



M2177.003100

Deep Learning

[4: Optimization]

Electrical and Computer Engineering
Seoul National University

© 2019 Sungroh Yoon. this material is for educational uses only. some contents are based on the material provided by other paper/book authors and may be copyrighted by them.

(last compiled at 12:49:00 on 2019/09/22)

Outline

Introduction

Gradient-Based Optimization

Additional Topics

Summary

References

- *Deep Learning* by Goodfellow, Bengio and Courville [▶ Link](#)
 - ▶ Chapter 8: Optimization for Training Deep Models
- online resources:
 - ▶ *Deep Learning Specialization (coursera)* [▶ Link](#)
 - ▶ *Stanford CS231n: CNN for Visual Recognition* [▶ Link](#)

Outline

Introduction

Gradient-Based Optimization

Additional Topics

Summary

Optimization in deep learning

- most difficult optimization task in DL: training
 - ▶ so important and so expensive \Rightarrow need specialized techniques
- mainstream: **stochastic gradient descent** (sgd) and its variants
- more complicated methods: _____
 - ▶ **second-order** methods
 - ▷ often too expensive to compute/store Hessian
 - ▷ memory-efficient techniques emerging
 - ▶ **convex optimization**
 - ▷ its importance greatly diminished
- for clarity: this lecture focuses on unregularized supervised case

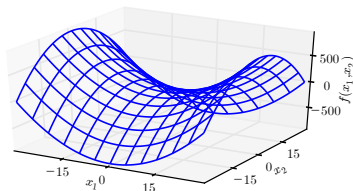
Derivatives and optimization order

- derivatives
 - ▶ first derivative (= gradient) \Rightarrow slope (Jacobian)
 - ▶ second derivative \Rightarrow _____ (Hessian)
- optimization
 - ▶ first-order algorithms
 - ▷ use only gradient (e.g. gradient descent)
 - ▶ second-order algorithms
 - ▷ also use Hessian matrix (e.g. Newton's method)

Critical points (= stationary points)

- points with zero slope: $\nabla_x f(\mathbf{x}) = \mathbf{0}$
 - ▶ derivative gives no info about which direction to move
 - ▶ three types: maxima (− curvature), minima (+ curvature), saddle points
- a **saddle point**: contains ____ positive and negative curvature

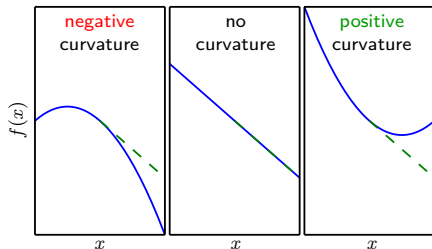
$$f(\mathbf{x}) = x_1^2 - x_2^2$$



- along x_1 axis: f curves upward
 - ▶ direction of $\text{eigenvec}(\mathbf{H})$ with $\lambda > 0$
 - ▶ local minimum
- along x_2 axis: f curves downward
 - ▶ direction of $\text{eigenvec}(\mathbf{H})$ with $\lambda < 0$
 - ▶ local maximum

Use of second derivative

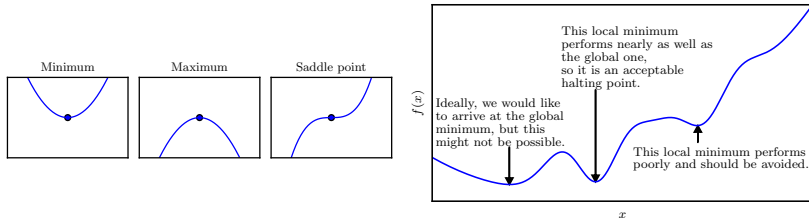
1. to characterize critical points
2. to measure curvature
3. to predict performance of an update in _____-based optimization



- **negative** curvature
 - ▶ f decreases faster than gradient predicts
- no curvature
 - ▶ gradient predicts the decrease correctly
- **positive** curvature
 - ▶ f decreases slower than gradient predicts (eventually increases)

In deep learning

- our objective function has
 - ▶ many local minima + many saddle points surrounded by very flat regions
 - ⇒ makes optimization very difficult (especially in high-dim space)
- we therefore usually settle for finding a **very low value of f**
 - ▶ **not necessarily minimal** in any formal sense
- recent research (Dauphin, 2014) reports
 - ▶ in high dim: _____ are much more common than **local minima**



