◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Adagrad, Adam and Online-to-Batch

Raunak Kumar

University of British Columbia

MLRG, 2017 Summer

July 4, 2017

Adagrad	Adam	Online-To-Batch
00000000000	0000000	00

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

Overview

Adagrad

- Motivation
- Algorithm
- Regret Bounds

2 Adam

- Motivation
- Algorithm
- Regret Bounds
- Comparison with Adagrad

3 Online-To-Batch

Brief Overview

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Adagrad

Adagrad	Adam	Online-To-Batch
•••••••		
Motivation		
Introduction		

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Let's begin by motivating Adagrad from 2 different viewpoints:

- Stochastic optimization (brief).
- Online convex optimization.

Stochastic Optimization

- We usually deal with sparse feature vectors.
- Infrequently occurring features are highly informative.
- Examples
 - Blue sky vs orange sky.
 - TF-IDF: Word w is important in document d if it occurs frequently in d but not in the entire corpus.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- Standard stochastic gradient algorithms follow a predetermined scheme.
- Ignore the characteristics of the observed data.
- Idea: Use a learning rate dependent on the frequency of the features.

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

- Frequent feature \rightarrow low learning rate.
- Infrequent feature \rightarrow high learning rate.

Adagrad	Adam	Online-To-Batch
00000000000		
Motivation		
FTI		

• In online convex optimization, Follow The Leader (FTL) chooses the next decision as the best one in hindsight

$$x_{t+1} = \arg\min_{x \in \mathcal{K}} \sum_{s=1}^{t} f_s(x).$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

- We've seen that regret $_{T} = O(T)$ because FTL is "unstable".
- Idea: Can we stabilize it by adding a regularizer?

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Regularization Functions

- We will consider regularization functions $R : \mathcal{K} \to \mathbb{R}$.
- R is strongly-convex, smooth and twice-differentiable.
- Thus, $\nabla^2 R(x) \succ 0$.

Adagrad	Adam	Online-To-Batch
00000000000		
Motivation		

Algorithm 1 Regularized Follow The Leader

- 1: Input: Stepsize $\eta >$ 0, regularization function R and a convex, compact set decision set $\mathcal K$
- 2: Let $x_1 = \operatorname{arg\,min}_{x \in \mathcal{K}} \{ R(x) \}.$
- 3: for t = 1 to T do
- 4: Predict x_t , observe f_t and let $\nabla_t = \nabla f_t(x_t)$.
- 5: Update

RFIL

$$x_{t+1} = \operatorname*{arg\,min}_{x \in \mathcal{K}} \{ \sum_{s=1}^{t} \nabla_s^T x + R(x) \}.$$

Adagrad	Adam	Online-To-Batch
00000000000		
Motivation		



Theorem (RFTL Regret)

$$regret_T \leq 2D_R G_R \sqrt{2T}.$$

Can also write

$$regret_{\mathcal{T}} \leq \max_{u \in \mathcal{K}} \sqrt{2\sum_{t} \|
abla_t \|_t^{*2} B_{\mathcal{R}}(u \| x_t)}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Adagrad	Adam	Online-To-Batch
0000000000		
Motivation		
Online Mirror Descent		

Recall OMD from 3 weeks ago:

Algorithm 2 Online Mirror Descent (Agile version)

- 1: Input: Stepsize $\eta > 0$, regularization function R
- 2: Let y_1 s.t. $\nabla R(y_1) = 0$.
- 3: Let $x_1 = \operatorname{arg\,min}_{x \in \mathcal{K}} B_R(x \| y_1)$.
- 4: for t = 1 to T do
- 5: Predict x_t , observe f_t and let $\nabla_t = \nabla f_t(x_t)$.
- 6: Update y_{t+1} as

$$\nabla R(y_{t+1}) = \nabla R(x_t) - \eta \nabla_t.$$

7: Project according to B_R

$$x_{t+1} = \argmin_{x \in \mathcal{K}} B_R(x \| y_{t+1}).$$

Adagrad	Adam	Online-To-Batch
00000000000		
Motivation		
Online Mirror Descent		

• Suppose we choose $R(x) = \frac{1}{2} ||x - x_0||_2^2$ for an arbitrary $x_0 \in \mathcal{K}$. What do we get?

