

Introduction to Deep Learning (Optimization Algorithms)

WM CS
Zeyi (Tim) Tao
11/01/2019

Topics

SGD

Robbins , Monro : A Stochastic Approximation Method

<https://projecteuclid.org> › euclid.aoms ▼

by H Robbins - 1951 - Cited by 6678 - Related articles

Let $M(x)$ denote the expected value at level x of the response to a certain experiment. $M(x)$ is assumed to be a monotone function of x but is unknown to the ...

SGDM

On the momentum term in gradient descent learning algorithms

<https://www.sciencedirect.com> › [science](#) › [article](#) › [pii](#)

by N Qian - 1999 - Cited by 855 - Related articles

The behavior of gradient descent near a local minimum is equivalent to a set of coupled and damped harmonic oscillators. Within a reasonable parameter range, the momentum term can improve the speed of convergence for most eigen components in the system by bringing them closer to critical damping.

AdaGrad

[PDF] Adaptive Subgradient Methods for Online Learning and ...

www.jmlr.org › [papers](#) › [volume12](#) ▼

by J Duchi - 2011 - Cited by 5323 - Related articles

Before introducing our adaptive gradient algorithm, which we term ADAGRAD, we establish notation. Vectors and scalars are lower case italic letters, such as x ...

AdaDelta

ADADELTA: An Adaptive Learning Rate Method

<https://arxiv.org> › [cs](#) ▼

by MD Zeiler - 2012 - Cited by 3495 - Related articles

Dec 22, 2012 - Abstract: We present a novel per-dimension learning rate method for gradient descent called ADADELTA. The method dynamically adapts over ...

Adam

Adam: A Method for Stochastic Optimization

<https://arxiv.org> › [cs](#) ▼

by DP Kingma - 2014 - Cited by 33005 - Related articles

Some connections to related algorithms, on which Adam was inspired, are ... We also analyze the theoretical convergence properties of the algorithm and ...

Cite as: [arXiv:1412.6980](https://arxiv.org/abs/1412.6980)

Adaptive LR

RMSprop

Neural Networks for Machine Learning

Lecture 6a

Overview of mini-batch gradient descent

Geoffrey Hinton

ERM problem Statement

- Given a training set of input-output pairs $(X_1, d_1), (X_2, d_2), \dots, (X_N, d_N)$
- Minimize the following function (w.r.t W)

The diagram illustrates the components of the ERM loss function. It features a central equation with four yellow boxes pointing to its parts: 'Total Loss' points to the entire equation; 'Sum over all training instances' points to the summation symbol; 'Output of net in response to input' points to the function $f(W_i; X)$; and 'Desired output (label)' points to d_i . Below the equation, two more yellow boxes point upwards: 'Measurement functions' points to the div operator, and 'Regularizer' points to $\gamma(w)$.

$$Loss(W) = \frac{1}{N} \sum_{i=1}^N div(f(W_i; X), d_i) + \gamma(w)$$

- This is problem of function minimization

Choice of div() functions

Mean Square Loss VS CrossEntropy

Prediction Mode: $f_w(x_i) = wx_i + b$

Square Loss: $Loss_w = \frac{1}{N} \sum_i^N (d_i - f_x(x_i))^2$

Suppose: $d_i = w_i x_i + b + n$ where $n \sim \text{Normal}(0,1)$

$$E[d_i] = E[w_i x_i + b + n] = w_i x_i + b$$

$$\text{Var}[d_i] = \text{Var}[w_i x_i + b + n] = 1$$

Choice of div() functions

Mean Square Loss VS CrossEntropy

The probability of observing a single (x_i, d_i)

$$p(d_i | x_i) = e^{-\frac{(d_i - (wx_i + b))^2}{2}}$$

The Max Likelihood:

$$Like(d, x) = \prod_{i=1}^N e^{-\frac{(d_i - (wx_i + b))^2}{2}}$$

Choice of div() functions

Mean Square Loss VS CrossEntropy

$$\text{Like}(d, x) = \prod_{i=1}^N e^{-\frac{(d_i - (wx_i + b))^2}{2}}$$

$$l(d, x) = -\frac{1}{2} \sum_{i=1}^N (d_i - (wx_i + b))^2 \quad (\text{MAX})$$

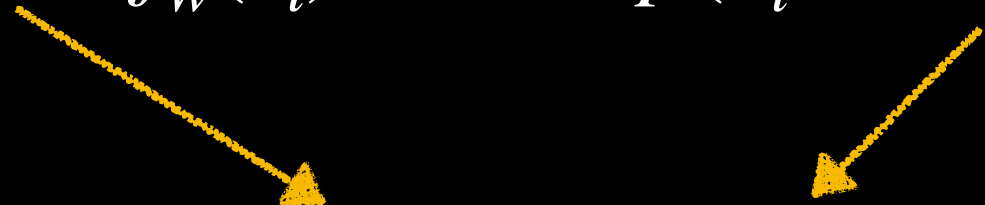
$$\text{MSE}(d, x) = \frac{1}{2} \sum_{i=1}^N (d_i - f_w(x_i))^2 \quad (\text{MIN})$$

Choice of div() functions

Mean Square Loss VS CrossEntropy

Prediction Mode: $f_w(x_i) = \sigma(wx_i + b)$ assume $\sigma()$ is softmax liked

$$p(d_i = 1 | x_i) = f_w(x_i) \quad p(d_i = 0 | x_i) = 1 - f_w(x_i)$$


$$p(d_i | x_i) = [f_w(x_i)]^{d_i} [1 - f_w(x_i)]^{(1-d_i)}$$

$$Like(d, x) = \prod_{i=1}^N [f_w(x_i)]^{d_i} [1 - f_w(x_i)]^{(1-d_i)}$$

Choice of div() functions

Mean Square Loss VS **CrossEntropy**

$$Like(d, x) = \prod_{i=1}^N [f_w(x_i)]^{d_i} [1 - f_w(x_i)]^{(1-d_i)}$$

$$like(d, x) = - \sum_{i=1}^N d_i \log(f_w(x_i)) + (1 - d_i) \log(1 - f_w(x_i))$$

(binary)

$$like(d, x) = - \sum_{i=1}^N \sum_{j=1}^N d_{ij} \log(f_w(x_i)_j) + (1 - d_{ij}) \log(1 - f_w(x_i)_j)$$

(multi-class)