Outline

Introduction

Gradient-Based Optimization

- Gradient Descent and its Limitations

- Exponentially Weighted Average

- Gradient Descent with Momentum

- Per-Parameter Adaptive Learning Rates

Additional Topics

Summary

Method of gradient descent

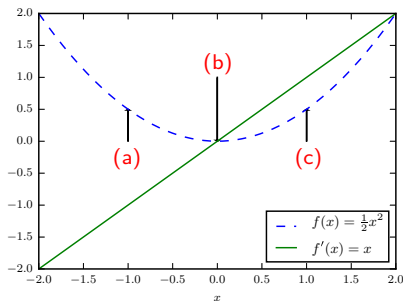
- derivative: useful for minimizing a function

▶ for small ϵ : $f(x - \epsilon \cdot \text{sign}[f'(x)]) < f(x)$

- we can thus reduce $f(x)$ by

▶ moving x in small steps with **opposite sign of derivative**

= **method of gradient descent**



$$\underbrace{x'}_{\text{new}} = \underbrace{x}_{\text{old}} - \underbrace{\epsilon}_{\text{learning rate}} \nabla_x f(x)$$

(a) $x < 0$: $f'(x) < 0$

⇒ can decrease f by moving rightward

(b) $x = 0$: $f'(x) = 0$

⇒ gd halts here (global min)

(c) $x > 0$: $f'(x) > 0$

⇒ can decrease f by moving leftward

Sgd and its variants

- probably the most used optimization algorithms for ML/DL

- ▶ can obtain an unbiased estimate of gradient

↑
by taking average gradient on a minibatch of m examples

Algorithm 1 gradient descent

- 1: initialize θ
 - 2: **while** stopping criterion not met **do**
 - 3: sample m examples: $\mathbb{X}_m = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$
 - 4: compute gradient estimate: $\hat{\mathbf{g}} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_{i=1}^m L(y^{(i)}, \hat{y}^{(i)})$ ▷ m forward props
 - 5: apply update: $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$ ▷ ϵ : learning rate
 - 6: **end while**
-

- three variants (N : total number of examples)
 - ▶ $m = 1$: stochastic gradient descent (sgd)
 - ▶ $1 < m < N$: _____ sgd (typical m : 64, 128, 256, 512)
 - ▶ $m = N$: batch gradient descent

Properties of sgd: good ones

- property #1:

computation time per update does not grow with # of training examples

- ▶ most important property of sgd/minibatch/online gradient-based optimization

⇒ allows convergence even when # of training examples becomes large

- property #2 (see textbook):

sgd works better in practice than its theoretical analysis says

- ▶ some benefits of sgd: obscured in asymptotic analysis

Properties of sgd: bad ones

- sgd may suffer in the following situations:

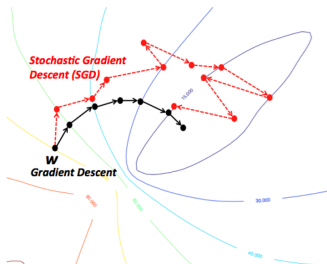
- ▶ local minima/saddle points



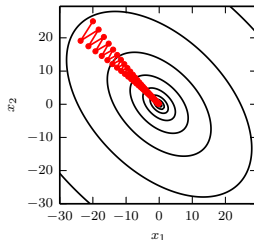
zero gradient

⇒ gradient descent gets stuck

- ▶ gradient noise



- ▶ poor _____ of H



(source: wikidocs, cs231n)

Ravine

- Chloe Kim (2018 Olympic Champion, Women's Snowboard Halfpipe) [▶ Clip](#)



Poor conditioning of H

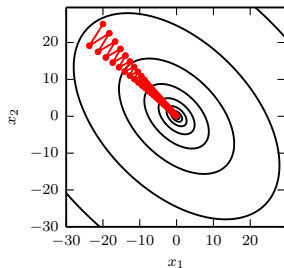
- consider a point x in multiple dimensions:
 - ▶ different second derivative for each direction
- condition number of Hessian H at x
 - ▶ measures how much the _____ derivatives differ from each other
 - ▶ recall: condition number of a matrix with eigenvalues $\{\lambda\}$:

$$\max_{i,j} \left| \frac{\lambda_i}{\lambda_j} \right|$$

- when H has a large condition number (“poorly conditioned”)
 1. gradient descent performs poorly
 - ∴ in a direction, derivative increases rapidly; in another, it increases slowly
 - ▶ gradient descent¹ is unaware of this change in the derivative
 2. it is difficult to choose a good step size ϵ