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

- Suppose we choose $R(x) = \frac{1}{2} ||x x_0||_2^2$ for an arbitrary $x_0 \in \mathcal{K}$. What do we get?
- Online Gradient Descent!
- Regret bound

$$regret_{T} \leq \frac{1}{\eta} D_{R}^{2} + 2\eta \sum_{t} \|\nabla_{t}\|_{t}^{*2}$$
$$\leq 2GD\sqrt{T}$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Online Mirror Descent

- $R(x) = \frac{1}{2} ||x x_0||_2^2$.
- $\nabla R(x) = x x_0$.
- Update

$$y_t = x_{t-1} - \eta \nabla_{t-1}$$
$$x_t = \Pi_{\mathcal{K}}(y_t)$$

• Projection with respect to *B_R* (exercises 2, 3)

$$\begin{aligned} x_t &= \operatorname*{arg\,min}_{x \in \mathcal{K}} B_R(x \| y_{t+1}) \\ &= \operatorname*{arg\,min}_{x \in \mathcal{K}} \frac{1}{2} (x - y_{t+1})^T \nabla^2 R(z) (x - y_{t+1}) \\ &= \operatorname*{arg\,min}_{x \in \mathcal{K}} \frac{1}{2} \| x - y_{t+1} \|_2^2 \end{aligned}$$

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Thus, OMD is equivalent to OGD with this choice of R.

Optimal Regularization

What have we seen so far?

- Using a regularization function to make FTL stable.
- OGD is a special case of OMD with $R(x) = \frac{1}{2} ||x x_0||_2^2$.
- The regularization function *R* affects our regret bound.

Question: Which regularization function should we choose to minimize regret?

Optimal Regularization

What have we seen so far?

- Using a regularization function to make FTL stable.
- OGD is a special case of OMD with $R(x) = \frac{1}{2} ||x x_0||_2^2$.
- The regularization function *R* affects our regret bound.

Question: Which regularization function should we choose to minimize regret? Depends on ${\cal K}$ and the cost functions.

So we'll learn the optimal regularizer online!

Adagrad	Adam	Online-To-Batch
0000000000000		
Algorithm		
Preliminaries		

- Goal: Learn a regularizer that adapts to the sequence of cost functions. In some sense, an optimal regularizer to use in hindsight.
- \bullet We will consider all strongly convex regularizers with a fixed and bounded Hessian in ${\cal H}$

$$\forall x \in \mathcal{K}. \ \nabla^2 R(x) = \nabla^2 \in \mathcal{H}$$

where

$$\mathcal{H} = \{X \in \mathbb{R}^{n \times n} : Tr(X) \le 1, X \succeq 0\}$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Adagrad	Adam	Online-To-Batch
000000000000		
Algorithm		
Adagrad		

Algorithm 3 AdaGrad: Adaptive (sub)Gradient Method (Full Matrix Version)

- 1: Input: Stepsize $\eta > 0$, $\delta > 0$, $x_1 \in \mathcal{K}$.
- 2: Initialize: $S_0 = G_0 = 0$.
- 3: for $t=1\$ to $\$ T do
- 4: Predict x_t , observe loss f_t and let $\nabla_t = \nabla f_t(x_t)$.

