

training nn

stats403_deep_learning
spring_2025
lecture_2

2.1 optimization methods

gradient descent

- consider the optimization task we face in deep learning

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$$

in the context of deep learning,
this is the vector of parameters

- the gradient descent update is

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta^t \nabla f(\mathbf{x}^t)$$

learning rate

gradient descent

- intuition:

$$f(\mathbf{x}^{t+1}) = f(\mathbf{x}^t - \eta^t \nabla f(\mathbf{x}^t))$$

gradient descent

- assume that the eigenvalues of the Hessian of $f(\mathbf{x})$ is uniformly bounded by β : $\lambda(\nabla^2(f(\mathbf{x}))) \leq \beta$ for any \mathbf{x}
- $f(\mathbf{x}^t)$ decreases with t for a proper choice of η^t

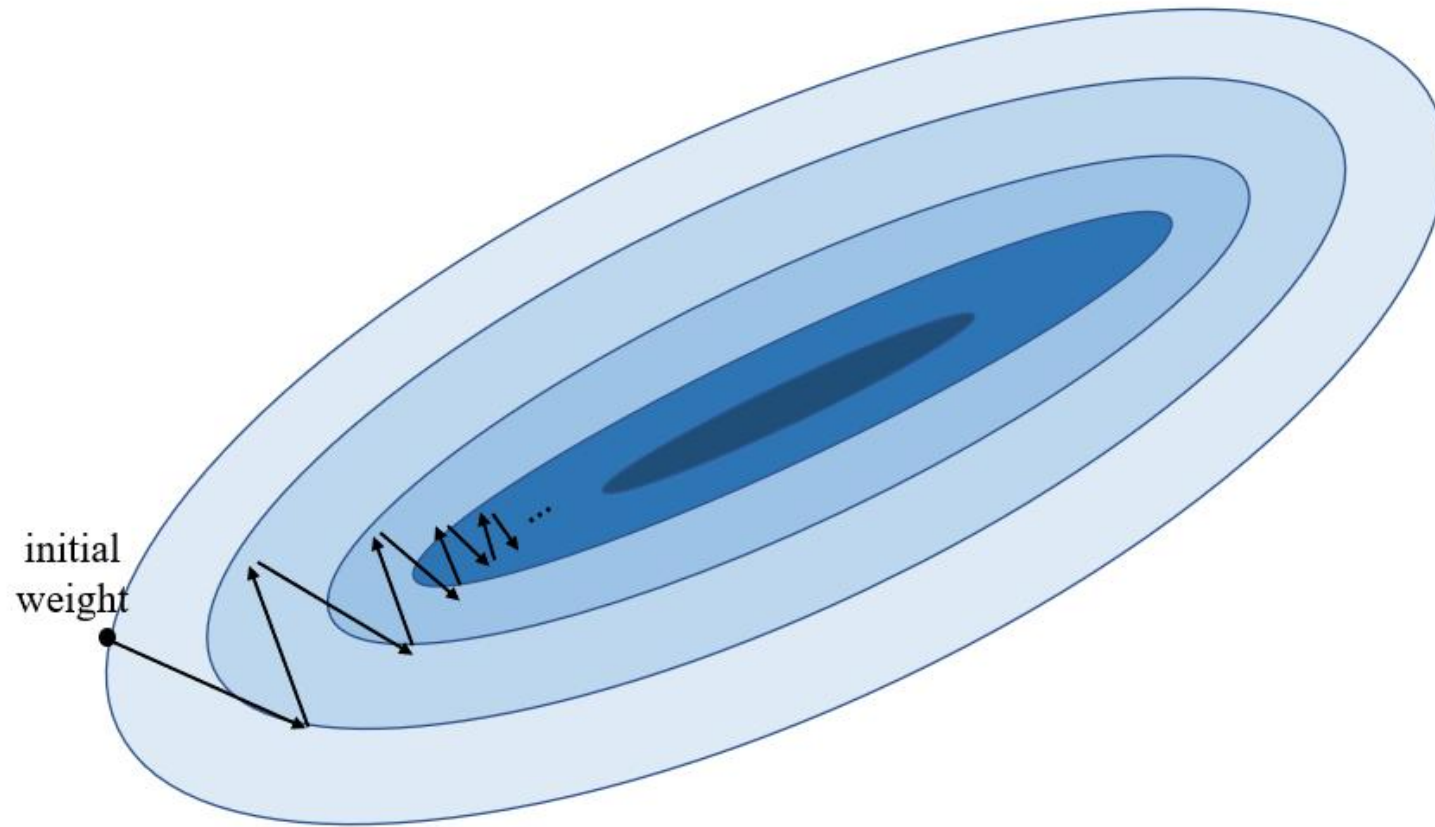
gradient descent

- further assume that f is strongly convex: the eigenvalues of the Hessian of $f(\mathbf{x})$ is uniformly bounded below by $\alpha > 0$:
 $\lambda(\nabla^2(f(\mathbf{x}))) \geq \alpha$ for any \mathbf{x}

- let \mathbf{x}^* be the minimizer of f ,

$$f(\mathbf{x}^T) - f(\mathbf{x}^*) \leq \left(1 - \frac{\alpha}{\beta}\right)^T (f(\mathbf{x}^0) - f(\mathbf{x}^*))$$

