

یادگیری عمیق با تنسورفلو و کراس در پایتون

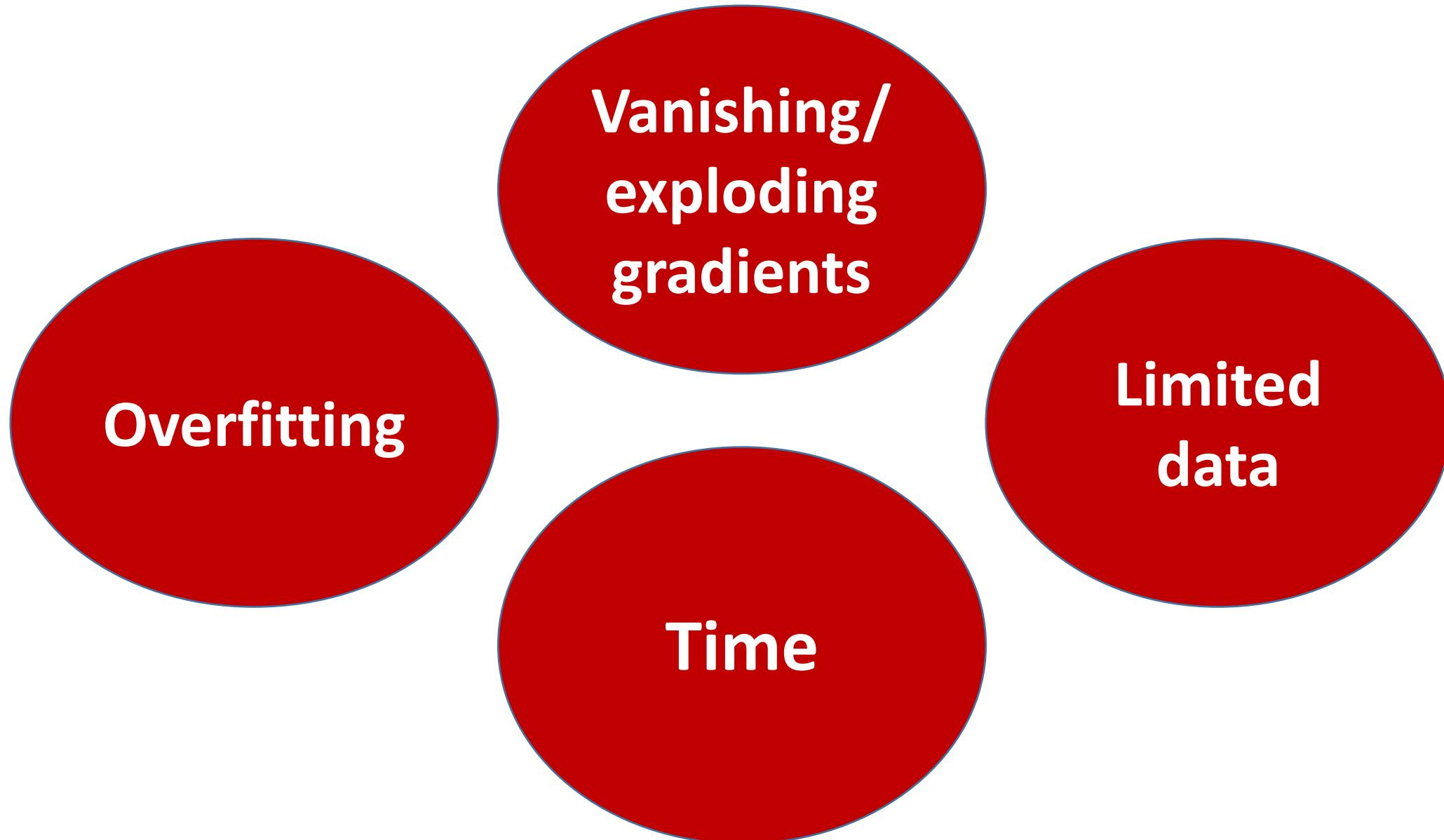
فصل سوم: آموزش شبکه های عصبی عمیق

پژمان اقبالی

PhD Student in Biomechanics

EPFL

چالش های آموزش شبکه های عصبی عمیق



چالش های آموزش شبکه های عصبی عمیق

Vanishing/exploding gradients

❑ Initialization

❑ Nonsaturating Activation Functions

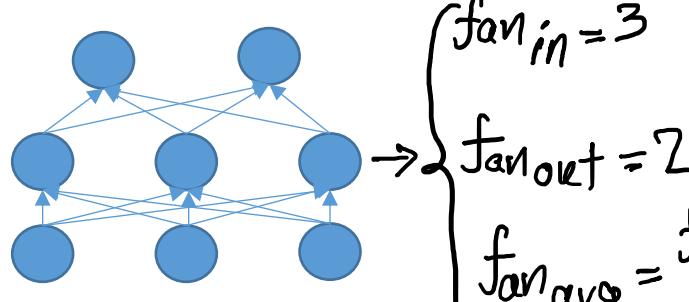
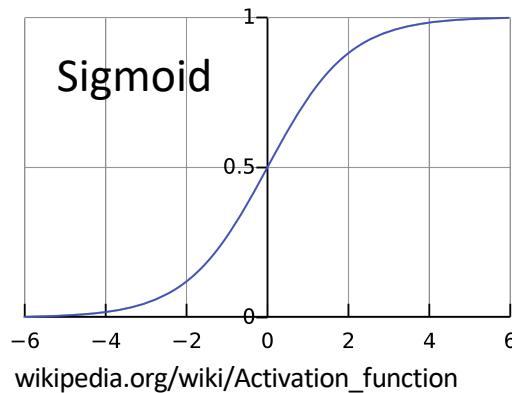
❑ Batch Normalization

❑ Gradient Clipping

چالش های آموزش شبکه های عصبی عمیق

Vanishing/exploding gradients

Initialization



Understanding the difficulty of training deep feedforward neural networks

Xavier Glorot

DIRO, Université de Montréal, Montréal, Québec, Canada

Yoshua Bengio

Glorot Initialization

- Normal distribution with mean 0 and $\sigma^2 = \frac{1}{\text{fan}_{avg}}$
- Uniform distribution between $-r$ and $+r$,
with $r = \sqrt{\frac{3}{\text{fan}_{avg}}}$

چالش های آموزش شبکه های عصبی عمیق

Vanishing/exploding gradients

□ Initialization

Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification

Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun

چالش های آموزش شبکه های عصبی عمیق

Vanishing/exploding gradients

□ Initialization

Activation functions		
Glorot	tanh, logistic, softmax	$\sigma^2(\text{Normal})$
He	RELU (and variants)	$1 / \text{fan}_{\text{avg}}$