¹it does not know it needs to explore preferentially in the direction where derivative remains negative for longer

example:



- assume: Hessian H has condition number 5
 - ▶ most curvature: 5 times more curvature than least (a long canyon)
 - most curvature: direction $[1, 1]$ ↗
 - least curvature: direction $[1, -1]$ ↘
- gradient descent (red lines): slow (zig-zag)
- by contrast: methods considering H
 - ▶ can predict: the steepest direction is not promising
(large $\lambda > 0 \Rightarrow$ large positive curvature \Rightarrow bad; see page 8)
- how to handle poor conditioning _____ directly considering H ?

Outline

Introduction

Gradient-Based Optimization

Gradient Descent and its Limitations

Exponentially Weighted Average

Gradient Descent with Momentum

Per-Parameter Adaptive Learning Rates

Additional Topics

Summary

Exponentially weighted moving average (EWMA)

- given: time series g_1, g_2, \dots
- EWMA defined as:

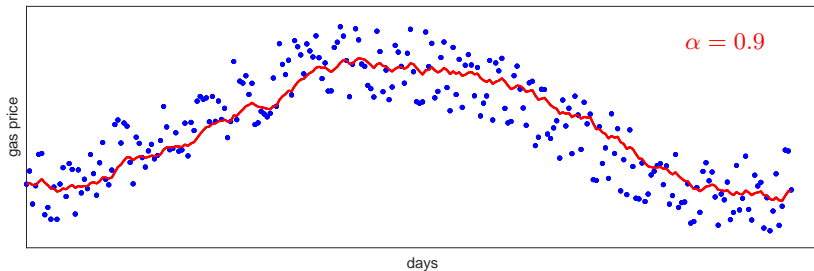
$$v_t = \begin{cases} g_1 & t = 1 \\ \alpha v_{t-1} + (1 - \alpha)g_t & t > 1 \end{cases}$$

- ▶ v_t : EWMA at time t
- ▶ g_t : observation at t
- ▶ $\alpha \in [0, 1]$: degree of weighting decrease (constant _____ factor)

- example: gas price over time

- ▶ blue dot: gas price g
- ▶ red curve: EWMA v

$$v_t = \alpha \cdot \underbrace{v_{t-1}}_{\text{previous EWMA}} + (1 - \alpha) \cdot \underbrace{g_t}_{\text{current observation}}$$



Properties of EWMA

- effective weighting **decreases exponentially** over time:

$$\begin{aligned}v_t &= \alpha v_{t-1} + (1 - \alpha)g_t \\&= \alpha [\alpha v_{t-2} + (1 - \alpha)g_{t-1}] + (1 - \alpha)g_t \\&\vdots \\&= \alpha^k v_{t-k} + (1 - \alpha) \underbrace{\left[g_t + \alpha g_{t-1} + \alpha^2 g_{t-2} + \cdots + \alpha^{k-1} g_{t-k+1} \right]}_{\text{weight exponentially decreases toward the past}}\end{aligned}$$

⇒ thus called “_____ weighted”

- approximation²

$$\begin{aligned}v_t &= (1 - \alpha)g_t + \alpha v_{t-1} \\&= (1 - \alpha) \left[g_t + \alpha g_{t-1} + \alpha^2 g_{t-2} + \alpha^3 g_{t-3} + \cdots \right] \\&= \frac{g_t + \alpha g_{t-1} + \alpha^2 g_{t-2} + \cdots}{1 + \alpha + \alpha^2 + \cdots} \Rightarrow \text{weighted average formula}\end{aligned}$$

- ▶ in such a formula, **denominator** = effective number of observations:

$$1 + \alpha + \alpha^2 + \cdots = \frac{1}{1 - \alpha}$$

- ▶ bottom line:

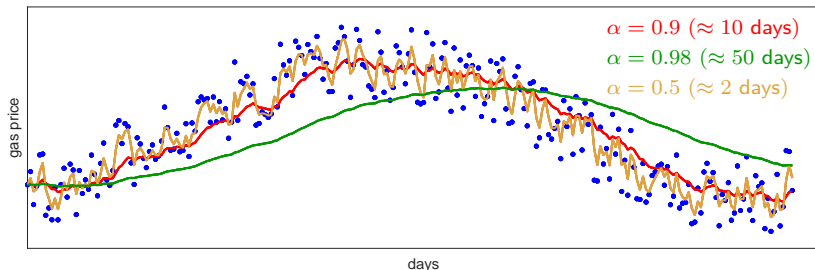
$v_t \approx \text{average over}$	last time points
-----------------------------------	---------------------------

e.g. $\alpha = 0.9 \Rightarrow \text{average over } 1/(1 - 0.9) = 10 \text{ points}$
 $\alpha = 0.98, 0.5 \Rightarrow \text{average over } 50, 2 \text{ points, respectively}$

²recall: $\frac{1}{1 - \alpha} = 1 + \alpha + \alpha^2 + \cdots$

Effect of smoothing factor α

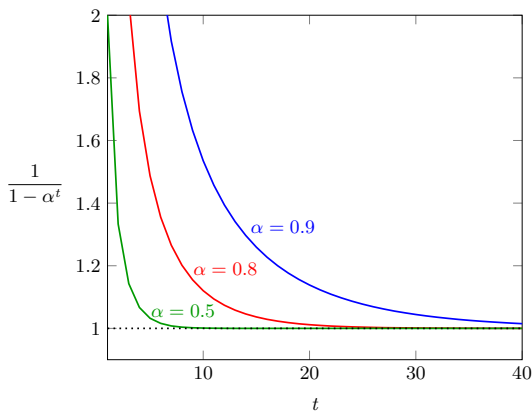
- higher α (= more weight to past, less weight to present)
 - ▶ smoother curve ← averaging over more days
 - ▶ shifted further ← averaging over larger window
 - ⇒ curve adapts more _____ to changes with more latency



Bias correction

- first few iterations: inaccurate average (have not seen enough samples)
 - ▶ instead of v_t , we thus use:

$$\frac{v_t}{1 - \alpha^t}$$



Outline

Introduction

Gradient-Based Optimization

Gradient Descent and its Limitations

Exponentially Weighted Average

Gradient Descent with Momentum

Per-Parameter Adaptive Learning Rates

Additional Topics

Summary

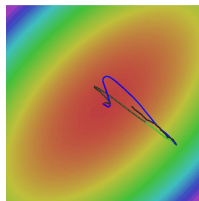
Method of momentum

- **sgd**: very popular but sometimes slow
- **method of momentum** (Polyak, 1964):
 - ▶ designed to accelerate learning, especially in the face of
 - ▷ high curvature
 - ▷ small but consistent gradients
 - ▷ noisy gradients
 - ▶ can be combined to existing sgd variants
- **common algorithm**
 - ▶ accumulates exponentially decaying moving average of past _____
 - ▶ then continues to move in their direction

Gradient descent with momentum

- idea: compute EWMA of gradients and use it to update weights
 - ▶ works almost always faster than standard gradient descent
- in physics
 - ▶ **momentum** = mass · velocity
 - ▶ for unit mass: momentum = _____
 - ▶ smoothing factor α : **friction**
- sometimes
 - ▶ smoothing factor α

⇒ called momentum (misnomer)



— sgd
— sgd + momentum
— Nesterov

(source: cs231n)

Three (equivalent) forms of sgd + momentum

- let $g \triangleq \nabla_{\theta} J(\theta)$ (θ represents W and b altogether)

form 1	$\begin{aligned} v &\leftarrow \alpha v - \epsilon g \\ \theta &\leftarrow \theta + v \end{aligned}$	\Rightarrow	$\theta \leftarrow \theta - \epsilon \left(g + \frac{\alpha}{-\epsilon} v \right)$
--------	--	---------------	---

form 2	$\begin{aligned} v &\leftarrow \alpha v + (1 - \alpha) g \\ \theta &\leftarrow \theta - \epsilon v \end{aligned}$	\Rightarrow^3	$\theta \leftarrow \theta - \tilde{\epsilon} \left(g + \frac{\alpha}{1 - \alpha} v \right)$
--------	---	-----------------	--

form 3	$\begin{aligned} v &\leftarrow \alpha v + g \\ \theta &\leftarrow \theta - \epsilon v \end{aligned}$	\Rightarrow	$\theta \leftarrow \theta - \epsilon (g + \alpha v)$
--------	--	---------------	--

