Linear Time Second Order Stochastic Optimization for Machine Learning

Naman Agarwal Brian Bullins Elad Hazan



ML optimization

• Single Neuron:



- Logistic / Ridge regression, SVMs, Deep Neural Nets
- training set size (m) & dimension of data (d) are very large
- Training time measured in weeks (speech recognition, images...)

First Order Methods



• Sample f_i randomly and take a step according to $-\nabla f_i$

Gradient Descent ++

- Accelerated Gradient Descent (Nesterov)
 - Use momentum from previous steps to take bigger steps
- Adaptive Regularization
 - Adagrad (Duchi, Hazan, Singer), etc...
- Variance Reduction
 - SVRG (Johnson and Zhang; Mahdavi, Zhang, Jin) / SAG, SAGA (Schmidt, LeRoux, Bach)
 - Dual Coordinate Ascent (Shalev-Shwartz & Zhang)
 - Are we at the limit ?
 - Beyond first order necessary?

Higher Order Optimization

- Gradient Descent Direction of Steepest Descent
- Second Order Methods Use Local Curvature





Higher Order Optimization

$$x_{t+1} = x_t - \eta \, [\nabla^2 f(x)]^{-1} \, \nabla f(x)$$

d³ time per iteration! Infeasible for ML!!

Till now 😊



Can the matrix inversion be sped-up??

- Spielman-Teng '04: diagonally dominant systems of equations in linear time!
 - 2015 Godel prize
 - Used by Daitch-Speilman for faster flow algorithms
- Erdogu-Montanari '15: low rank approximation & inversion by Sherman-Morisson
 - Allow stochastic information
 - Still prohibitive: rank * d²
- Concurrently: Agarwal, Langford, Luo: faster, linear-time low-rank sketching (online setting)

Our results

• LiSSA

- Natural Stochastic Newton Method
- Every iteration in O(d) time. Linear in Input Sparsity
- $\log\left(\frac{1}{\epsilon}\right)$ iterations (Linear Convergence)
- Better dependence on the condition number
- better local convergence than known FO methods empirically

• LiSSA ++

- better condition number
- Couple with Matrix Sampling/Sketching techniques Best known running time for $m \gg d$

Overview of Running Times

Algorithm	Running Time
Gradient Descent	$md\kappa\lograc{1}{\epsilon}$
Accelerated Gradient Descent	$md\sqrt{\kappa}\lograc{1}{\epsilon}$
SVRG/SAGA/SDCA	$(m + O(\kappa))d\log\frac{1}{\epsilon}$
LiSSA	$(m + O(\kappa')V) d\log \frac{1}{\epsilon}$
Accelerated SDCA / Catalyst	$\left(m + O\left(\sqrt{\kappa m}\right)\right) d \log \frac{1}{\epsilon}$
LiSSA ++	$O\left(m + O\left(\sqrt{\kappa d}\right)\right) d \log^2 \frac{1}{\epsilon}$

Making Second Order Stochastic

• Easy to get an unbiased estimator of the Hessian

$$\widetilde{\nabla^2 f} = \nabla^2 f_i \qquad i \sim uniform \ [1,m]$$
$$E\left[\widetilde{\nabla^2 f}\right] = \nabla^2 f$$

• But

$$E\left[\widetilde{\nabla^2 f^{-1}}\right] \neq \nabla^2 f^{-1}$$

Nevertheless, used recently in NewSamp (Erdogdu & Montanari)

Our Approach – circumvent inversion!

• The Neumann Series/ Taylor Expansion of the Inverse

• For
$$M \ge 0$$
, $||M|| \le 1$
 $M^{-1} = \sum_{i=0 \text{ to } \infty} (I - M)^i$

• Thus:

$$E\left[\widetilde{\nabla^{2}f^{-1}}\nabla f\right] = \sum_{i} (I - \nabla^{2})^{i}\nabla = E_{k \sim N} \prod_{k=1 \text{ to } q} \left(I - \widetilde{\nabla_{k}^{2}}\right)\nabla$$

For any distribution on naturals $k \sim N$ Vector-vector products only

Improved Estimator

- Previously, Estimate a single term in one estimate
- Recursive Reformulation of the series

$$M_{S_2}^{-1} = I + (I - M)(M + ((I + M))))$$

Ind. Sample
Recursive estimate $M_{S_2-1}^{-tp \ S_2}$

- Truncate after S_2 steps. Typically $S_2 \sim \kappa$
- $E\left[\widetilde{M_{S_2}^{-1}}\right] \to M^{-1} \text{ as } S_2 \to \infty$
- Repeat and average to reduce the variance

Linear time step

• Need to compute the Newton direction

$$\widetilde{M^{-1}} \nabla f = \left(I + \left(I - \widetilde{M} \right) (\dots) \right) \nabla f$$

- Reduces to computing $\widetilde{M}\nabla f$ quickly
- O(d) if loss is of the form $l(w^T x, y)$
 - Hessian is of the form g(x, y) xx^T matrix-vector product → vector-vector product input sparsity time !!