LeCun	SELU	$2 / \text{fan}_{\text{in}}$

چالش های آموزش شبکه های عصبی عمیق

Vanishing/exploding gradients

❑ Nonsaturating Activation Functions

Empirical Evaluation of Rectified Activations in Convolutional Network

Bing Xu, Naiyan Wang, Tianqi Chen, Mu Li

- ❑ Leaky RELU (α)
- ❑ Randomized leaky RELU (RRELU)
- ❑ Parametric leaky RELU (PRELU)

چالش های آموزش شبکه های عصبی عمیق

Vanishing/exploding gradients

□ Nonsaturating Activation Functions

Fast and Accurate Deep Network Learning by Exponential Linear Units (ELUs)

Djork-Arné Clevert, Thomas Unterthiner, Sepp Hochreiter

چالش های آموزش شبکه های عصبی عمیق

Vanishing/exploding gradients

❑ Nonsaturating Activation Functions

Self-Normalizing Neural Networks

Günter Klambauer, Thomas Unterthiner, Andreas Mayr, Sepp Hochreiter

❑ Scaled ELU (SELU)

چالش های آموزش شبکه های عصبی عمیق

Vanishing/exploding gradients

❑ Nonsaturating Activation Functions

SELU > ELU > leaky RELU > RELU > tanh > logistic

Not self normalizing: **ELU**

Speed: **RELU**

چالش های آموزش شبکه های عصبی عمیق

Vanishing/exploding gradients

□ Batch Normalization

Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift

Sergey Ioffe, Christian Szegedy

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$;
Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