Choice of div() functions

Mean Square Loss VS CrossEntropy

$$f_w(x_i) = \sigma(wx_i + b) = \hat{d}_i$$


$$Loss = \frac{1}{2} \sum_{i=1}^N (\hat{d}_i - d_i)^2$$

$$Loss = - \sum_{i=1}^N d_i \log(\hat{d}_i) + (1 - d_i) \log(1 - \hat{d}_i)$$

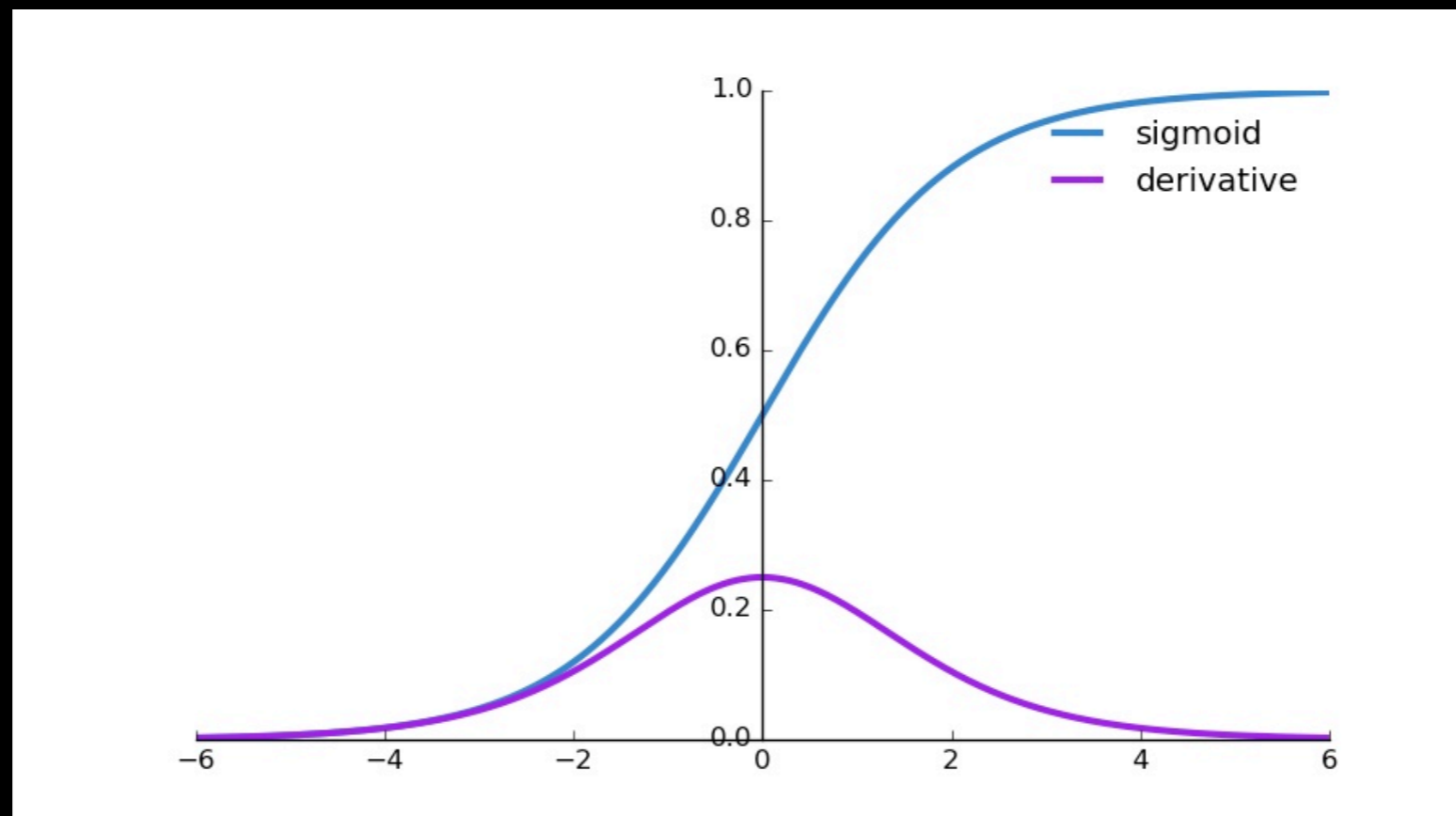
$$\frac{dLoss}{dw} = \sum_{i=1}^N (\hat{d}_i - d_i) \sigma'(wx_i + b) x_i$$

$$\frac{dLoss}{dw} = \sum_{i=1}^N x_i (\sigma(wx_i + b) - d_i)$$

Choice of div() functions

Mean Square Loss VS CrossEntropy

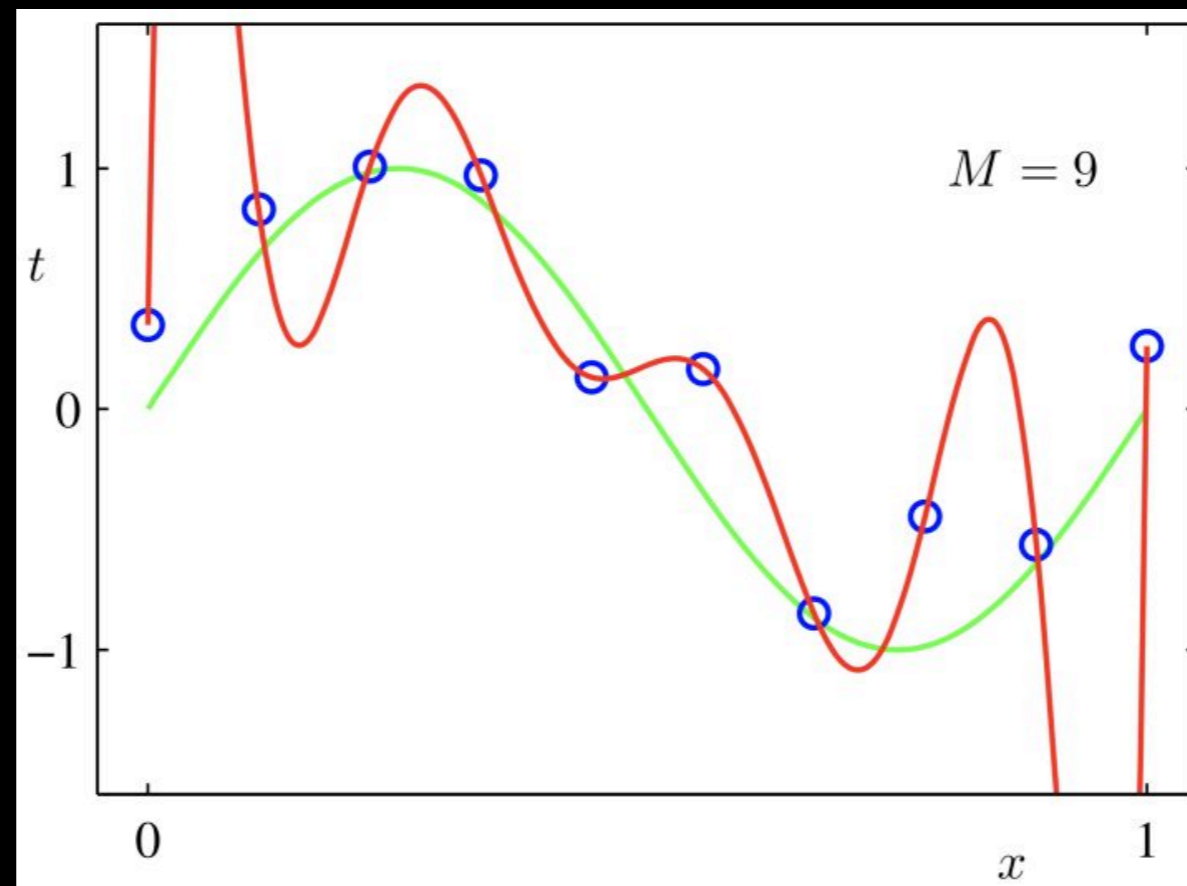
$$\frac{dLoss}{d_w} = \sum_{i=1}^N (\hat{d}_i - d_i) \sigma'(wx_i + b) x_i \quad \frac{dLoss}{d_w} = \sum_{i=1}^N x_i (\sigma(wx_i + b) - d_i)$$