5: Update

$$S_t = S_{t-1} + \nabla_t \nabla_t^T$$
$$G_t = S_t^{\frac{1}{2}} + \delta I$$
$$y_{t+1} = x_t - \eta G_t^{-1} \nabla_t$$
$$x_{t+1} = \operatorname*{arg min}_{x \in \mathcal{K}} ||x - y_{t+1}||_{G_t}^2$$

Adagrad	Adam	Online-To-Batch
000000000000		
Algorithm		
Adagrad		

Algorithm 4 AdaGrad: Adaptive (sub)Gradient Method (Diagonal Version)

- 1: Input: Stepsize $\eta > 0$, $\delta > 0$, $x_1 \in \mathcal{K}$.
- 2: Initialize: $S_0 = G_0 = 0$.
- 3: Initialize: $S_{1:0} = 0$.
- 4: for t = 1 to T do
- 5: Predict x_t , observe loss f_t and let $\nabla_t = \nabla f_t(x_t)$.
- 6: Update

$$S_{1:t} = [S_{1:t-1}\nabla_t]$$

$$H_{t,i} = \|S_{1:t,i}\|_2$$

$$G_t = diag(H_t) + \delta I$$

$$y_{t+1} = x_t - \eta G_t^{-1}\nabla_t$$

$$x_{t+1} = \operatorname*{arg\,min}_{x \in \mathcal{K}} \|x - y_{t+1}\|_{G_t}^2$$

Adagrad	Adam	Online-To-Batch
0000000000000		
Regret Bounds		
Regret		

We will show the following regret bound

Theorem

Let{ x_t } be the sequence of decisions defined by AdaGrad. Let $D = \max_{u \in \mathcal{K}} ||u - x_1||_2$. Choose $\delta = \frac{1}{2nD}$, $\eta = D$. Then, for any $x^* \in \mathcal{K}$

$$\mathsf{regret}_T(\mathsf{AdaGrad}) \leq 2D \sqrt{\min_{H \in \mathcal{H}} \sum_t \|
abla_t \|_H^{*2} + 1}$$

Adagrad	Adam	Online-To-Bat
00000000000000		
Regret Bounds		

Optmial Regularizer

Proposition

(Exercise 11) Let $A \succeq 0$. Consider the following problem

minimize: $Tr(X^{-1}A)$ subject to: $X \succeq 0$ $Tr(X) \le 1$

Then, minimizer is $X = A^{\frac{1}{2}} / Tr(A^{\frac{1}{2}})$ and the minimum objective value is $Tr^{2}(A^{\frac{1}{2}})$.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Adagrad	Adam	Online-To-Batch
000000000000000000000000000000000000000		
Regret Bounds		
Optmial Regularizer		

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Corollary

$$\sqrt{\min_{H\in\mathcal{H}}\sum_{t} \|\nabla_{t}\|_{H}^{*2}} = Tr(G_{T} - \delta I) = Tr(G_{T}) - \delta n.$$

2 Optimial regularizer: $(G_T - \delta I)/Tr(G_T - \delta I)$.

Adam

Optmial Regularizer

Proof.

$$\begin{split} \|\nabla_t\|_{H}^{*2} &= \nabla_t^T H^{-1} \nabla_t \\ \sum_{t} \|\nabla_t\|_{H}^{*2} &= \sum_{t} \nabla_t^T H^{-1} \nabla_t \\ &= \sum_{t} Tr(\nabla_t^T H^{-1} \nabla_t) \\ &= \sum_{t} Tr(H^{-1} \nabla_t \nabla_t^T) \\ &= Tr(\sum_{t} H^{-1} \nabla_t \nabla_t^T) \\ &= Tr(H^{-1} \sum_{t} \nabla_t \nabla_t^T) \\ &= Tr(H^{-1} (G_T - \delta I)^2 \end{split}$$

Adagrad	Adam
00000000000000	
Regret Bounds	

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Optmial Regularizer

Proof.

Thus, we can rewrite our objective as

$$\min_{H\in\mathcal{H}}\sum_{t}\|\nabla_{t}\|_{H}^{*2}=\min_{H\in\mathcal{H}}Tr(H^{-1}(G_{T}-\delta I)^{2}).$$

Plugging in $A = (G_T - \delta I)^2$ in the previous proposition yields the result.