چالش های آموزش شبکه های عصبی عمیق

Vanishing/exploding gradients

□ Gradient Clipping

On the difficulty of training Recurrent Neural Networks

Razvan Pascanu, Tomas Mikolov, Yoshua Bengio

چالش های آموزش شبکه های عصبی عمیق

Limited data



لایه های از قبل آموزش دیده

- Transfer learning
- What Else?
- Unsupervised pretraining

چالش های آموزش شبکه های عصبی عمیق

Time



های سریعتر Optimizer

- Momentum optimization
- Nesterov
- AdaGrad
- RMSProp
- Adam and Nadam
- Learning rate scheduling

چالش های آموزش شبکه های عصبی عمیق

Time



های سریعتر Optimizer

❑ Stochastic Gradient Descent (SGD)

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

Require: Initial parameter θ

$k \leftarrow 1$

while **not** stoping criteria:

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

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

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

$k \leftarrow k + 1$

چالش های آموزش شبکه های عصبی عمیق

Time

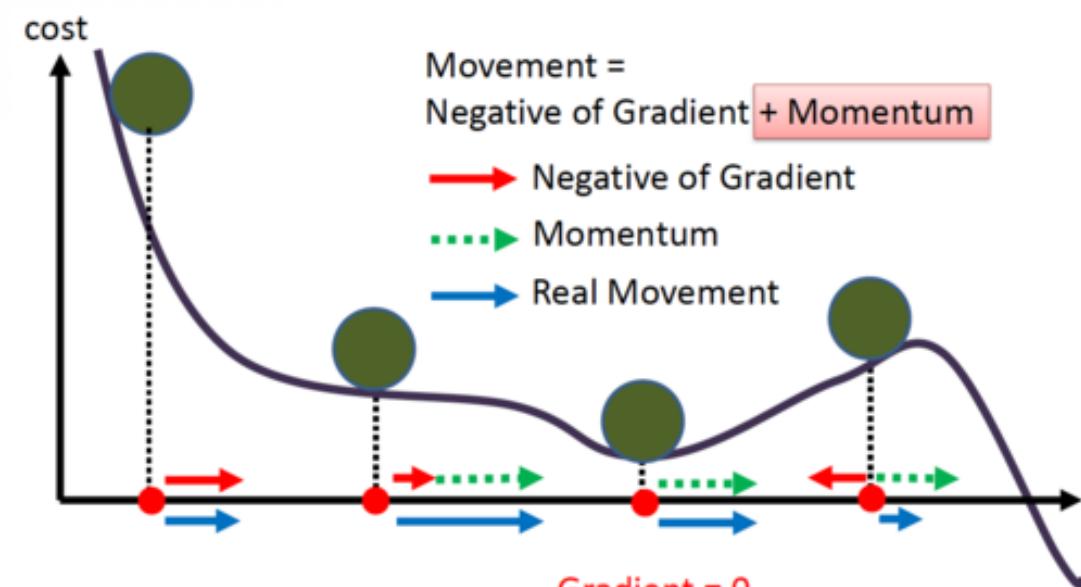
های سریعتر Optimizer

❑ Momentum optimization

Some methods of speeding up the convergence of iteration methods ☆

B.T. Polyak

[https://doi.org/10.1016/0041-5553\(64\)90137-5](https://doi.org/10.1016/0041-5553(64)90137-5)



چالش های آموزش شبکه های عصبی عمیق

□ Momentum

Require: Learning rate schedule ϵ , momentum parameter α

Require: Initial parameter θ , initial velocity v

while not stoping criteria:

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

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

 Compute velocity update: $v \leftarrow \alpha v - \epsilon \hat{g}$

 update: $\theta \leftarrow \theta + v$

چالش های آموزش شبکه های عصبی عمیق

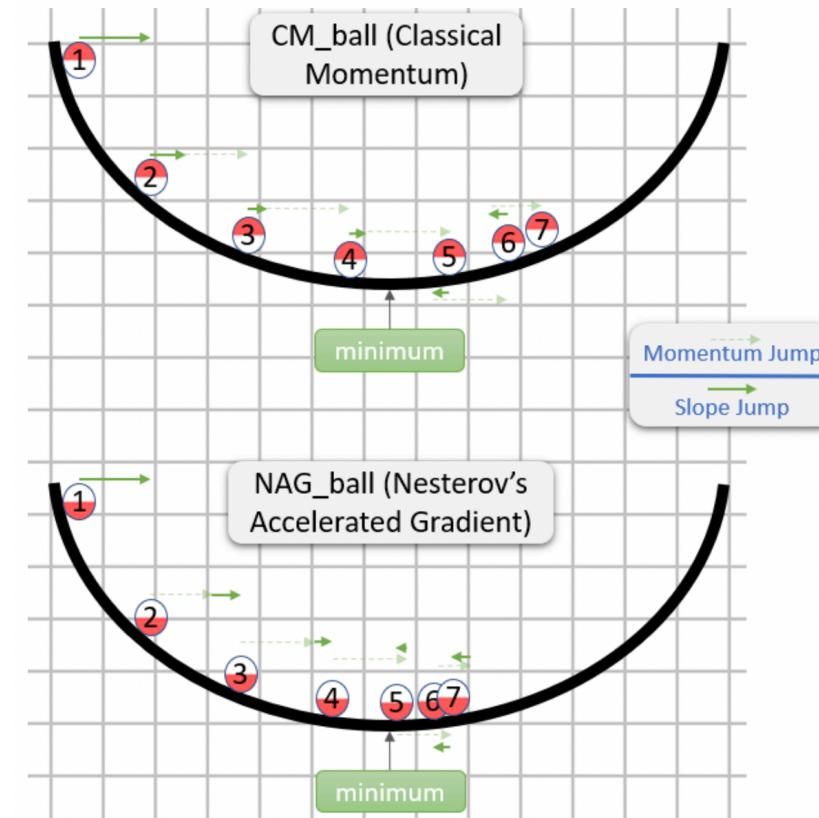
Time

های سریعتر Optimizer

Nesterov

A method for unconstrained convex minimization problem with the rate of convergence $\mathcal{O}(1/k^2)$

Yuriii Nesterov



چالش های آموزش شبکه های عصبی عمیق

❑ Nesterov Momentum

Require: Learning rate schedule ϵ , momentum parameter α

Require: Initial parameter θ , initial velocity v

while not stoping criteria:

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

 Apply interim update: $\tilde{\theta} \leftarrow \theta + \alpha v$

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

 Compute velocity update: $v \leftarrow \alpha v - \epsilon g$

 update: $\theta \leftarrow \theta + v$

چالش های آموزش شبکه های عصبی عمیق

Time



های سریعتر Optimizer

□ AdaGrad

Adaptive Subgradient Methods for Online Learning and Stochastic Optimization



Authors:  John Duchi,  Elad Hazan,  Yoram Singer [Authors Info & Claims](#)

The Journal of Machine Learning Research, Volume 12 • 2/1/2011 • pp 2121–2159

چالش های آموزش شبکه های عصبی عمیق

□ AdaGrad

Require: Global learning rate ϵ , Initial parameter θ , Small constant δ (for numerical instability = 1e-7)

Initialize gradient accumulation variable $r=0$

while **not** stoping criteria:

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

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

 Accumulate squared gradient : $r \leftarrow r + g \odot g$

 Compute update: $\Delta\theta \leftarrow -\frac{\epsilon}{\delta+\sqrt{r}} \odot g$

 update: $\theta \leftarrow \theta + \Delta\theta$

چالش های آموزش شبکه های عصبی عمیق

Time



های سریعتر Optimizer

☐ RMSProp

Neural Networks for Machine Learning

Lecture 6a
Overview of mini-batch gradient descent

Geoffrey Hinton
with
Nitish Srivastava
Kevin Swersky

چالش های آموزش شبکه های عصبی عمیق

□ RMSProp

Require: Global learning rate ϵ , decay rate ρ

Require: Initial parameter θ , Small constant δ (1e-6)