Regularization

$$Loss(W) = \frac{1}{N} \sum_{i=1}^N div(f(W_i; X), d_i) + \lambda \gamma(w)$$

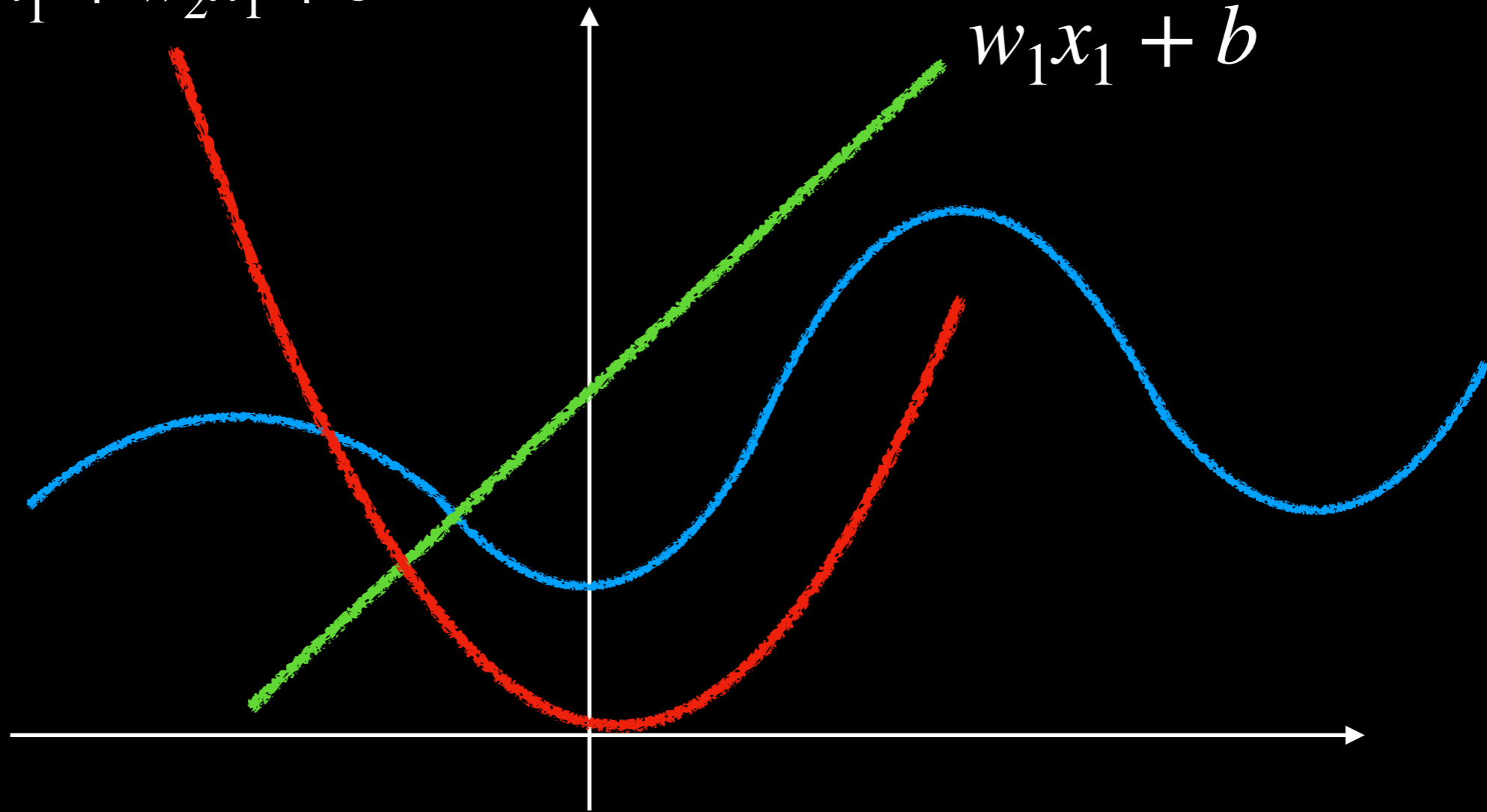
$$\gamma(w) = \frac{1}{2} ||w||^2 = \frac{1}{2} ||w - 0||^2$$