problem with gradient descent



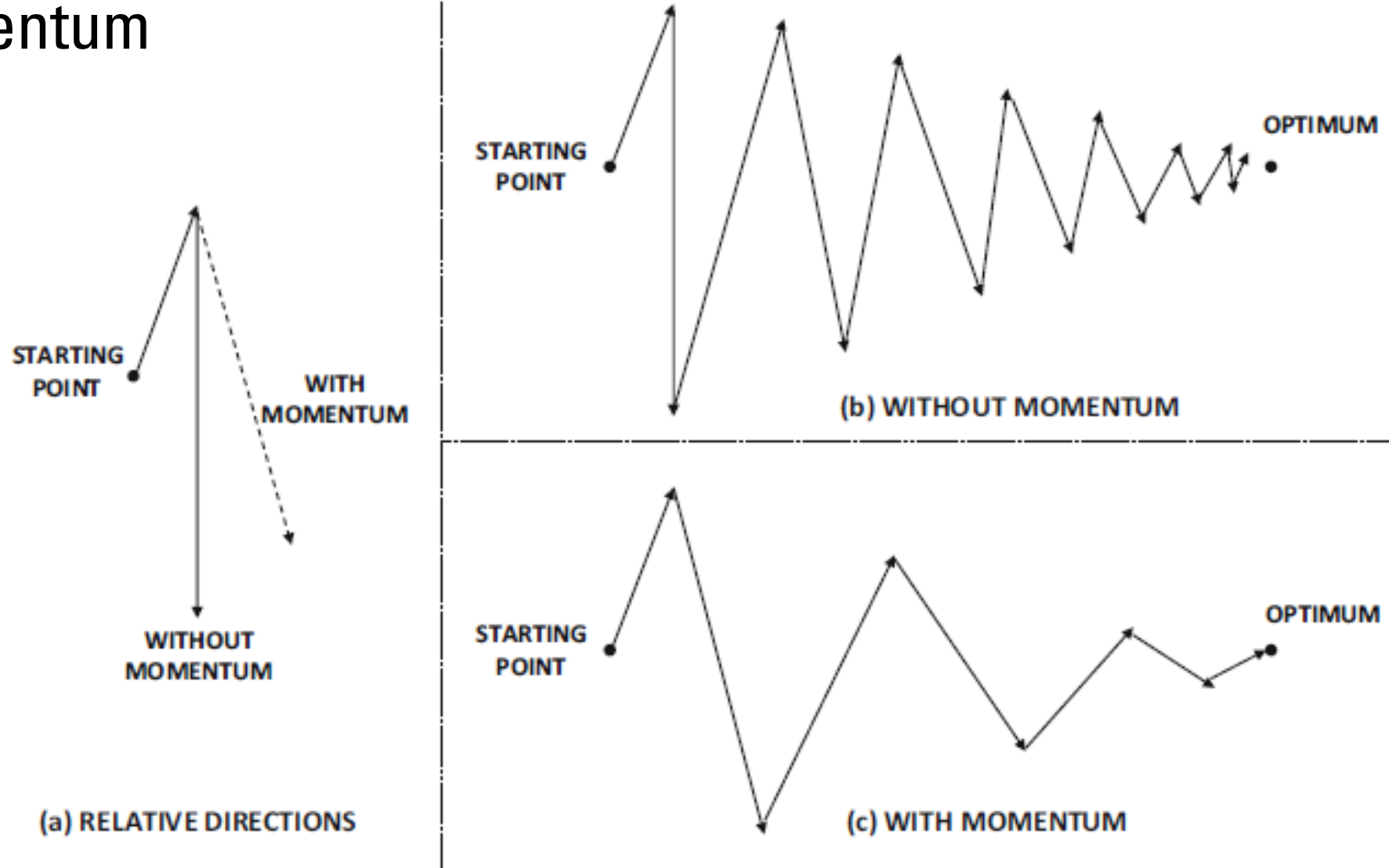
momentum

new direction old direction target direction

$$\mathbf{u}^{t+1} = \mathbf{u}^t - \beta^t \nabla f(\mathbf{x}^t)$$

$$\mathbf{x}^{t+1} = \mathbf{x}^t + \gamma^t \mathbf{u}^{t+1}$$

momentum



Nesterov's acceleration

$$\mathbf{u}^{t+1} = \mathbf{x}^t - \beta^t \nabla f(\mathbf{x}^t)$$

$$\mathbf{x}^{t+1} = \mathbf{u}^{t+1} + \frac{t}{t+3}(\mathbf{u}^{t+1} - \mathbf{u}^t)$$

Gradient descent

$$x_{k+1} = x_k - \tau \nabla f(x_k)$$

$$\tau \rightarrow 0 \quad \downarrow \quad k\tau \rightarrow t$$

$$\frac{dx(t)}{dt} = -\nabla f(x(t))$$

Nesterov's acceleration

$$\begin{aligned} x_{k+1} &= y_k - \tau \nabla f(y_k) \\ y_{k+1} &= x_{k+1} + \frac{k}{k+3} (x_{k+1} - x_k) \end{aligned}$$

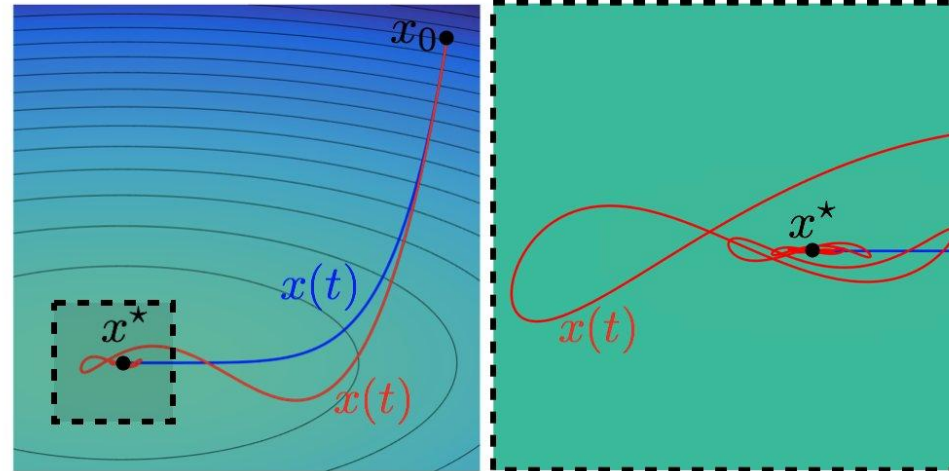
$$\tau \rightarrow 0 \quad \downarrow \quad k\sqrt{\tau} \rightarrow t$$

$$\frac{d^2 x(t)}{dt^2} + \frac{3}{t} \frac{dx(t)}{dt} = -\nabla f(x(t))$$

Theorem:

$$f(x_k) - f(x^*) = O(1/k)$$

$$f(x_k) - f(x^*) = O(1/k^2)$$



Yurii
Nesterov

image credit:



Gabriel Peyré @gabrielpeyre · 23 Dec 2017

Gradient descent is consistent with a first order ODE. Nesterov's acceleration is consistent with a second order ODE. Nesterov is better in the worse case. Gradient is better for strongly convex. [#NoFreeLunch](https://arxiv.org/abs/1503.01243) arxiv.org/abs/1503.01243

preconditioner

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta^t \mathbf{H}_t^{-1} \nabla f(\mathbf{x}^t)$$

preconditioner $\mathbf{x}^{t+1} = \mathbf{x}^t - \eta^t \mathbf{H}_t^{-1} \nabla f(\mathbf{x}^t)$