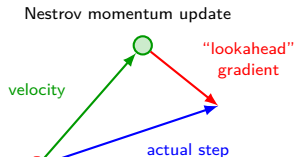
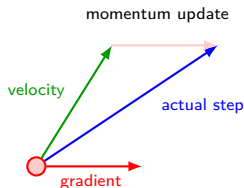
- bottom line: θ is updated by linear _____ of **gradient** and **velocity**

$$\theta \leftarrow \theta - \epsilon \left(\underbrace{g}_{\text{gradient}} + \text{constant} \cdot \underbrace{v}_{\text{velocity}} \right)$$

$^3 \tilde{\epsilon} \triangleq \epsilon(1 - \alpha)$

Nesterov momentum

- difference from standard momentum:
 - ▶ _____ gradient $g = \nabla_{\theta} J$ is evaluated



standard momentum

g evaluated at current position θ
(red circle)

Nesterov momentum

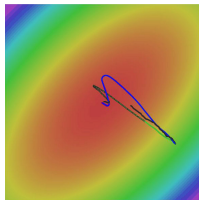
g evaluated at "lookahead" position
 $\theta + \alpha v$ (green circle)

- rationale: momentum is about to carry us to a new position
 - ▶ make sense to evaluate g at new position instead of "old/stale" position

- Nesterov momentum update rule:

$$\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta} + \alpha \mathbf{v})$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mathbf{v}$$



— sgd
— sgd + momentum
— Nesterov

(source: cs231n)

- Nesterov momentum
 - ▶ gradient is evaluated after the current velocity is applied
 - ⇒ interpreted as adding a _____ factor to standard momentum
- advantages
 - ▶ stronger theoretical converge guarantees for convex functions
 - ▶ consistently works slightly better than standard momentum

Outline

Introduction

Gradient-Based Optimization

- Gradient Descent and its Limitations

- Exponentially Weighted Average

- Gradient Descent with Momentum

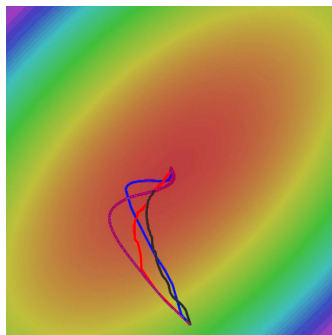
- Per-Parameter Adaptive Learning Rates

Additional Topics

Summary

Per-parameter adaptive learning rates

- optimization methods explained so far
 - ▶ set learning rate ϵ globally and equally for all parameters
- methods presented now: AdaGrad, RMSProp, Adam
 - ▶ adaptively tune ϵ for ____ parameter

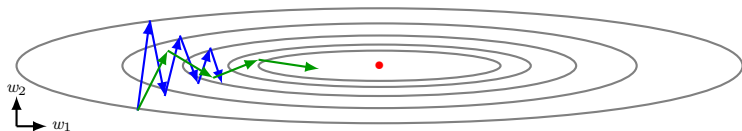


- sgd
- sgd + momentum
- RMSProp
- Adam

(source: cs231n)

Motivation

- recall: limitation of gradient descent



- ▶ **goal**: move horizontally
- ▶ **problem**: huge vertical oscillations
- **solution**: we want to
 - ▶ slow down learning vertically
 - ▶ speed up (or at least not slow down) learning horizontally
- how to implement this idea (without relying on __ explicitly)?

AdaGrad

- individually adapts learning rate of each direction (*i.e.* each parameter)
 - ▶ **steep** direction (large $\frac{\partial J}{\partial \theta_j}$): **slow down** learning
 - ▶ **gently sloped** direction (small $\frac{\partial J}{\partial \theta_j}$): **speed up** learning
- adjusts learning rates per parameter:

$$\epsilon_j = \frac{\epsilon}{\sqrt{\sum_{\text{all previous iterations}} (g_j \cdot g_j)}}$$

- ▶ ϵ : global learning rate
 - ▶ ϵ_j : learning rate of dimension j (parameter θ_j)
 - ▶ $g_j = \frac{\partial J(\theta)}{\partial \theta_j}$: gradient wrt dimension j
- net effect:
 - ▶ **greater progress** in more **sloped** directions

- downside (esp in deep learning)
 - ▶ monotonically decreasing ϵ : too aggressive
 - ⇒ stops learning **too** _____
- TF: AdagradOptimizer
 - ▶ but do not use it for neural nets
- **Adadelta**: an extension of Adagrad
 - ▶ restricts the window of accumulated past gradients to some **fixed size**
 - ⇒ reduces aggressive, monotonically decreasing learning rate