Regularization

$$w_1 x_1^2 + w_2 x_1 + b$$

$$w_1 x_1 + b$$

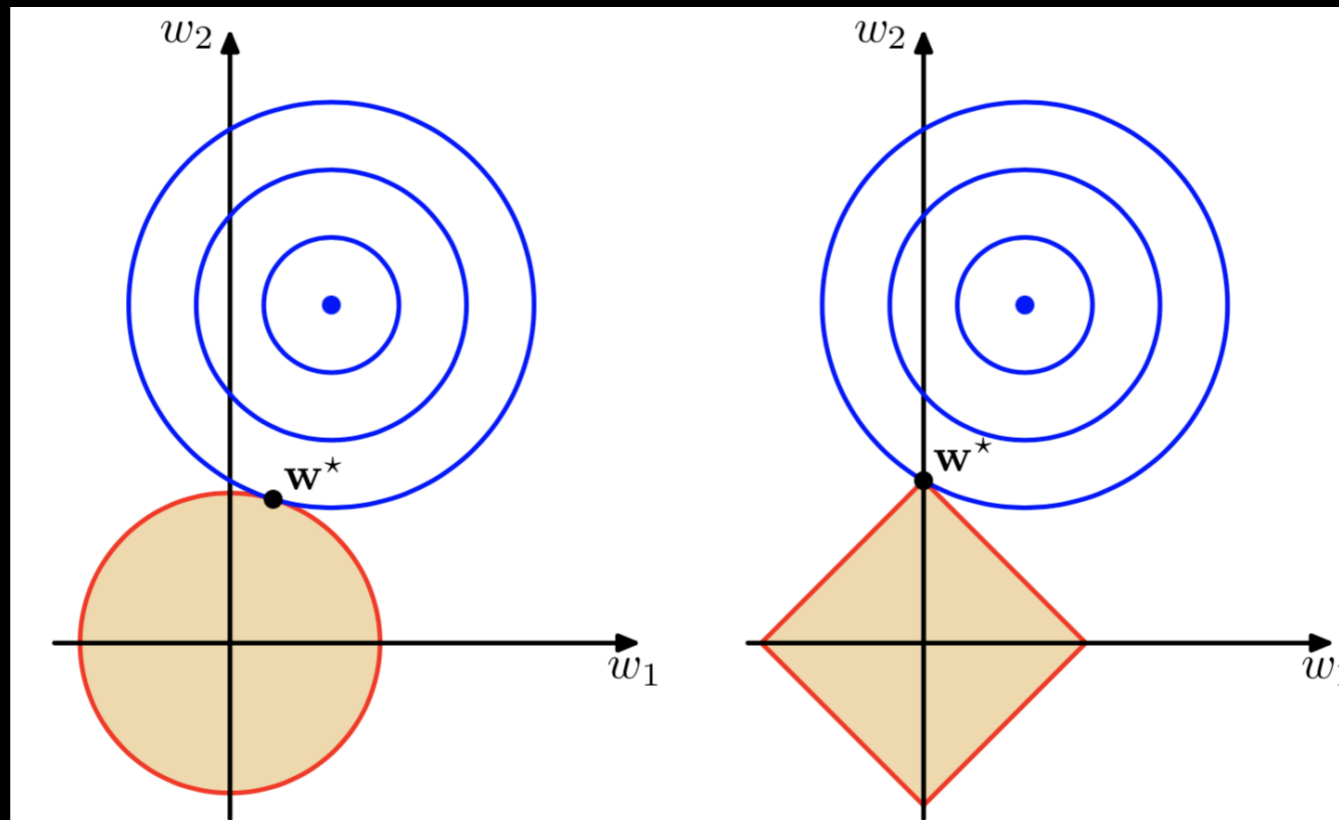


$$w_1 x^n + w_2 x^{n-1} + \dots + w_{n-1} x^2 + w_n x_1 + b$$

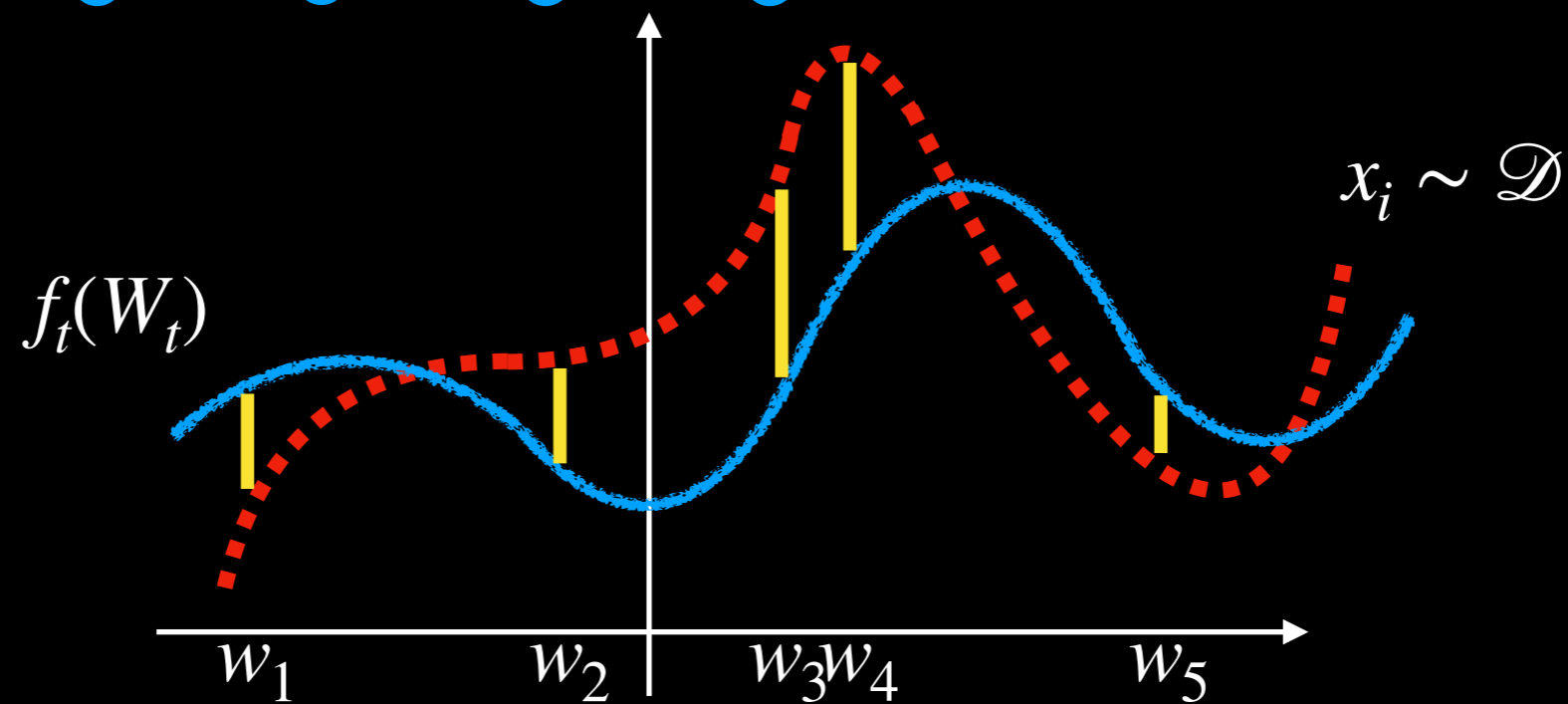
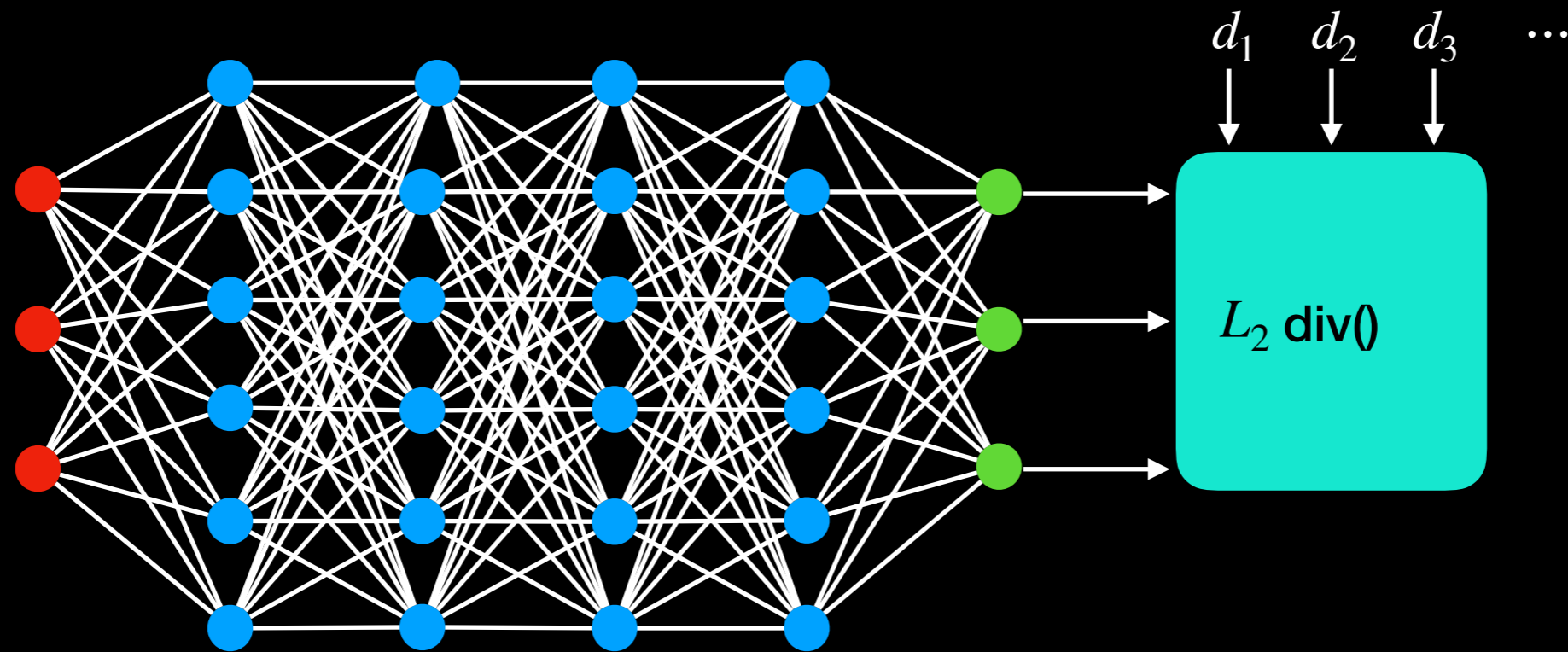
Regularization

$$\min_w \frac{1}{N} \sum_{i=1}^N \text{div}(f(W_i; X), d_i) + \lambda \gamma(w)$$

$$\min_w \frac{1}{N} \sum_{i=1}^N \text{div}(f(W_i; X), d_i) + 100w_n^2 + 100w_{n-1}^2 + \dots$$



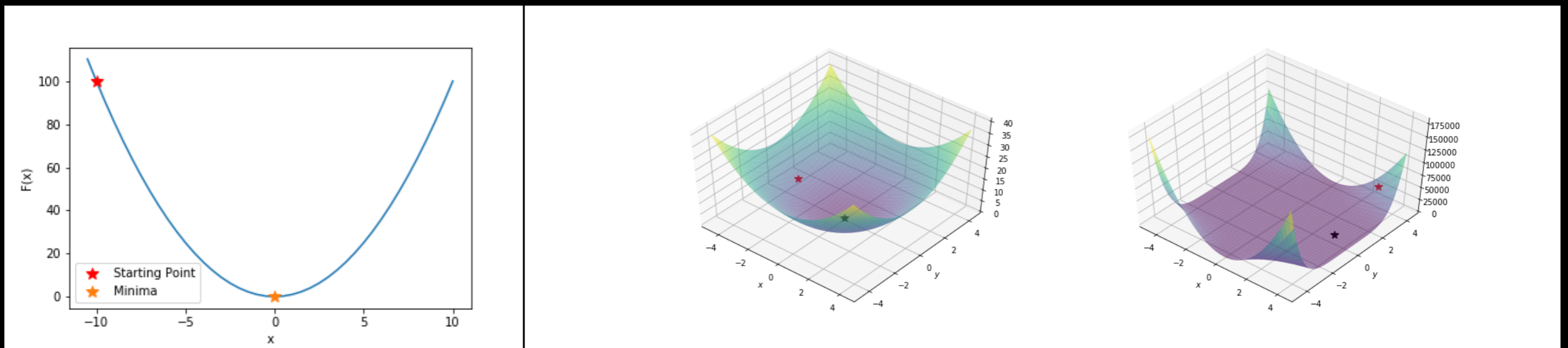
Deep neuron network



Unconstrained First-order Optimization

- For real-valued output vectors, the (scaled) L_2 divergence is popular

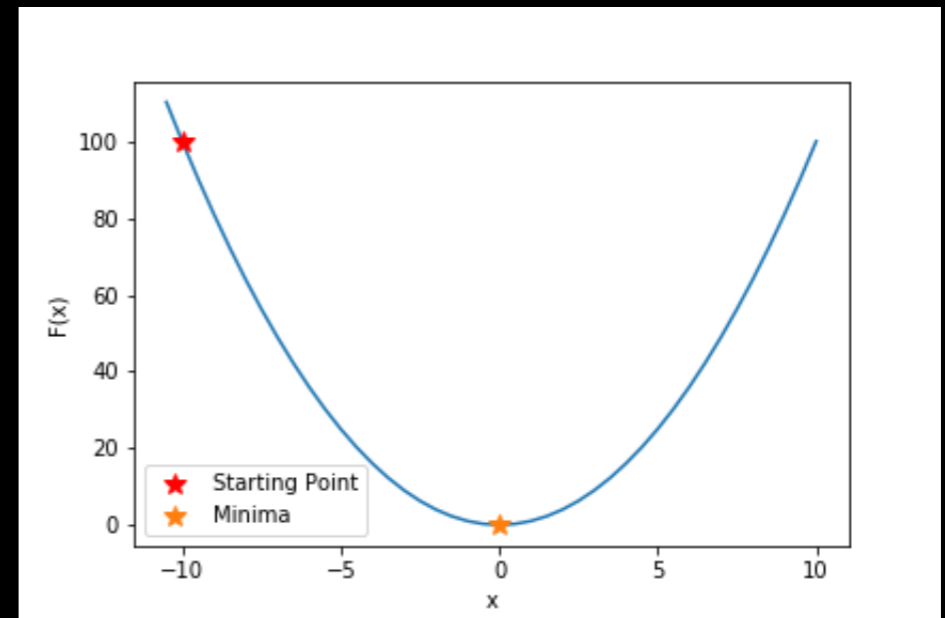
$$\text{Div}(\hat{d}, d) = \frac{1}{2} \|\hat{d} - d\|^2 = \frac{1}{2} \sum_i (\hat{d}_i - d_i)^2$$



