Methods for Stochastic Optimisation

AdaGrad, RMSProp & Adam

Bigi Varghese Philip

Contents

- Why Stochastic?
- Improvement on Gradient Descent
 - GD with Momentum
 - AdaGrad
 - RMS Prop
- ADAM: Adaptive Moment Estimation

Why Stochastic?

- Possible reasons for stochasticity :
 - Inherent function noise
 - Evaluation at random subsamples
 - Stochastic/mini-batch Vs Batch GD
 - Suitable for Online Learning
- $Min E[F(\theta)]$ w.r.to θ
- Perspective Algorithms should average gradients over Mini-batches



• The Problem with GD:

$$\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - \alpha_t \nabla F_{t-1}$$

- Oscillations
- Stuck at flat portions
- Solutions :
 - Adaptive learning rates
 - Perspective Algorithms should average gradients over Mini-batches
 - General Pattern:
 - $\boldsymbol{\theta}_{t} = \boldsymbol{\theta}_{t-1} \alpha G(\nabla F_{t-1} \dots \nabla F_{0}, higher moments)$
- E.g.
 - Gradient Descent With Momentum
 - AdaGrad
 - RMS Prop
 - Adam and Many more
- Why no simple/ Universal soln.?

 $\boldsymbol{\theta} = [b; w]$

b

• The Problem with GD:

 $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - \alpha_t \nabla F_{t-1}$

- Gradient Descent With Momentum:
 - Motivation: to average out the gradient to reduce oscillations and deal with flat regions

$$\boldsymbol{m}_{t} = \beta \boldsymbol{m}_{t-1} + (1 - \beta) \nabla \boldsymbol{F}_{t-1} \\ \boldsymbol{\theta}_{t} = \boldsymbol{\theta}_{t-1} - \alpha \boldsymbol{m}_{t}$$

- Initialization: $m_0=0, <eta=0.9>$



• The Problem with GD:

$$\theta_t = \theta_{t-1} - \alpha_t \nabla F_{t-1}$$

- Gradient Descent With Momentum
- AdaGrad
 - Motivation: Adaptively scaling the gradient $v_t = \sum \nabla F_i^2 \mid i = 0 \text{ to } t - 1 \rightarrow \text{Element wise}$

$$\theta_t = \theta_{t-1} - \frac{\alpha \{ \nabla F_{t-1} \}}{\sqrt{v_t} + \epsilon}$$

Ref: Deep Learning, Andrew Ng (video Lectures)

Duchi, John, Elad Hazan, and Yoram Singer. "Adaptive subgradient methods for online learning and stochastic optimization." Journal of machine learning research 12.7 (2011).

 $\theta = [b; w]$

b

• The Problem with GD:

 $\theta_t = \theta_{t-1} - \alpha \nabla F_{t-1}$

- Gradient Descent With Momentum
- AdaGrad
- RMSProp

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) \nabla F_{t-1}^2$$

$$\theta_t = \theta_{t-1} - \frac{\alpha \{ \nabla F_{t-1} \}}{\sqrt{\nu_t}}$$



Ref: Deep Learning, Andrew Ng (video Lectures)

Lecture 6.5 — Rmsprop: normalize the gradient [Neural Networks for Machine Learning]

- Gradient Descent With Momentum: Simple, Suitable for online algorithm
- AdaGrad: Works well with sparse gradients, popular for linear, less complex models
- RMSProp : Works fine for online and Non-stationary settings

Adam: Adaptive Moment Estimation

Require: α : Stepsize **Require:** $\beta_1, \beta_2 \in [0, 1)$: Exponential decay rates for the moment estimates **Require:** $f(\theta)$: Stochastic objective function with parameters θ **Require:** θ_0 : Initial parameter vector $m_0 \leftarrow 0$ (Initialize 1st moment vector) $v_0 \leftarrow 0$ (Initialize 2nd moment vector) $t \leftarrow 0$ (Initialize timestep) while θ_t not converged do $t \leftarrow t+1$ $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ (Get gradients w.r.t. stochastic objective at timestep t) $m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (Update biased first moment estimate) $v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ (Update biased second raw moment estimate) $\widehat{m}_t \leftarrow m_t/(1-\beta_1^t)$ (Compute bias-corrected first moment estimate) $\hat{v}_t \leftarrow v_t/(1-\beta_2^t)$ (Compute bias-corrected second raw moment estimate) $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \widehat{m}_t / (\sqrt{\widehat{v}_t} + \epsilon)$ (Update parameters) end while **return** θ_t (Resulting parameters)

- Initialisations: $m_0 = \mathbf{0}, v_0 = \mathbf{0}$
- Bias correction to avoid bias towards 0
- ϵ to avoid divide by zero

Adam: Adaptive Moment Estimation

- Recommended Parameters
 - $\begin{array}{l} \beta_1=0.9\\ \beta_2=0.999\\ \epsilon=10^{-8}\\ \alpha \text{ Needs tuning.} \end{array}$
- Update Rule: (Step-size upper bounds) $|\Delta t| = \alpha \cdot \frac{m_t}{\sqrt{v_t}}$ $|\Delta t| \le \alpha \cdot \frac{1-\beta_1}{\sqrt{1-\beta_2}}, \text{ iff } \frac{1-\beta_1}{\sqrt{1-\beta_2}} > 1$ $\le \alpha, Otherwise$

Require: α : Stepsize **Require:** $\beta_1, \beta_2 \in [0, 1)$: Exponential decay rates for the moment estimates **Require:** $f(\theta)$: Stochastic objective function with parameters θ **Require:** θ_0 : Initial parameter vector $m_0 \leftarrow 0$ (Initialize 1st moment vector) $v_0 \leftarrow 0$ (Initialize 2nd moment vector) $t \leftarrow 0$ (Initialize timestep) while θ_t not converged do $t \leftarrow t + 1$ $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ (Get gradients w.r.t. stochastic objective at timestep t) $m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (Update biased first moment estimate) $v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ (Update biased second raw moment estimate) $\widehat{m}_t \leftarrow m_t/(1-\beta_1^t)$ (Compute bias-corrected first moment estimate) $\hat{v}_t \leftarrow v_t/(1-\beta_2^t)$ (Compute bias-corrected second raw moment estimate) $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \widehat{m}_t / (\sqrt{\widehat{v}_t} + \epsilon)$ (Update parameters) end while **return** θ_t (Resulting parameters)

Thank You !!

Let us discuss ...