Adagrad	Adam	Online-To-Batch
00000000000		
Regret Bounds		
Proof of Regret		

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Based on the corollary, it suffices to prove

Lemma

$$regret_T \leq 2DTr(G_T) = 2D\sqrt{\min_{H \in \mathcal{H}} \sum_t \|\nabla_t\|_H^{*2}} + 2nD\delta$$

Adagrad	Adam	Online-To-Batch
000000000000		
Regret Bounds		
Proof of Regret		

Proof.

Using first order definition of convexity,

$$f(x_t) - f(x^*) \leq \nabla_t^T (x_t - x^*)$$

Thus, we can bound the total regret as

$$egin{aligned} \mathsf{regret}_{\mathcal{T}} &= \sum_t f_t(x_t) - f_t(x^*) \ &\leq \sum_t
abla_t^{\mathcal{T}}(x_t - x^*) \end{aligned}$$

Goal: Split this bound further into 2 terms and bound each of them separately.

Adagrad	Adam	Online-To-Batch
00000000000		
Regret Bounds		

Proof of Regret

Proof.

Using definition of y_{t+1} , we get

$$y_{t+1} - x^* = x_t - x^* - \eta G_t^{-1} \nabla_t$$

$$G_t(y_{t+1} - x^*) = G_t(x_t - x^*) - \eta \nabla_t$$

$$(y_{t+1} - x^*)^T G_t(y_{t+1} - x^*) = (x_t - x^*)^T G_t(x_t - x^*)$$

$$- 2\eta \nabla_t^T (x_t - x^*) + \eta^2 \nabla_t^T G_t^{-1} \nabla_t$$

Using Pythagoras' Theorem

$$\|y_{t+1} - x^*\|_{\mathcal{G}_t}^2 \ge \|x_{t+1} - x^*\|_{\mathcal{G}_t}^2$$

in the above, we get

$$\nabla_t^T (x_t - x^*) \leq \frac{\eta}{2} \nabla_t^T G_t^{-1} \nabla_t + \frac{1}{2\eta} (\|x_t - x^*\|_{G_t}^2 - \|x_t - x^*\|_{G_t}^2)$$

Proof of Regret

Proof.

Thus, we have

$$\sum_{t} f_{t}(x_{t}) - f_{t}(x^{*}) \leq \frac{\eta}{2} \sum_{t=1}^{T} \nabla_{t}^{T} G_{t}^{-1} \nabla_{t}$$
$$+ \frac{1}{2\eta} \sum_{t=1}^{T} (x_{t} - x^{*})^{T} (G_{t} - G_{t-1}(x_{t} - x^{*}))$$
$$+ \frac{\delta}{2\eta} D^{2}$$

Now, we proceed by bounding the first 2 terms in the inequality to get our result.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Proof of Regret

Lemma

$$\sum_{t=1}^{T} \nabla_t^T G_t^{-1} \nabla_t \leq 2 \operatorname{Tr}(G_T)$$

Proof.

Proof Sketch: Induction and use that for $A \succeq B \succ 0$,

$$2Tr((A-B)^{\frac{1}{2}}) + Tr(A^{\frac{1}{2}}B) \le 2Tr(A^{\frac{1}{2}})$$

Adagrad	Adam	Online-To-Batch
00000000000		
Regret Bounds		

Proof of Regret

Lemma

$$\sum_{t=1}^{T} (x_t - x^*)^T (G_t - G_{t-1})(x_t - x^*) \leq D^2 \operatorname{Tr}(G_T)$$

Proof.

Proof sketch: Use that for $A \succeq 0$, $\lambda_{\max}(A) \leq Tr(A)$ and linearity of trace.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Proof of Regret

Proof.