Unconstrained First-order Optimization

ERM: convex

Local optimal = Global optimal



First Order: Gradient Descent

$$f(w + \Delta w) = f(w) + f'(w) \Delta w \quad \text{Linear approximation}$$

Second Order: Newton Method

$$f(w + \Delta w) = f(w) + f'(w) \Delta w + \frac{1}{2} f''(w) (\Delta x)^2$$

Gradient descent variants

Batch gradient descent

$$w_t = w_{t-1} - \eta^k \nabla_w f(W; X = \{(x_1, d_1), (x_2, d_2), \dots\})$$

Stochastic gradient descent

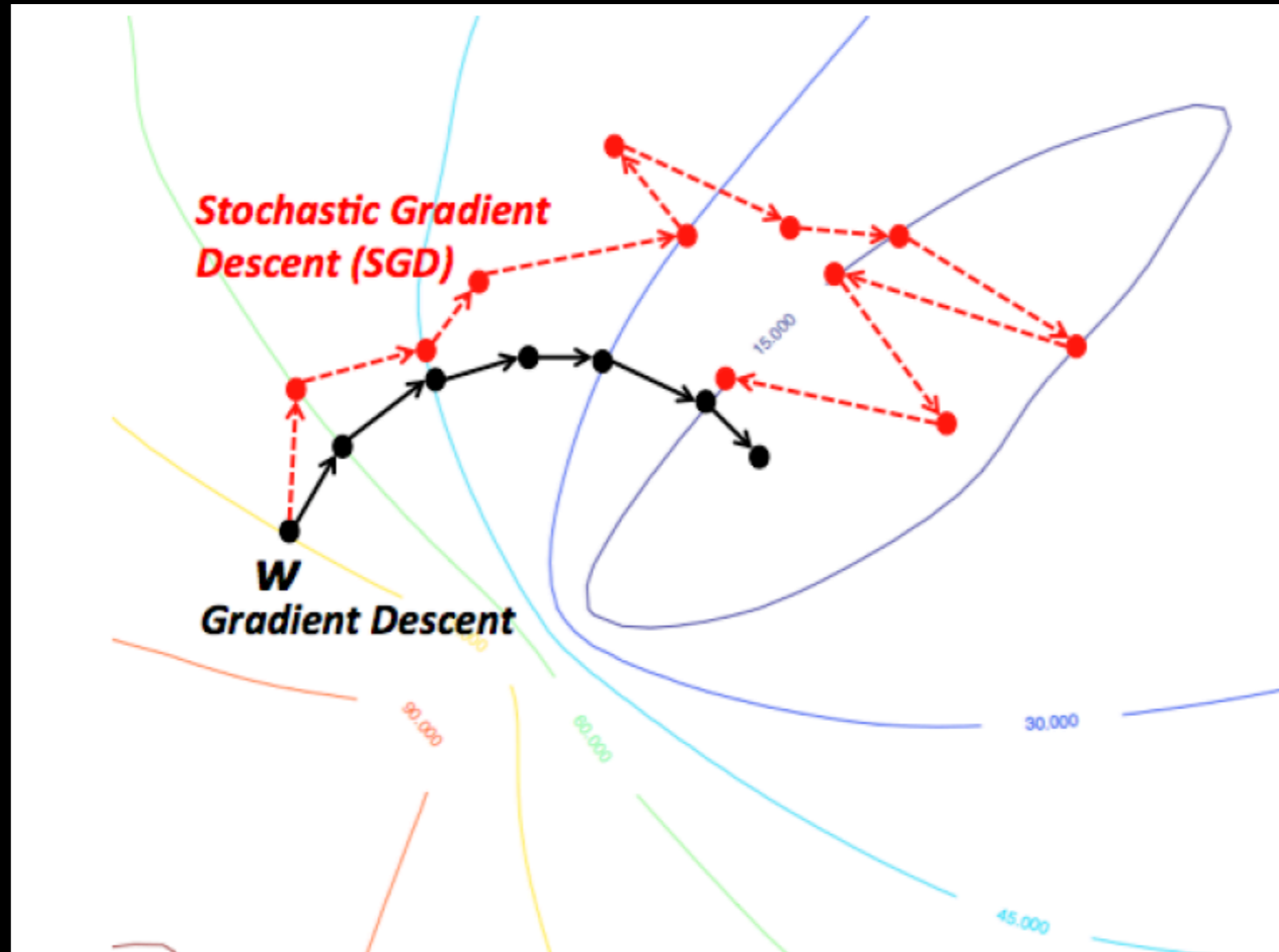
$$w_t = w_{t-1} - \eta^k \nabla_w f(W; (x_i, d_i))$$

Mini-batch gradient descent

$$w_t = w_{t-1} - \eta^k \nabla_w f(W; X = \{(x_i, d_i) \mid i = 1, 2, \dots, b\})$$

Shuffle

Batch Gradient descent VS SGD



Gradient descent variants

Batch gradient descent

- **Pro**
- **Less oscillations and noise**
- **Vectorization**
- **Stable**

- **Con**
- **Local optimal**
- **Not memory friendly**

Stochastic gradient descent

- **Pro**
- **Computationally fast**
- **Fast convergence (large dataset)**
- **Fit into memory**

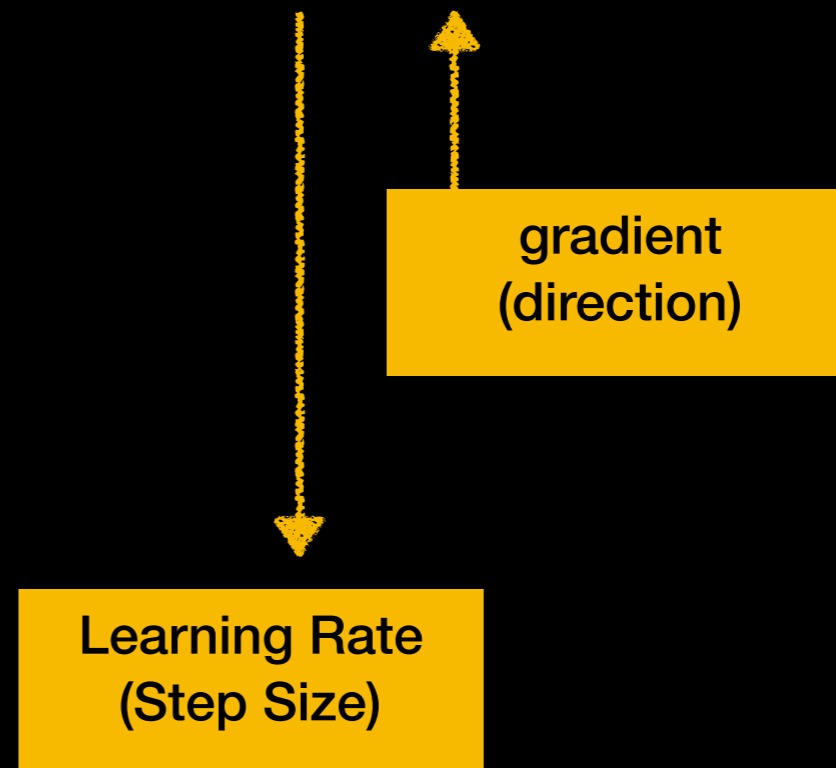
- **Con**
- **Noisy**
- **Longer training time**
- **No vectorization**

