

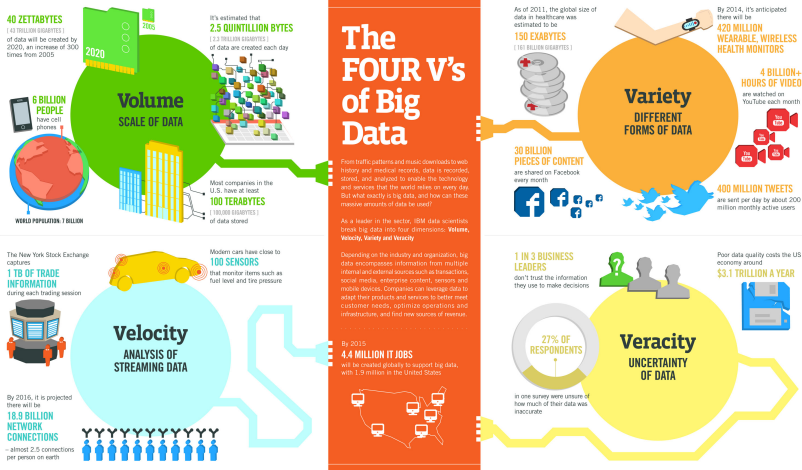
# Online MM Algorithms for Machine Learning

Josh Day

NC State University  
*jtday2@ncsu.edu*

December 21, 2016

## Motivation



Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, HP/ETEC, SAS

IBM

- Statisticians don't have the tools to handle all of this
- Adapting our favorite algorithms is often nontrivial
- Note: If we can handle Velocity, we can handle Volume

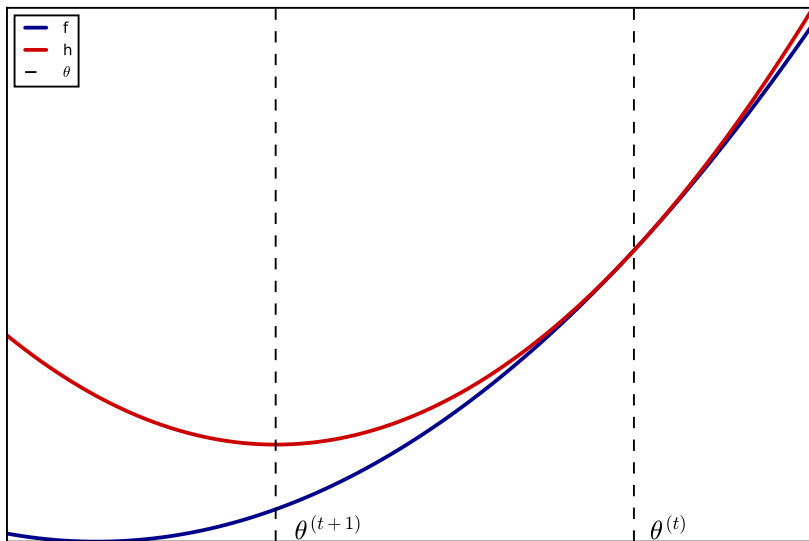
# What is a Majorization Minimization Algorithm?

- $h$  is said to *majorize*  $f$  at  $\theta^{(t)}$  if:

$$\begin{aligned}h(\theta|\theta^{(t)}) &\geq f(\theta|\theta^{(t)}) \\h(\theta^{(t)}|\theta^{(t)}) &= f(\theta^{(t)}|\theta^{(t)})\end{aligned}$$

- MM update:  $\theta^{(t+1)} = \arg \min_{\theta} h(\theta|\theta^{(t)})$

## MM Algorithms



# What is an Online Algorithm?

Velocity can be handled by *online algorithms*

- Input enters piece by piece
- Entire dataset is not available at once
- A well known example: Stochastic Gradient Descent

$$\theta^{(t)} = \theta^{(t-1)} - \gamma_t \nabla f_t(\theta^{(t-1)})$$

- The MM principle is common in offline optimization
  - Coordinate Descent, Proximal Gradient Method, ADMM, ...
- Where does it occur in online optimization?