- a simple quadratic example: $f(\mathbf{x}) = c + \mathbf{b}^T \mathbf{x} + \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x}$
- $\mathbf{x}^t \in \mathbb{R}^d$
- set $\nabla f(\mathbf{x}^{t+1}) = 0$

AdaGrad: adaptive preconditioner

$$\mathbf{G}_t = \left(\sum_{i=1}^t \nabla f(\mathbf{x}^i) \nabla f(\mathbf{x}^i)^{\mathsf{T}} \right)^{\frac{1}{2}}$$

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \mathbf{G}_t^{-1} \nabla f(\mathbf{x}^t)$$

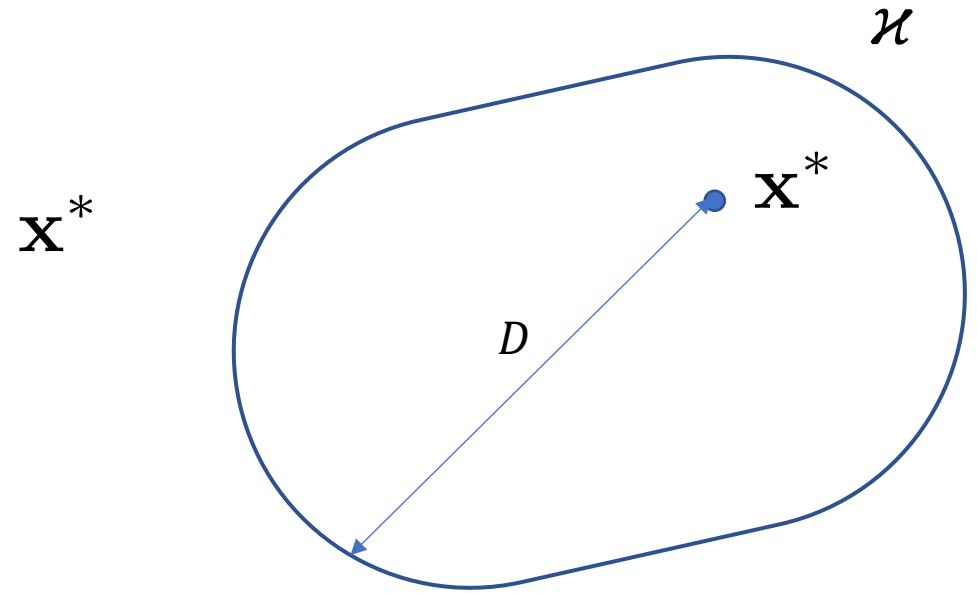
AdaGrad: convergence* (see draft book!)

assume

- f : convex objective function
- \mathcal{K} : bounded convex set for searching
- $D = \max_{\mathbf{x} \in \mathcal{K}} \|\mathbf{x} - \mathbf{x}^*\|$

define $\|\mathbf{x}\|_{\mathbf{A}}^2 = \mathbf{x}^T \mathbf{A} \mathbf{x}$

examining $\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_{\mathbf{G}_t}^2$ yields a convergence bound (see draft book!)



AdaGrad (element-wise version)

$$\mathbf{g}^{(t)} = \left(\sum_{i=0}^t \nabla f(\mathbf{x}^{(i)}) \odot \nabla f(\mathbf{x}^{(i)}) \right)$$
$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta \nabla f(\mathbf{x}^{(t)}) \oslash \sqrt{\mathbf{g}^{(t)}},$$

Algorithm 8.4 The AdaGrad algorithm

Require: Global learning rate ϵ

Require: Initial parameter $\boldsymbol{\theta}$

Require: Small constant δ , perhaps 10^{-7} , for numerical stability

Initialize gradient accumulation variable $\mathbf{r} = \mathbf{0}$

while stopping criterion not met **do**

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

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

 Accumulate squared gradient: $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$.

