sum_i l_i [\mathbf{x}_i, \mathbf{y}_i] + \lambda \cdot g(\Theta) \right],$$

where $g(\Theta)$ is a function that returns a (scalar) large value when the weights are large.

Regularization

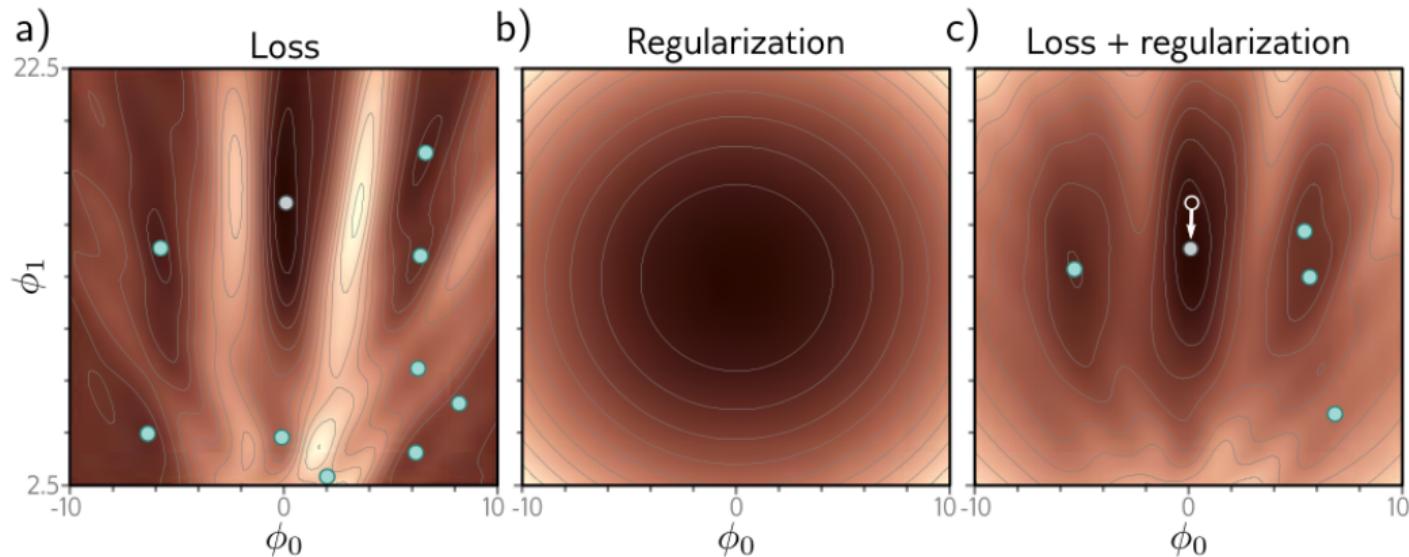


Figure: Effect of regularization objective.

$$\hat{\Theta} = \arg \min_{\Theta} \left[\sum_i l_i [\mathbf{x}_i, \mathbf{y}_i] + \lambda \cdot g(\Theta) \right],$$

B

Thank you!

florin.gogianu@gmail.com



bit-ml.github.io



References i