Mini-batch gradient descent

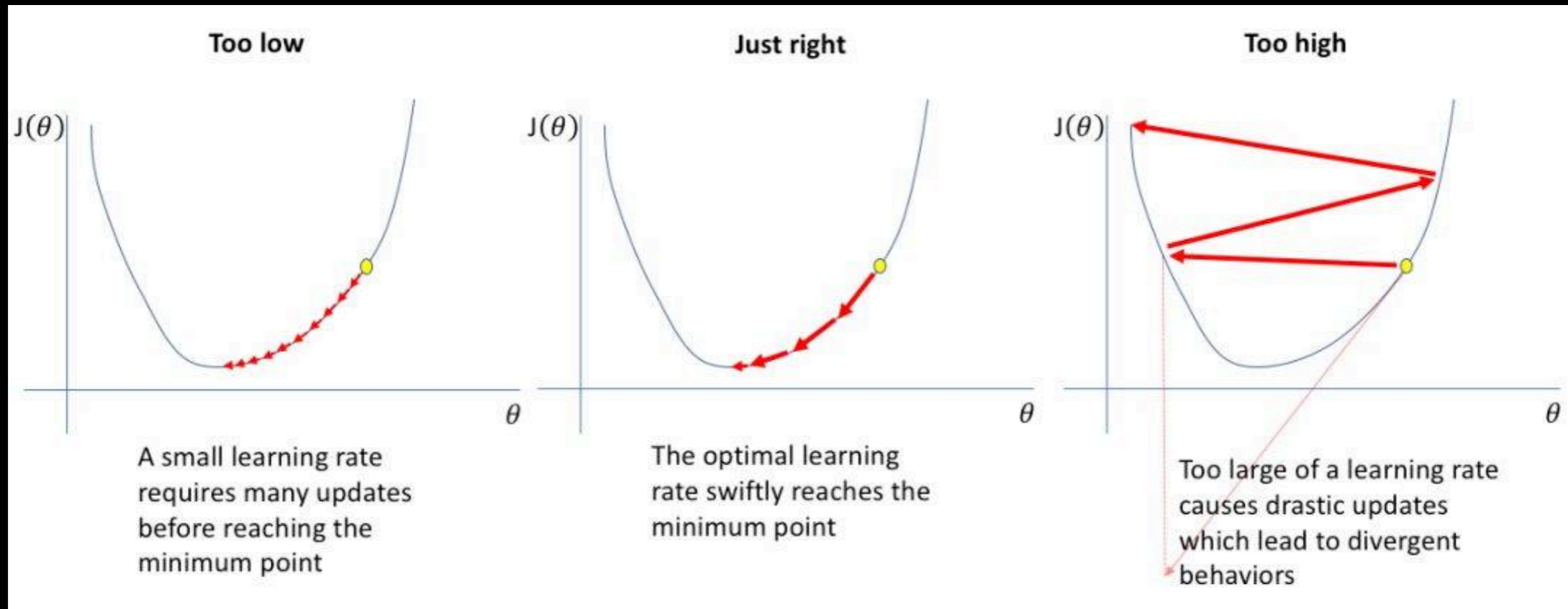
Mini-batch gradient descent (Vanilla SGD)