 Compute update: $\Delta \boldsymbol{\theta} \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$. (Division and square root applied element-wise)

 Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta \boldsymbol{\theta}$.

end while

RMSProp: modified AdaGrad

$$\mathbf{G}^{(t)} = \text{diag} \left(\sum_{i=0}^t \nabla f(\mathbf{x}^{(i)}) \odot \nabla f(\mathbf{x}^{(i)}) \right)$$

$$\mathbf{R}^{(t)} = \alpha \mathbf{R}^{(t-1)} + (1 - \alpha) \mathbf{G}^{(t)}$$

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta (\mathbf{R}^{(t)})^{-1/2} \nabla f(\mathbf{x}^{(t)}).$$

Algorithm 8.5 The RMSProp algorithm

Require: Global learning rate ϵ , decay rate ρ

Require: Initial parameter $\boldsymbol{\theta}$

Require: Small constant δ , usually 10^{-6} , used to stabilize division by small numbers

Initialize accumulation variables $\mathbf{r} = 0$

while stopping criterion not met **do**

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

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

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

 Compute parameter update: $\Delta \boldsymbol{\theta} = -\frac{\epsilon}{\sqrt{\delta + \mathbf{r}}} \odot \mathbf{g}$. ($\frac{1}{\sqrt{\delta + \mathbf{r}}}$ applied element-wise)

 Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta \boldsymbol{\theta}$.

end while

changing the gradient accumulation into an exponentially weighted moving average

Adam

$$\mathbf{g}^{(t)} = \nabla f(\mathbf{x}^{(t)})$$

$$\mathbf{m}^{(t)} = \beta_1 \mathbf{m}^{(t-1)} + (1 - \beta_1) \mathbf{g}^{(t)}$$

$$\mathbf{s}^{(t)} = \beta_2 \mathbf{s}^{(t-1)} + (1 - \beta_2) \mathbf{g}^{(t)} \odot \mathbf{g}^{(t)}$$

$$\hat{\mathbf{m}}^{(t)} = \mathbf{m}^{(t)} / (1 - \beta_1^t)$$

$$\hat{\mathbf{s}}^{(t)} = \mathbf{s}^{(t)} / (1 - \beta_2^t)$$

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta \hat{\mathbf{m}}^{(t)} \oslash \sqrt{\hat{\mathbf{s}}^{(t)}},$$

Adam

Algorithm 8.7 The Adam algorithm

Require: Step size ϵ (Suggested default: 0.001)

Require: Exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1)$.
(Suggested defaults: 0.9 and 0.999 respectively)

Require: Small constant δ used for numerical stabilization (Suggested default: 10^{-8})

Require: Initial parameters θ

Initialize 1st and 2nd moment variables $\mathbf{s} = \mathbf{0}$, $\mathbf{r} = \mathbf{0}$

Initialize time step $t = 0$

while stopping criterion not met **do**

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

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

$t \leftarrow t + 1$

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

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

 Correct bias in first moment: $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}$

 Correct bias in second moment: $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$

 Compute update: $\Delta \theta = -\epsilon \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \delta}}$ (operations applied element-wise)

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

end while

Let's compare the performance of these optimization methods!

- see [colab notebook](#)

SGDW and AdamW

Algorithm 1 SGD with L₂ regularization and SGD with decoupled weight decay (SGDW), both with momentum

- 1: **given** initial learning rate $\alpha \in \mathbb{R}$, momentum factor $\beta_1 \in \mathbb{R}$, weight decay/L₂ regularization factor $\lambda \in \mathbb{R}$
 - 2: **initialize** time step $t \leftarrow 0$, parameter vector $\theta_{t=0} \in \mathbb{R}^n$, first moment vector $m_{t=0} \leftarrow \mathbf{0}$, schedule multiplier $\eta_{t=0} \in \mathbb{R}$
 - 3: **repeat**
 - 4: $t \leftarrow t + 1$
 - 5: $\nabla f_t(\theta_{t-1}) \leftarrow \text{SelectBatch}(\theta_{t-1})$ ▷ select batch and return the corresponding gradient
 - 6: $\mathbf{g}_t \leftarrow \nabla f_t(\theta_{t-1}) + \lambda \theta_{t-1}$
 - 7: $\eta_t \leftarrow \text{SetScheduleMultiplier}(t)$ ▷ can be fixed, decay, be used for warm restarts
 - 8: $\mathbf{m}_t \leftarrow \beta_1 \mathbf{m}_{t-1} + \eta_t \alpha \mathbf{g}_t$
 - 9: $\theta_t \leftarrow \theta_{t-1} - \mathbf{m}_t - \eta_t \lambda \theta_{t-1}$
 - 10: **until** *stopping criterion is met*
 - 11: **return** optimized parameters θ_t
-