Initialize accumulation variable $r=0$

while **not** stoping criteria:

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

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

 Accumulate squared gradient : $r \leftarrow \rho r + (1 - \rho) g \odot g$

 Compute update: $\Delta \theta = -\frac{\epsilon}{\sqrt{\delta+r}} \odot g$

 update: $\theta \leftarrow \theta + \Delta \theta$

چالش های آموزش شبکه های عصبی عمیق

□ RMSProp with Nesterov

Require: Global learning rate ϵ , decay rate ρ , momentum coefficient α

Require: Initial parameter θ , initial velocity v

Initialize accumulation variable $r=0$

while **not** stoping criteria:

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

 Compute interim update: $\tilde{\theta} \leftarrow \theta + \alpha v$

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

 Accumulate squared gradient: $r \leftarrow \rho r + (1 - \rho) g \odot g$

 Compute velocity update: $v \leftarrow \alpha v - \frac{\epsilon}{\sqrt{r}} \odot g$

 update: $\theta \leftarrow \theta + v$

چالش های آموزش شبکه های عصبی عمیق

Time



های سریعتر Optimizer

□ Adam

Adam: A Method for Stochastic Optimization

Diederik P. Kingma, Jimmy Ba

چالش های آموزش شبکه های عصبی عمیق

□Adam

Require: Step size ϵ (suggested default: 0.001), Exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1]$ (suggested defaults: 0.9 and 0.999), Small constant δ (1e-8)

Require: Initial parameters θ

Initialize first and second moment variables $s=0$ and $r = 0$ and time step $t = 0$

while **not stoping criteria**:

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

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

$t \leftarrow t + 1$

 Update biased first moment estimate: $s \leftarrow \rho_1 s + (1 - \rho_1)g$

 Update biased second moment estimate: $r \leftarrow \rho_2 r + (1 - \rho_2)g \odot g$

 Correct bias in first and second moment: $\hat{s} \leftarrow \frac{s}{1 - \rho_1^t}, \hat{r} \leftarrow \frac{r}{1 - \rho_2^t}$

 Compute update: $\Delta\theta = -\epsilon \frac{\hat{s}}{\sqrt{\hat{r}} + \delta}$

 update: $\theta \leftarrow \theta + \Delta\theta$

چالش های آموزش شبکه های عصبی عمیق

Time



های سریعتر Optimizer

□ Adam -> AdaMax

while **not** stoping criteria:

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

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

$t \leftarrow t + 1$

Update biased first moment estimate: $s \leftarrow \rho_1 s + (1 - \rho_1)g$

Update the exponentially weighted infinity norm : $r \leftarrow \max(\rho_2 r, |g|)$