# Revisit SGD

- SGD minimizes quadratic approximations of noisy functions

$$h_t(\boldsymbol{\theta}) = f_t(\boldsymbol{\theta}^{(t-1)}) + \nabla f_t(\boldsymbol{\theta}^{(t-1)})^T (\boldsymbol{\theta} - \boldsymbol{\theta}^{(t-1)}) + \frac{1}{2\gamma_t} \|\boldsymbol{\theta} - \boldsymbol{\theta}^{(t-1)}\|_2^2$$

$$\begin{aligned}\boldsymbol{\theta}^{(t)} &= \boldsymbol{\theta}^{(t-1)} - \gamma_t \nabla f_t(\boldsymbol{\theta}^{(t-1)}) \\ &= \arg \min_{\boldsymbol{\theta}} h_t(\boldsymbol{\theta})\end{aligned}$$



- $L$ -Lipschitz continuous gradient

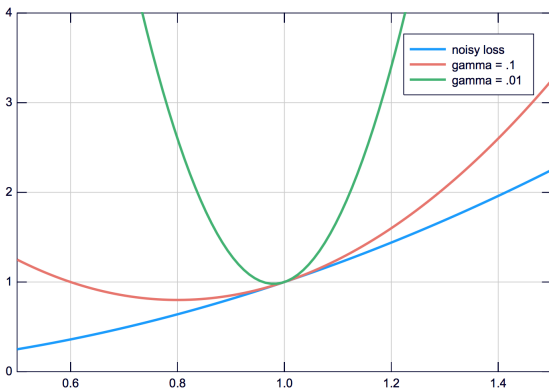
$$f(\boldsymbol{\theta}) \leq f(\boldsymbol{u}) + \nabla f(\boldsymbol{u})^T (\boldsymbol{\theta} - \boldsymbol{u}) + \frac{L}{2} \|\boldsymbol{\theta} - \boldsymbol{u}\|_2^2$$

- i.e., there exists a quadratic upper bound
- RHS is a majorizing function

Suppose  $f_t(\theta)$  has  $L_t$ -Lipschitz continuous gradient

- SGD updates are MM updates anytime  $\gamma_t^{-1} \leq L_t$

- Majorizations get worse over time ( $\gamma_t \rightarrow 0$ ), forcing the parameter to remain near its current state



# Online MM

- Define an Online MM Algorithm as one which follows these steps for each observation/minibatch of data:

- 1 observe  $f_t(\theta)$
- 2 create majorization  $h_t(\theta)$
- 3 Update surrogate objective function  $Q_t(\theta)$  with  $h_t(\theta)$
- 4  $\theta^{(t+1)} = \operatorname{argmin}_{\theta} Q_t(\theta)$

- Algorithms differ at steps 2 (the majorization) and 3 (how  $Q_t$  is updated)
- Ex: SGD
  - $h_t = \text{quadratic upper bound } (m_t\text{-strongly convex, } m_t \rightarrow \infty)$
  - $Q_t(\theta) = h_t(\theta)$

# Online MM Type 1

- Rare case when there exists “sufficient statistics” for majorizing  $\frac{1}{t} \sum_{\tau=1}^t f_{\tau}(\boldsymbol{\theta})$  analytically
- Online majorization = offline majorization
- However, to get same estimate as offline algorithm, we need multiple iterations per update

## Online MM Type 2

- $Q_t(\theta) = (1 - \gamma_t)Q_{t-1}(\theta) + \gamma_t h_t(\theta)$
- Weighted average of surrogate objective and noisy majorization

## Online MM Type 3

- $Q_t(\theta) = h_t(\theta)$ ,  $h_t$  is  $m_t$  strongly convex,  $m_t \rightarrow \infty$
- SGD-like algorithms (ADAGRAD, ADAM)
- Are there non-SGD algorithms?