RMSProp (root-mean-square prop)

- modifies AdaGrad to perform better in non-convex setting
 - ▶ changes gradient accumulation to **EWMA**
- use of **exponentially decaying average** allows RMSprop to
 - ▶ discard history from _____ past
 - ⇒ converge rapidly after finding a convex bowl
- comparison (r : accumulation variable)

AdaGrad	RMSprop
$r \leftarrow r + g \odot g$	$r \leftarrow \rho r + (1 - \rho) g \odot g$
$\Delta \theta \leftarrow -\frac{\epsilon}{\sqrt{\delta + r}} \odot g$	$\Delta \theta \leftarrow -\frac{\epsilon}{\sqrt{\delta + r}} \odot g$
$\theta \leftarrow \theta + \Delta \theta$	$\theta \leftarrow \theta + \Delta \theta$

- ▶ **decay rate** ρ : hyperparameter (typically 0.9, 0.99, 0.999)

Adam (adaptive moment estimation)

- idea: **RMSPprop** + **momentum**

- ▶ with ____ correction

- for each iteration:

- ① compute gradient \mathbf{g}

- ② update first moment: $\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}$

← “momentum”

- ③ update second moment: $\mathbf{r} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$

← “RMSPprop”

- ④ bias correction:

$$\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}, \quad \hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$$

- ⑤ update parameter:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \epsilon \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \delta}}$$

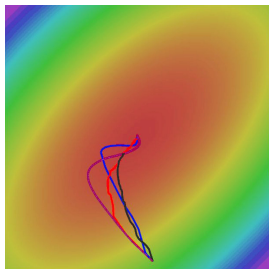
Algorithm 2 Adam optimizer

Require:

- ▶ step size ϵ
- ▶ exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1)$
- ▶ small constant δ used for numerical stabilization
- ▶ initial parameters θ

- 1: initialize 1st and 2nd moment variables $s = \mathbf{0}$, $r = \mathbf{0}$
 - 2: initialize time step $t = 0$
 - 3: **while** stopping criterion not met **do**
 - 4: sample a minibatch $\{x^{(1)}, \dots, x^{(m)}\}$ with corresponding targets $y^{(i)}$
 - 5: compute gradient: $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$
 - 6: $t \leftarrow t + 1$
 - 7: update biased first moment estimate: $s \leftarrow \rho_1 s + (1 - \rho_1)g$
 - 8: update biased second moment estimate: $r \leftarrow \rho_2 r + (1 - \rho_2)g \odot g$
 - 9: correct bias in first moment: $\hat{s} \leftarrow \frac{s}{1 - \rho_1^t}$
 - 10: correct bias in second moment: $\hat{r} \leftarrow \frac{r}{1 - \rho_2^t}$
 - 11: compute update: $\Delta\theta = -\epsilon \frac{\hat{s}}{\sqrt{\hat{r} + \delta}}$ (operations applied element-wise)
 - 12: apply update: $\theta \leftarrow \theta + \Delta\theta$
 - 13: **end while**
-

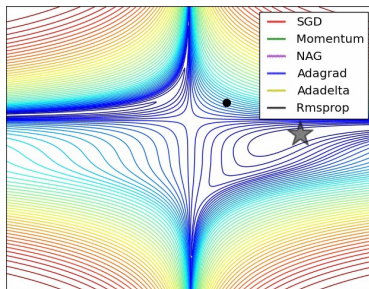
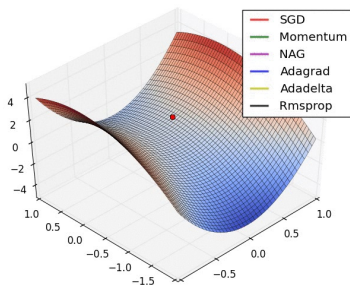
- recommended values in the paper
 - ▶ learning rate ϵ : needs tuning (suggested default: 0.001)
 - ▶ for momentum ρ_1 : 0.9
 - ▶ for RMSProp ρ_2 : 0.999
 - ▶ for stability δ : 10^{-8}
- Adam: often works better than RMSProp
 - ▶ recommended as the default algorithm to use
 - ▶ alternative to Adam worth trying: **sgd + Nesterov momentum**



- sgd
- sgd + momentum
- RMSProp
- Adam

(source: cs231n)