Compute update: $\Delta \theta = -\frac{\epsilon}{1-\rho_1^t} \frac{s}{r}$

update: $\theta \leftarrow \theta + \Delta \theta$

چالش های آموزش شبکه های عصبی عمیق

Time



های سریعتر Optimizer

Adam -> Nadam

Workshop track - ICLR 2016

INCORPORATING NESTEROV MOMENTUM INTO ADAM

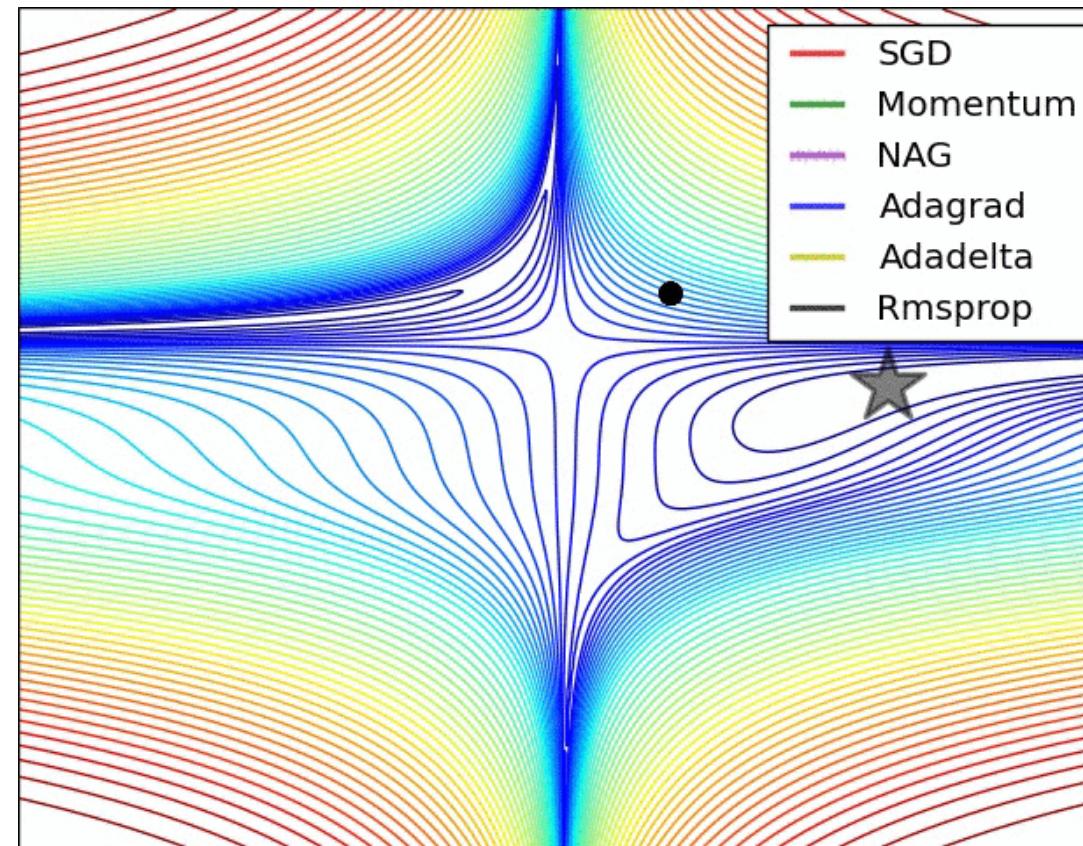
Timothy Dozat

چالش های آموزش شبکه های عصبی عمیق

Time



های سریعتر Optimizer



چالش های آموزش شبکه های عصبی عمیق

Time



های سریعتر Optimizer

کیفیت همگرایی	سرعت همگرایی	Optimizer
3	1	SGD
3	2	SGD(momentum)
3	2	SGD(momentum, nesterov)
1	3	Adagrad
2-3	3	RMSprop
2-3	3	Adam
2-3	3	Nadam
2-3	3	AdaMax