## Which types can you use?

- Type 1? Only if majorization depends on  $O(1)$  sufficient statistics
- Type 2? Always
- Type 3? Only if majorization can get worse as  $t \rightarrow \infty$

# Can you Make a Majorization "Worse"?

- ✓ Quadratic upper bound

$$h_t(\boldsymbol{\theta}) = f_t(\mathbf{u}) + \nabla f_t(\mathbf{u})^T (\boldsymbol{\theta} - \mathbf{u}) + \frac{L}{2\gamma_t} \|\boldsymbol{\theta} - \mathbf{u}\|_2^2$$

- ✗ Definition of convexity with linear predictor

$$h_t(\mathbf{x}_t^T \boldsymbol{\beta}) = \sum_j \alpha_{tj} f_t \left[ \frac{x_{tj}}{\alpha_{tj}} (\beta_j - \beta_j^{(t-1)}) + \mathbf{x}_t^T \boldsymbol{\beta}^{(t-1)} \right]$$

Online **MM** with **Q**uadratic upper bound using a **D**iagonal Hessian approximation

- If  $\mathbf{H} - d^2f(\boldsymbol{\theta})$  is nonnegative definite, then

$$h(\boldsymbol{\theta}) = f(\mathbf{u}) + \nabla f(\mathbf{u})^T(\boldsymbol{\theta} - \mathbf{u}) + \frac{1}{2}(\boldsymbol{\theta} - \mathbf{u})\mathbf{H}(\boldsymbol{\theta} - \mathbf{u})$$

is a majorizing function

- However, parameters only split if  $\mathbf{H}$  is diagonal

- Type 2 Online MM with Quadratic Upper Bound is

$$\boldsymbol{\theta}^{(t)} = \mathbf{A}_t^{-1} \mathbf{b}_t,$$

where

$$\mathbf{A}_t = (1 - \gamma_t) \mathbf{A}_{t-1} + \gamma_t \mathbf{H}_t$$

$$\mathbf{b}_t = (1 - \gamma_t) \mathbf{b}_{t-1} + \gamma_t [\mathbf{H}_t \boldsymbol{\theta}^{(t-1)} - \nabla f_t(\boldsymbol{\theta}^{(t-1)})]$$

# Applying regularization

For machine learning problems

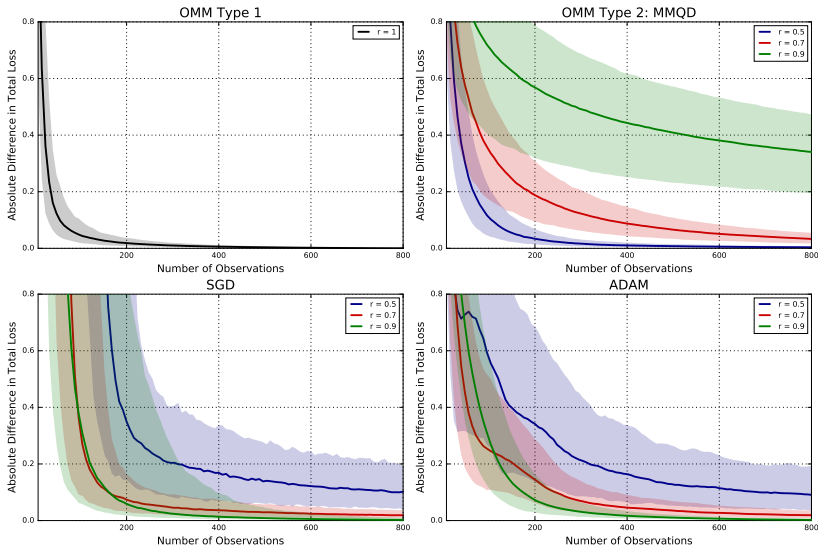
$$\frac{1}{n} \sum f_i(\boldsymbol{\theta}) + g(\boldsymbol{\theta}),$$

the MMQD element-wise updates simplify to

$$\theta_j^{(t)} = \mathbf{prox}_{h_{jj}^{-1}g} \left( \frac{b_{tj}}{h_{jj}} \right)$$

- The following plots display the average loss of various online algorithms (using several different learning rates) relative to the offline solution for 1000 replications of simulated data.
- Solid line: median relative loss
- Ribbons: 0.05 and 0.95 quantiles

## Linear Regression

■ Learning Rate:  $\gamma_t = t^{-r}, r \in [.5, .7, .9]$ 



# Thank You



- Software available:
  - OnlineStats.jl
  - OnlineStatsModels.jl
- Questions?