Comparison



(source: Radford)

- more information: [▶ Link](#)

Outline

Introduction

Gradient-Based Optimization

Additional Topics

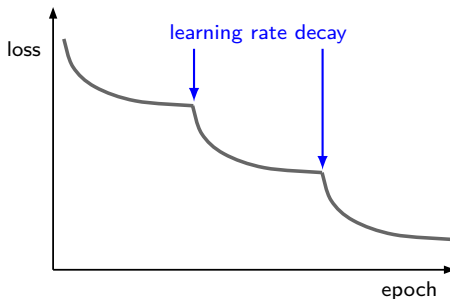
Learning Rate Scheduling

Second-Order Optimization

Summary

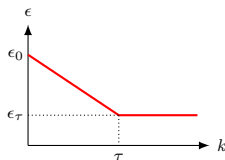
Learning rate

- hyperparameter for many gradient-based algorithms
 - ▶ sgd, sgd + momentum, AdaGrad, RMSProp, Adam
- need to **gradually** _____ learning rate over time
 - ⇒ now denote ϵ_k : learning rate at iteration k (ϵ_0 : initial)
 - ▶ more critical with sgd + momentum (less common with Adam)



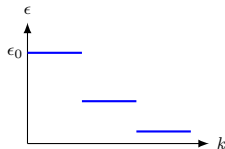
(source: cs231n)

How to decay learning rate

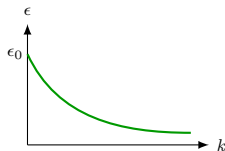


- linear decay (until τ , and then constant)

$$\epsilon_k = \left(1 - \frac{k}{\tau}\right) \epsilon_0 + \frac{k}{\tau} \epsilon_\tau$$



- step decay
 - ▶ discrete staircase



- exponential decay: *e.g.* $\epsilon = \epsilon_0(0.95)^k$
- $1/k$ or $1/\sqrt{k}$ decay:

- also popular: _____ decay (by trial-and-error or monitoring learning curve)

How to set initial learning rate

ϵ_0	if too large:	if too low:
	<ul style="list-style-type: none">• violently oscillating learning curve• cost function often increases significantly	<ul style="list-style-type: none">• learning proceeds slowly• learning may stuck with a high cost value

- typically:

▶ optimal $\epsilon_0 > \underbrace{\epsilon_{\sim 100}^*}_{\uparrow}$

learning rate that yields best performance after first ____ iterations or so

- advice: monitor the first several iterations and
 - ▶ use a learning rate that is
 - ▷ higher than best-performing ϵ at this time
 - ▷ but not so high that it causes severe instability

Outline

Introduction

Gradient-Based Optimization

Additional Topics

Learning Rate Scheduling

Second-Order Optimization

Summary

Idea behind Newton's method

- consider a second-order Taylor series approximation
 - ▶ to function $f(\mathbf{x})$ around the current point $\mathbf{x}^{(0)}$:

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{g} + \frac{1}{2}(\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{H}(\mathbf{x} - \mathbf{x}^{(0)})$$

- ▶ $\mathbf{g} \triangleq \nabla_{\mathbf{x}} f(\mathbf{x}^{(0)})$: gradient of f at $\mathbf{x}^{(0)}$
- ▶ $\mathbf{H} \triangleq \mathbf{H}(f)(\mathbf{x}^{(0)})$: Hessian of f at $\mathbf{x}^{(0)}$

- solving for the critical point of f gives Newton's update rule:

$$\mathbf{x}^* = \mathbf{x}^{(0)} - \mathbf{H}^{-1} \mathbf{g}$$

- ▶ **pros**: (in theory) no hyperparameter (*i.e.* learning rate)
- ▶ **cons**: _____ (\mathbf{H} has $O(n^2)$ elements and takes $O(n^3)$ for inverting)

※ **Levenberg-Marquardt** algorithm

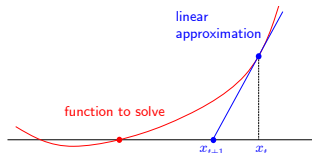
- ▶ switches between Newton's and gradient descent

Comparison (1D)