SGDW and AdamW

Algorithm 2 Adam with L₂ regularization and Adam with decoupled weight decay (AdamW)

- 1: **given** $\alpha = 0.001, \beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}, \lambda \in \mathbb{R}$
 - 2: **initialize** time step $t \leftarrow 0$, parameter vector $\boldsymbol{\theta}_{t=0} \in \mathbb{R}^n$, first moment vector $\mathbf{m}_{t=0} \leftarrow \mathbf{0}$, second moment vector $\mathbf{v}_{t=0} \leftarrow \mathbf{0}$, schedule multiplier $\eta_{t=0} \in \mathbb{R}$
 - 3: **repeat**
 - 4: $t \leftarrow t + 1$
 - 5: $\nabla f_t(\boldsymbol{\theta}_{t-1}) \leftarrow \text{SelectBatch}(\boldsymbol{\theta}_{t-1})$ ▷ select batch and return the corresponding gradient
 - 6: $\mathbf{g}_t \leftarrow \nabla f_t(\boldsymbol{\theta}_{t-1}) + \lambda \boldsymbol{\theta}_{t-1}$
 - 7: $\mathbf{m}_t \leftarrow \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$ ▷ here and below all operations are element-wise
 - 8: $\mathbf{v}_t \leftarrow \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2$
 - 9: $\hat{\mathbf{m}}_t \leftarrow \mathbf{m}_t / (1 - \beta_1^t)$ ▷ β_1 is taken to the power of t
 - 10: $\hat{\mathbf{v}}_t \leftarrow \mathbf{v}_t / (1 - \beta_2^t)$ ▷ β_2 is taken to the power of t
 - 11: $\eta_t \leftarrow \text{SetScheduleMultiplier}(t)$ ▷ can be fixed, decay, or also be used for warm restarts
 - 12: $\boldsymbol{\theta}_t \leftarrow \boldsymbol{\theta}_{t-1} - \eta_t \left(\alpha \hat{\mathbf{m}}_t / (\sqrt{\hat{\mathbf{v}}_t} + \epsilon) + \lambda \boldsymbol{\theta}_{t-1} \right)$
 - 13: **until** *stopping criterion is met*
 - 14: **return** optimized parameters $\boldsymbol{\theta}_t$
-

initialization

- good initialization is important in iterative optimization methods for both convergence and efficiency
- this is especially true when it comes to training neural networks since the loss landscape is very complicated
- more importantly, since gradients propagate from the output to the input, a good initialization must avoid vanishing or exploding gradients in this process

initialization

- to achieve this, a good heuristic for initializing the weights of a dense network is to keep the variance of each layer the same after applying activation functions
- for instance, one popular initialization method is the Xavier normal initialization, where each weight is initialized randomly according to $\mathcal{N}(0, \sigma^2)$ where $\sigma^2 = \text{gain} \times \frac{2}{\text{\#in} + \text{\#out}}$
- “gain” is a scalar normalization factor which takes into account the activation function used (<https://pytorch.org/docs/stable/nn.init.html>).

Thank you!

Reference

- Ch 6, 7.12, 8.3-8.5, 8.7.1, *Deep Learning*.
- Ch 20.3, *Understanding Machine Learning*.
- Cybenko, G. (1989). Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4), 303-314.
- Arora, S., Lectures 3 & 5, *Towards Theoretical Understanding of Deep Learning (COS597G)*.
- <https://pytorch.org/docs/stable/optim.html>
- <https://github.com/jettify/pytorch-optimizer>
- Défossez, A., Bottou, L., Bach, F., & Usunier, N. (2020). *A Simple Convergence Proof of Adam and Adagrad*.
<https://arxiv.org/pdf/2003.02395.pdf>
- <https://blogs.princeton.edu/imabandit/orf523-the-complexities-of-optimization/>
- Hazan, *Introduction to Online Convex Optimization*.
<https://arxiv.org/pdf/1909.05207.pdf>