FINALLY(!), we get $\sum_{t} f_t(x_t) - f_t(x^*) \leq \sum_{t} \nabla_t^T (x_t - x^*) \qquad (1)$ $\leq \eta \operatorname{Tr}(G_T) + \frac{1}{2\eta} D^2 \operatorname{Tr}(G_T) + \frac{1}{2\eta} \delta D^2 \qquad (2)$ $\leq 2D \operatorname{Tr}(G_T) \qquad (3)$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Adam



• Adagrad accumulates the squared gradients and G_T increases monotonically.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- The learning rate eventually becomes too small.
- Idea: Use a decaying average!
- Examples: Adadelta, RMSProp, etc.

Adagrad	Adam	Online-To-Batch
	000000	
Algorithm		
Adam		

- Adam: Adaptive Moment Estimation.
- It estimates the mean (first moment) and variance (second moment) of the gradients.
- Keeps track of
 - Exponentially decaying average of past squared gradients v_t

- Exponentially decaying average of past gradients m_t

Adagrad	Adam	Online-To-Batch
	000000	
Algorithm		
Adam		

Algorithm 5 Adam: Adaptive Moment Estimation

- 1: Input: Stepsize $\eta > 0$, decay rates β_1 , β_2 , $\epsilon > 0$, and $x_1 \in \mathcal{K}$.
- 2: Intialize: First moment $m_0 = 0$, second moment $v_0 = 0$.
- 3: for t = 1 to T do
- Predict x_t , observe loss f_t and let $\nabla_t = \nabla f_t(x_t)$. 4:
- Update 5:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) \nabla_t$$
$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) \nabla_t^2$$
$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$
$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$
$$x_{t+1} = x_t - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$$

t

Adagrad	Adam	Online-To-Batch
	0000000	
Algorithm		
Bias Correction		

- Let $\nabla_t = \nabla f_t(x_t)$ and assume each ∇_t is a draw from some underlying distribution.
- We estimate its second raw moment as

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) \nabla_t^2 = (1 - \beta_2) \sum_{i=1}^t \beta_2^{t-i} \nabla_i^2$$

 We want to compare estimate E[v_t] with true E[∇²_t] and correct for the mismatch (bias correction).

Adagrad	Adam	Online-To-Batch
	0000000	
Algorithm		
Bias Correction		

$$E[v_t] = E[(1 - \beta_2) \sum_{i=1}^t \beta_2^{t-i} \nabla_i^2]$$

= $E[\nabla_t^2](1 - \beta_2) \sum_{i=1}^t \beta_2^{t-i} + c$
= $E[\nabla_t^2](1 - \beta_2^t) + c$

- If $E[\nabla_i^2]$ is stationary, then c = 0 and we divide by $1 \beta_2^t$ to correct for the bias.
- Otherwise, β_1 should be chosen to assign small weights to gradients too far in the past.

Adagrad	Adam	Online-To-Batch
	0000000	
Regret Bounds		

Regret

Theorem

Assume

• f_t has bounded gradients: $\|f_t(x)\|_2 \leq G$ and $\|f_t(x)\|_{\infty} \leq G_{\infty}$

Distance between x_t generated by Adam is bounded: \(\forall n, m \in [T]]\). \(\|x_n - x_m \|_2 \le D\) and \(\|x_n - x_m \|_\infty \le D_\infty \)
\(\beta_1, \beta_2 \in [0, 1]\) and \(\forall \frac{\beta_1^2}{\sqrt{\beta_2}} < 1\)
\(\alpha_t = \frac{\alpha}{\sqrt{t}}\)
\(\beta_{1,t} = \beta_1 \lambda^{t-1}\) for \(\lambda \in (0, 1)\)
Then \(\forall T > 1\), \(\text{regret}_T = O(\sqrt{T})\)

▲ロト ▲帰 ト ▲ ヨ ト ▲ ヨ ト ・ ヨ ・ の Q ()

Average Regret

Corollary

Assume

$${f 0}$$
 f_t has bounded gradients: $\|f_t(x)\|_2 \leq G$ and $\|f_t(x)\|_\infty \leq G_\infty$

Oistance between x_t generated by Adam is bounded: ∀n, m ∈ [T]. ||x_n - x_m||₂ ≤ D and ||x_n - x_m||_∞ ≤ D_∞

Then

$$\forall T \geq 1, \; rac{regret_T}{T} = O(rac{1}{\sqrt{T}})$$

Comparison with Adagrad

• Choose $\beta_1 = 0$, $\beta_2 \rightarrow 1$ and stepsize $\eta_t = t^{\frac{-1}{2}}$.