$$w_t = w_{t-1} - \eta^k \nabla_w f(W; (x_i, d_i))$$

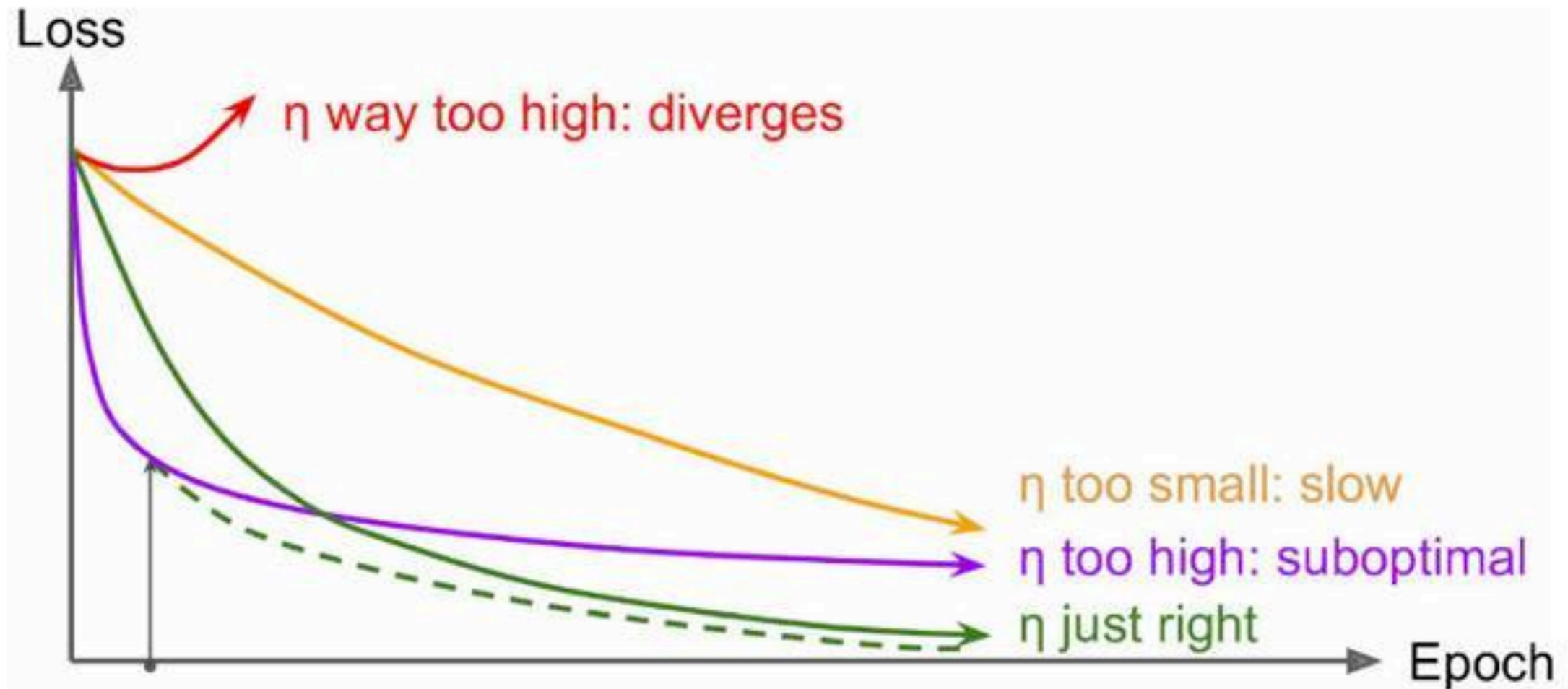
$$w_t = w_{t-1} - \eta_k g_t$$



Learning Rate I

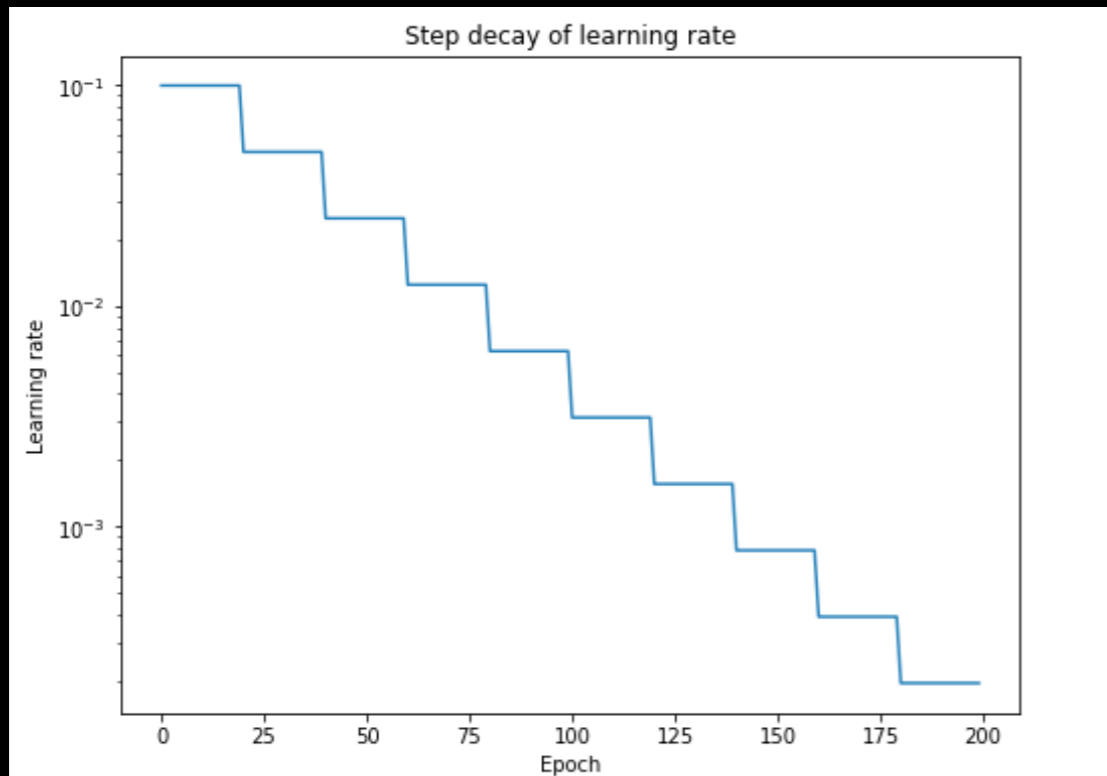


Learning Rate II



Start with a high learning rate then reduce it: perfect! <https://www.youtube.com/watch?v=JNingWei>

Learning Rate III: Strategies



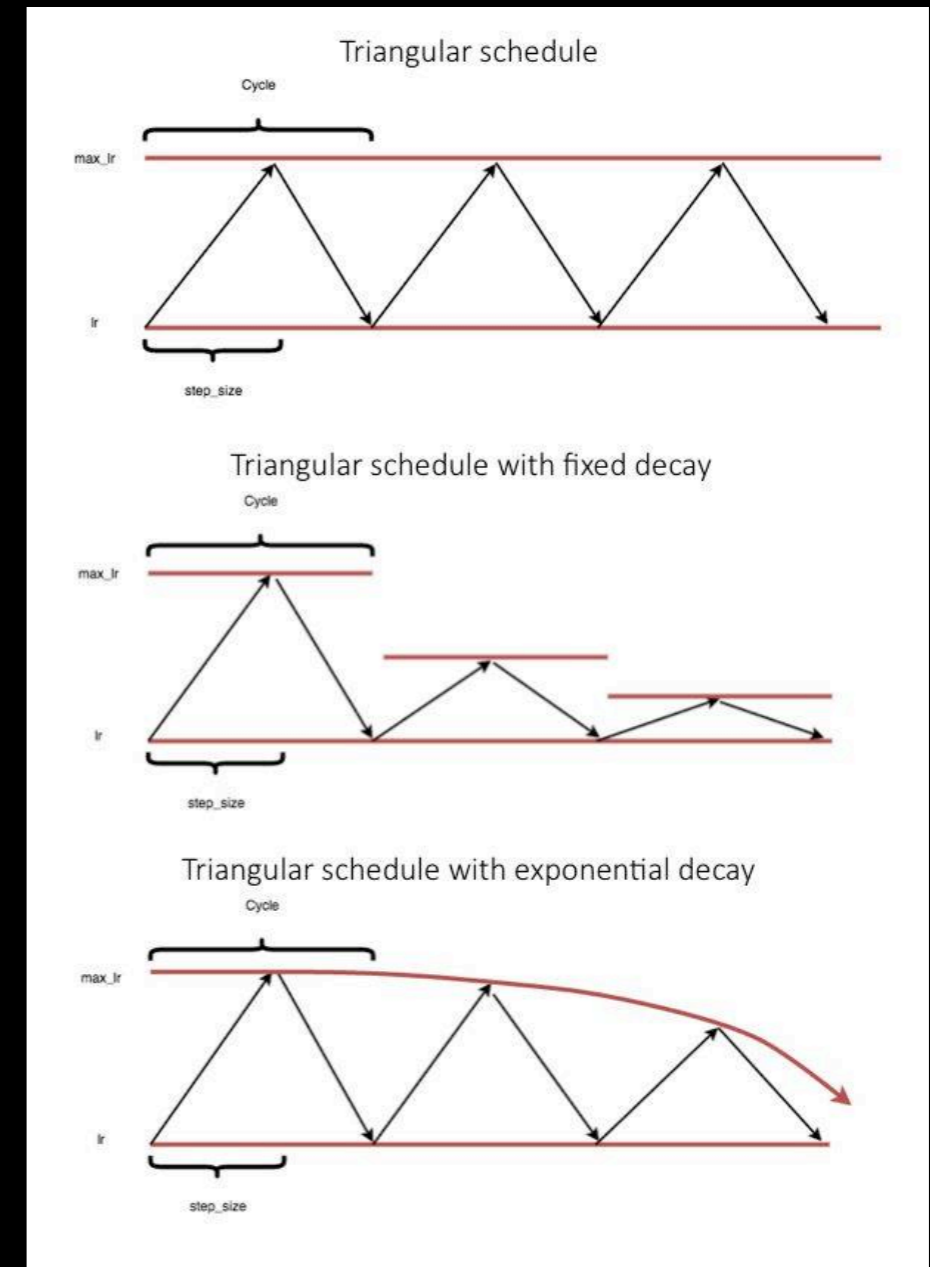
Step decay

CLASS `torch.optim.lr_scheduler.StepLR(optimizer, step_size, gamma=0.1, last_epoch=-1)` [\[SOURCE\]](#)

Sets the learning rate of each parameter group to the initial lr decayed by gamma every step_size epochs. When last_epoch=-1, sets initial lr as lr.

Parameters

- **optimizer** (*Optimizer*) – Wrapped optimizer.
- **step_size** (*int*) – Period of learning rate decay.
- **gamma** (*float*) – Multiplicative factor of learning rate decay. Default: 0.1.
- **last_epoch** (*int*) – The index of last epoch. Default: -1.



Cyclical Learning Rate

SGD to SGD with Momentum(SGDM)

On the momentum term in gradient descent learning algorithms

<https://www.sciencedirect.com/science/article/pii>

by N Qian - 1999 - Cited by 855 - Related articles

The behavior of gradient descent near a local minimum is equivalent to a set of coupled and damped harmonic oscillators. Within a reasonable parameter range, the momentum term can improve the speed of convergence for most eigen components in the system by bringing them closer to critical damping.

$$v_t = \beta v_{t-1} + \eta g_t$$

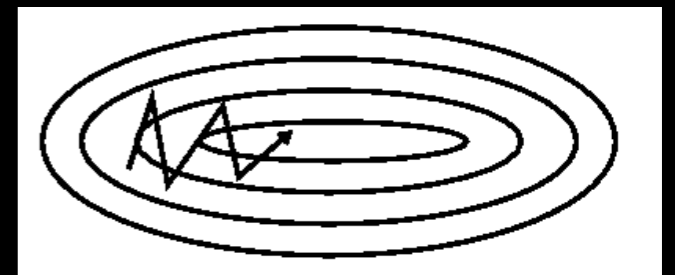
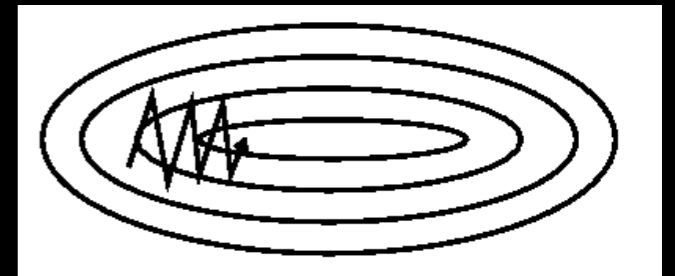
$$w_t = w_{t-1} - v_t$$

Pro:

- Accelerate SGD
- Overcome local minima
- Dampen oscillations (begin)

Con:

- oscillations (end)



The momentum term **increases** for dimensions whose gradients point in the **same directions** and reduces updates for dimensions whose gradients change directions.