چالش های آموزش شبکه های عصبی عمیق

Time



برنامه های زمانی نرخ یادگیری

☐ Learning Rate

☐ Power scheduling

$$\eta(t) = \eta_0 / (1 + t/s)^c$$

☐ Exponential scheduling

$$\eta(t) = \eta_0 0.1^{t/s}$$

☐ Piecewise constant scheduling

☐ Performance scheduling

☐ 1cycle scheduling

N. Smith, 2018
arXiv:1801.09821

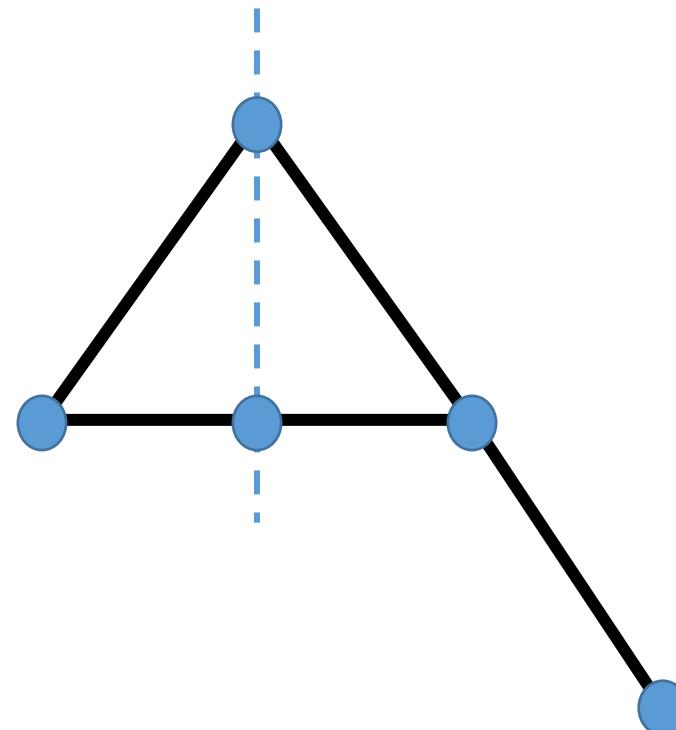
چالش های آموزش شبکه های عصبی عمیق

Time

برنامه های زمانی نرخ یادگیری

Learning Rate

1cycle scheduling



چالش های آموزش شبکه های عصبی عمیق

Overfitting



Regularization

- l1 and l2 regularization
- Dropout
- Max-Norm

چالش های آموزش شبکه های عصبی عمیق

Overfitting



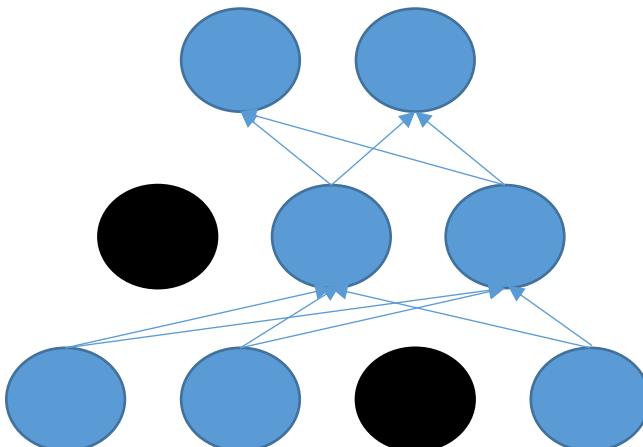
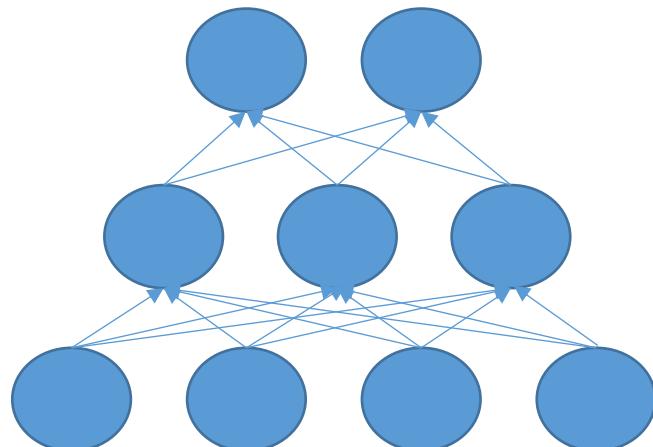
Regularization

□ Dropout

JMLR

Dropout: A Simple Way to Prevent Neural Networks from Overfitting

Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, Ruslan Salakhutdinov; 15(56):1929–1958, 2014.



چالش های آموزش شبکه های عصبی عمیق

Overfitting



Regularization

□ Dropout -> Monte Carlo

Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning

Yarin Gal, Zoubin Ghahramani

چالش های آموزش شبکه های عصبی عمیق

Overfitting



Regularization

Max-Norm

$$w \leftarrow w * c / \| w \|_2$$

چالش های آموزش شبکه های عصبی عمیق

Default Config

Regular Network

Hyperparameter	Default
Kernel initializer	He initialization
Activation function	ELU
Normalization	Batch (None for shallow)
Regularization	Early stopping (l2)
Optimizer	Momentum (RMSProp or Nadam)
Learning rate schedule	1cycle

Self-normalizing Network

Hyperparameter	Default
Kernel initializer	LeCun
Activation function	SELU
Normalization	None
Regularization	Alpha dropout
Optimizer	Momentum (RMSProp or Nadam)
Learning rate schedule	1cycle