• Then,
$$\lim_{\beta_2 \to 1} \hat{v}_t = \frac{1}{t} \sum_{t=1}^T \nabla_t^2$$
.

• Thus, the update is

$$x_{t+1} = x_t - \frac{\eta t^{\frac{-1}{2}}}{\sqrt{t^{-1}\sum_{t=1}^T \nabla_t^T}} \nabla_t$$

• Note that bias correction is important for this comparison.

Adagrad	Adam	Online-To-Batch
	0000000	
Comparison with Adagrad		

Experiments



Figure 1: Logistic regression training negative log likelihood on MNIST images and IMDB movie reviews with 10,000 bag-of-words (BoW) feature vectors.

◆□▶ ◆圖▶ ◆臣▶ ◆臣▶ ─ 臣

- Adagrad performs well in case sparse features and gradients.
- Adam performs well always?

Experiments



Figure 3: Convolutional neural networks training cost. (left) Training cost for the first three epochs. (right) Training cost over 45 epochs. CIFAR-10 with c64-c64-c128-1000 architecture.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- Adam reduces minibatch variance through first moment.
- This is important in the case of CNNs.

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Online-To-Batch

Adagrad	Adam	Online-To-Batch
		•0
Brief Overview		
Introduction		

Given:

- Training set $S = ((x_1, y_1), (x_2, y_2), \dots, (x_T, y_T)).$
- Online learner \mathcal{A} which constructs a sequence of models h_1, h_2, \ldots, h_T .

(日) (日) (日) (日) (日) (日) (日) (日)

Goal: Construct a model h_S based on $\{h_t\}_{t=1}^T$.

Adagrad	Adam	Online-To-Batch
		00
Brief Overview		
Ideas		

Possible conversion strategies:

• Return the final model h_T .



Adagrad	Adam	Online-To-Batch
		00
Brief Overview		
Ideas		

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

- Return the final model h_T .
- Return the longest surviving model.

Adagrad	Adam	Online-To-Batch
		00
Brief Overview		
ldeas		

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

- Return the final model h_T .
- Return the longest surviving model.
- Seturn h_t which performs best on a validation set.

Adagrad	Adam	Online-To-Batch
		00
Brief Overview		
Ideas		

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

- Return the final model h_T .
- Return the longest surviving model.
- Seturn h_t which performs best on a validation set.
- Solution Randomly choose h_t from $\{h_t : t \in [T]\}$.

Adagrad	Adam	Online-To-Batch
		00
Brief Overview		
ldeas		

- Return the final model h_T .
- Return the longest surviving model.
- Seturn h_t which performs best on a validation set.
- Solution Randomly choose h_t from $\{h_t : t \in [T]\}$.
- S Average over the models in some way.
- Outoff-Averaging.

References

Elad Hazan

Introduction to Online Convex Optimization

Foundations and Trends in Optimization, vol 2, no. 3 - 4, pp. 157 - 325, 2015.



Duchi, Hazan, Singer

Adaptive Subgradient Methods for Online Learning and Stochastic Optimization *Journal of Machine Learning Research*, vol 12(Jul), pp. 2121 – 2159, 2011.



Sebastian Ruder

An Overview of Gradient Descent Optimization Algorithms.

arXiv 1609.04747



Kingma, Ba

Adam: A Method For Stochastic Optmization *ICLR*, 2015.



Shivani Agarwal

Online to Batch Conversions

E0: 370 Statistical Learning Theory, 2013.



Shivani Agarwal

From Online to Batch Learning with Cutoff-Averaging

NIPS, 2008.

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

The End