SGD to **Adaptive** Subgradient Methods (AdaGrad)

[PDF] Adaptive Subgradient Methods for Online Learning and ...

www.jmlr.org > papers > volume12 ▾

by J Duchi - 2011 - Cited by 5323 - Related articles

Before introducing our adaptive gradient algorithm, which we term ADAGRAD, we establish notation. Vectors and scalars are lower case italic letters, such as x ...

$$g_{t,i} = \nabla_{w_{t,i}} f$$

$$G_{t,i} = \sum_{i=1}^t g_{t,i}^2$$

$$w_{t+1,i} = w_{t,i} - \frac{\eta_0}{\sqrt{G_{t,i} + \epsilon}} \odot g_{t,i}$$

Pro:

- **Suit for dealing with sparse data**

Con:

- **Monotonically decreasing LR**

Performing smaller updates (i.e. low learning rates) for parameters associated with frequently occurring features, and larger updates (i.e. high learning rates) for parameters associated with infrequent features.

AdaGrad to AdaDelta

ADADELTA: An Adaptive Learning Rate Method

<https://arxiv.org> > cs ▼

by MD Zeiler - 2012 - Cited by 3495 - Related articles

Dec 22, 2012 - Abstract: We present a novel per-dimension learning rate method for gradient descent called ADADELTA. The method dynamically adapts over ...

$$g_{t,i} = \nabla_{w_{t,i}} f$$

$$v_{t,i} = \beta v_{t-1,i} + (1 - \beta) g_{t,i}^2$$

$$w_{t+1,i} = w_{t,i} - \frac{\eta_0}{\sqrt{v_{t,i} + \epsilon}} \odot g_{t,i}$$

Pro:

- **Suit for dealing with sparse data**
- **Using a sliding window**

Instead of accumulating all past squared gradients, Adadelta **restricts the window of accumulated past gradients** to some fixed size

AdaGrad to Adam

Adam: A Method for Stochastic Optimization

<https://arxiv.org> > cs ▾

by DP Kingma - 2014 - Cited by 33005 - Related articles

Some connections to related algorithms, on which Adam was inspired, are ... We also analyze the theoretical convergence properties of the algorithm and ...

Cite as: [arXiv:1412.6980](https://arxiv.org/abs/1412.6980)

$$g_{t,i} = \nabla_{w_{t,i}} f$$

$$m_{t,i} = \beta_1 m_{t-1,i} + (1 - \beta_1) g_{t,i}$$

$$v_{t,i} = \beta_2 v_{t-1,i} + (1 - \beta_2) g_{t,i}^2$$

$$w_{t+1,i} = w_{t,i} - \frac{\eta_0}{\sqrt{v_{t,i} + \epsilon}} \odot m_{t,i}$$

Pro:

- Super fast (current)
- Exponential Moving exponential (EMA)

Con:

- Suboptimal solution
- EMA diminishes the changes of gradients

Adam to ?

[\[PDF\] on the convergence of adam and beyond - OpenReview](#)

<https://openreview.net> › [pdf](#) ▼

by SJ Reddi - 2018 - [Cited by 425](#) - [Related articles](#)

ON THE CONVERGENCE OF ADAM AND BEYOND. Sashank J. Reddi, Satyen Kale & Sanjiv Kumar.

Google New York. New York, NY 10011, USA. {sashank ...

You've visited this page many times. Last visit: 10/21/19

[The Marginal Value of Adaptive Gradient Methods in Machine ...](#)

<https://papers.nips.cc> › [paper](#) › [7003-the-marginal-value-of-adaptive-gradi...](#) ▼

by AC Wilson - 2017 - [Cited by 288](#) - [Related articles](#)

The Marginal Value of Adaptive Gradient Methods in Machine Learning. Part of: Advances in Neural Information Processing Systems 30 (NIPS 2017).

You've visited this page 3 times. Last visit: 11/9/19

Improving Generalization Performance by Switching from Adam to SGD

[Nitish Shirish Keskar](#), [Richard Socher](#)

(Submitted on 20 Dec 2017)

Decoupled Weight Decay Regularization

[Ilya Loshchilov](#), [Frank Hutter](#)

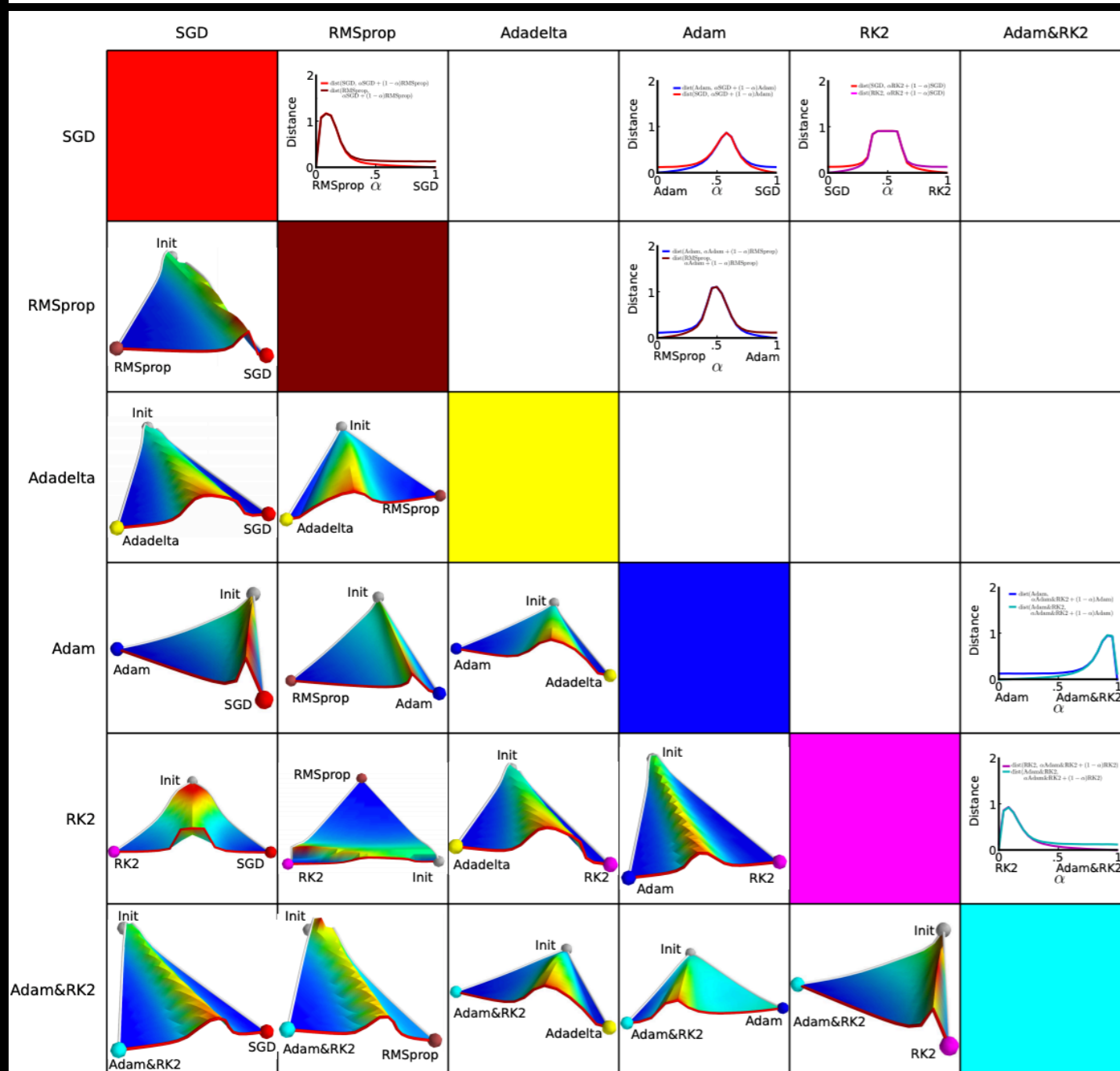
(Submitted on 14 Nov 2017 (v1), last revised 4 Jan 2019 (this version, v3))

Which to use?

An empirical analysis of the optimization of deep network loss surfaces

Daniel Jiwoong Im, Michael Tao, Kristin Branson

(Submitted on 13 Dec 2016 (v1), last revised 7 Dec 2017 (this version, v4))



- SGD + SGDM
- Familiar with
- Knowing your data
- Test on small batch
- Adam + SGD
- Shuffle
- Choosing prop LR