LiSSA –

Linear-time Second-order Stochastic Algorithm

- Use the estimator $\widetilde{\nabla^{-2}f}$ defined previously
- Compute a full (large batch) gradient ∇f
- Move in the direction $\widetilde{\nabla^{-2}f} \nabla f$
- Start with a few Gradient Descent Steps

Main Theorem – (not ++)

Theorem 1

For large t, LiSSA returns a point in the parameter space w_t s.t.

 $f(w_t) \le f(w^*) + \epsilon$

In total time
$$\log\left(\frac{1}{\epsilon}\right) d \left(m + O(\kappa) V\right)$$

- V is a bound on the variance of the estimator
 - In Practice a small constant (e.g. 1)
 - In Theory $V \leq \kappa^2$

Analysis: Newton ~ Iterative Mirrored Descent



Local vs Global Condition Number

 For GD++, condition number is defined by the smoothness and strong convexity parameters

$$\begin{aligned} || \nabla^2 f(x) || &\leq \beta \\ || \nabla^2 f(x) || &\geq \alpha \end{aligned} \qquad Must \ hold \ \forall x \qquad \qquad \kappa_{global} = \frac{\beta}{\alpha} \end{aligned}$$

• For LiSSA, condition number of the local quadratic

$$\kappa_{local} = \max_{x} \kappa(\nabla^2 f(x)) \le \kappa_{global}$$

•In practice – A difference factor of up to 10 between $\kappa_{local} vs \kappa_{global}$

Experimental Results - convex



Comparison b/w Second Order Methods



LiSSA: Alternative Interpretation

- Variance adjusted SGD on the local quadratic approximation
- Quadratic Approximation of f(x) around a point

$$Q(y) = \nabla f(x)^T y + y^T \nabla^2 f(x) y$$

LiSSA ⇔ SGD on a Quadratic

tic

$$\nabla Q(y) = \nabla f(x) + \underbrace{\nabla^2 f(x)}_{estimate} yy$$

$$V - \nabla Q(y) = \left(I - \underbrace{\nabla^2 f(x)}_{estimate}\right) yy - \underbrace{\nabla f(x)}_{Fixed}$$

Acts as Variance Reduction

LiSSA+

- Use FO algorithms (SVRG/SDCA) for intermediate quadratic
- Need to solve each quadratic upto error $\boldsymbol{\epsilon}$
- Total of $\log \log \frac{1}{\epsilon}$ quadratic sub problems
- Total time $Time(FO) \log \log \frac{1}{\epsilon}$

- Better conditioned than the original
- Faster running time in practice

LiSSA ++ (Motivation)

- Why run FO on the quadratic approximation when full f(x) is available ?
- Better Conditioning
- Can use special quadratic structure of the problem
- Make use of Matrix Sampling and Sketching Techniques



- $O(d \log d)$
- Can use FO on those systems

LiSSA++ Theorem

Theorem 2

For large t, LiSSA++ returns a point in the parameter space w_t s.t.

 $f(w_t) \le f(w^*) + \epsilon$

In total time
$$\tilde{O}\left(\log^2\left(\frac{1}{\epsilon}\right) d\left(m + \sqrt{\kappa d}\right)\right)$$

• Fastest over all running time for when $\kappa, m \gg d$

2nd order information: new phenomena?

- escape saddle points?
- "Computational lens for deep nets": experiment with 2nd order information:
 - Trust region
 - Cubic regularization
 - Eigenvalue methods....



Follow-up work (w. Naman, Brian & Tengyu)

- Apply to train deep nets non-convex optimization
- Can implement variants that are reasonable (trust region etc.) using LiSSA techniques
- Theory works (w. the Pearlmuter trick)
- practice doesn't improve upon GD++ 😕

Bottou "say what doesn't work"



Summary

- LiSSA
 - A natural, practical and efficient Second Order Stochastic Algorithm
 - Empirically faster than GD++ on real world datasets
- LiSSA++
 - Fastest running time known theoretically
 - Better dependence on the condition number
- 2nd order information doesn't seem to improve performance for deep nets (so far)