- Newton's method: second-order

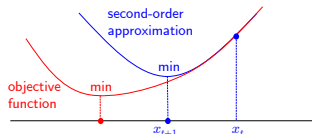
- ▶ zero-finding

$$x_{t+1} = x_t - \frac{f(x_t)}{f'(x_t)}$$



- ▶ minimization/maximization

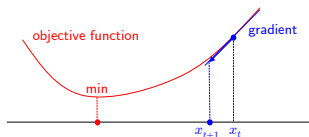
$$x_{t+1} = x_t - \frac{f'(x_t)}{f''(x_t)}$$



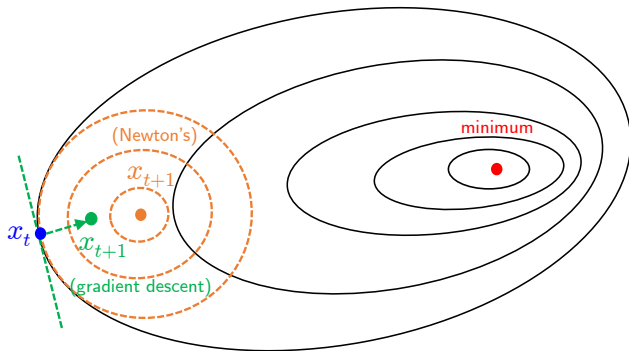
- gradient descent: first-order

- ▶ minimization/maximization

$$x_{t+1} = x_t - \epsilon f'(x_t)$$



Comparison (2D)



- Newton's method for optimization
 - ▶ idea: get a _____-order approximation and minimize it
 - ⇒ faster than gradient descent

Quasi-Newton methods

- idea: avoid directly inverting \mathbf{H}
 - ▶ approximate \mathbf{H}^{-1} with matrix \mathbf{M}_t
 - ▷ \mathbf{M}_t : iteratively refined by low-rank updates
 - ▶ determine direction of descent by $\boldsymbol{\rho}_t = \mathbf{M}_t \mathbf{g}_t$ and update:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \epsilon \boldsymbol{\rho}_t$$

- Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm
 - ▶ most popular quasi-Newton method
 - ▶ still requires $O(n^2)$ memory to store \mathbf{H}^{-1}
- L-BFGS (limited memory BFGS): does not form/store full \mathbf{H}^{-1}
 - ▶ usually works very well in full batch/deterministic mode
 - ▶ but performs poorly in _____/stochastic setting (research topic)

Practical advice on choosing optimizer in DL

- Adam
 - ▶ a good default choice in many cases
- **sgd + momentum + learning rate decay**
 - ▶ often outperforms Adam
 - ▶ but requires more _____
- L-BFGS
 - ▶ try it if you can afford to do **full batch** updates
 - ▶ but should disable all sources of noise

(source: cs231n)

Outline

Introduction

Gradient-Based Optimization

Additional Topics

Summary

Summary

- optimization in deep learning
 - ▶ mostly sgd and its variants

- gradient estimate

$$\hat{g} = \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$$

- stochastic gradient descent (sgd)

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

- method of momentum ($\alpha \in [0, 1)$)

$$v \leftarrow \alpha v - \epsilon \hat{g}$$

$$\theta \leftarrow \theta + v$$

- Nesterov momentum (corrected momentum)

$$v \leftarrow \alpha v - \epsilon \nabla_{\theta} \left[\frac{1}{m} \sum_{i=1}^m L(f(x^{(i)}; \theta + \alpha v), y^{(i)}) \right]$$

$$\theta \leftarrow \theta + v$$

- AdaGrad (r: for gradient accumulation)

$$r \leftarrow r + \hat{g} \odot \hat{g}$$

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

$$\theta \leftarrow \theta + \Delta \theta$$

- RMSProp (gradient accumulation by EWMA)

$$r \leftarrow \rho r + (1 - \rho) \hat{g} \odot \hat{g}$$

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

$$\theta \leftarrow \theta + \Delta \theta$$

- Adam (a reasonable default choice)

$$s \leftarrow \rho_1 s + (1 - \rho_1) \hat{g} \quad (\text{momentum})$$

$$r \leftarrow \rho_2 r + (1 - \rho_2) \hat{g} \odot \hat{g} \quad (\text{RMSProp})$$

$$\hat{s} \leftarrow \frac{s}{1 - \rho_1^t}, \quad \hat{r} \leftarrow \frac{r}{1 - \rho_2^t} \quad (\text{bias correction})$$

$$\theta \leftarrow \theta - \epsilon \frac{\hat{s}}{\sqrt{\hat{r} + \delta}}$$