# Stochastic Newton and quasi-Newton Methods for Large-Scale Convex Optimization

Donald Goldfarb

Department of Industrial Engineering and Operations Research
Columbia University

Joint with Robert Gower and Peter Richtárik

Optimization Without Borders 2016
Les Houches, February 7-12, 2016

- Newton-like and quasi-Newton methods for convex stochastic optimization problems using limited memory block BFGS updates.
- In the class of problems of interest, the objective functions can be expressed as the sum of a huge number of functions of an extremely large number of variables.
- We present preliminary numerical results on problems from machine learning.

# Related work on L-BFGS for Stochastic Optimization

P1 N.N. Schraudolph, J. Yu and S.Günter. A stochastic quasi-Newton method for online convex optim. Int'l. Conf. AI & Stat., 2007

P2 A. Bordes, L. Bottou and P. Gallinari. SGD-QN: Careful quasi-Newton stochastic gradient descent. JMLR vol. 10, 2009

P3 R.H. Byrd, S.L. Hansen, J. Nocedal, and Y. Singer. A stochastic quasi-Newton method for large-scale optim. arXiv1401.7020v2, 2014

P4 A. Mokhtari and A. Ribeiro. RES: Regularized stochastic BFGS algorithm. IEEE Trans. Signal Process., no. 10, 2014.

P5 A. Mokhtari and A. Ribeiro. Global convergence of online limited memory BFGS. to appear in J. Mach. Learn. Res., 2015.

P6 P. Moritz, R. Nishihara, M.I. Jordan. A linearly-convergent stochastic L-BFGS Algorithm, 2015 arXiv:1508.02087v1

P7 X. Wang, S. Ma, D. Goldfarb and W. Liu. Stochastic quasi-Newton methods for nonconvex stochastic optim. 2015, submitted.

(the first 6 papers are for strongly convex problems, the last one is for nonconvex problems)

# Stochastic optimization

- Stochastic optimization

$$\min \ f(x) = \mathbb{E}[f(x, \xi)], \quad \xi \text{ is random variable}$$

- Or finite sum (with $f_i(x) \equiv f(x, \xi_i)$ for $i = 1, \ldots, n$ and very large $n$)

$$\min \ f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x)$$

- $f$ and $\nabla f$ are very expensive to evaluate; e.g., SGD methods randomly choose a random subset $\mathcal{S} \subset [n]$ and evaluate

$$f_{\mathcal{S}}(x) = \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} f_i(x) \quad \text{and} \quad \nabla f_{\mathcal{S}}(x) = \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \nabla f_i(x)$$

- Essentially, only noisy info about $f$, $\nabla f$ and $\nabla^2 f$ is available
- Challenge: how to design a method that takes advantage of noisy 2nd-order information?

# Using 2nd-order information

- Assumption: $f(x) = \frac{1}{n}\sum_{i=1}^{n} f_i(x)$ is strongly convex and twice continuously differentiable.

- Choose (compute) a sketching matrix $S_k$ (the columns of $S_k$ are a set of directions).

- Following Byrd, Hansen, Nocedal and Singer, we do not use differences in noisy gradients to estimate curvature, but rather compute the action of the sub-sampled Hessian on $S_k$. i.e.,

- compute $Y_k = \frac{1}{|\mathcal{T}|}\sum_{i\in\mathcal{T}} \nabla^2 f_i(x)S_k$, where $\mathcal{T} \subset [n]$.

- We choose $\mathcal{T} = \mathcal{S}$

# block BFGS

Given $H_k = B_k^{-1}$, the block BFGS method computes a "least change" update to the current approximation $H_k$ to the inverse Hessian matrix $\nabla^2 f(x)$ at the current point $x$, by solving

$$\begin{aligned} \min \quad & \|H - H_k\| \\ \text{s.t.,} \quad & H = H^\top, \quad HY_k = S_k. \end{aligned}$$

This gives the updating formula (analgous to the updates derived by Broyden, Fletcher, Goldfarb and Shanno).

$$H_{k+1} = (I - S_k[S_k^\top Y_k]^{-1} Y_k^\top) H_k (I - Y_k[S_k^\top Y_k]^{-1} S_k^\top) + S_k[S_k^\top Y_k]^{-1} S_k^\top$$

or, by the Sherman-Morrison-Woodbury formula:

$$B_{k+1} = B_k - B_k S_k[S_k^\top B_k S_k]^{-1} S_k^\top B_k + Y_k[S_k^\top Y_k]^{-1} Y_k^\top$$

# Limited Memory Block BFGS

After $M$ block BFGS steps starting from $H_{k+1-M}$, one can express $H_{k+1}$ as

$$
\begin{aligned}
H_{k+1} &= V_k H_k V_k^T + S_k \Lambda_k S_k^T \\
&= V_k V_{k-1} H_{k-1} V_{k-1}^T V_k + V_k S_{k-1} \Lambda_{k-1} S_{k-1}^T V_k^T + S_k \Lambda_k S_k^T \\
&\vdots \\
&= V_{k:k+1-M} H_{k+1-M} V_{k:k+1-M}^T + \sum_{i=k}^{k+1-M} V_{k:i+1} S_i \Lambda_i S_i^T V_{k:i+1}^T,
\end{aligned}
$$

where

$$
V_k = (I - S_k \Lambda_k Y_k^T) \tag{1}
$$

and $\Lambda_k = (S_k^T Y_k)^{-1}$ and $V_{k:i} = V_k \cdots V_i$.

# Limited Memory Block BFGS

- Hence, when the number of variables $d$ is large, instead of storing the $d \times d$ matrix $H_k$, we store the previous $M$ block curvature pairs

$$\left(S_{k+1-M}, Y_{k+1-M}\right), \ldots, \left(S_k, Y_k\right),$$

  and the Cholesky factors of the matrices $(S_i^T Y_i) = \Lambda_i^{-1}$ for $i = k + 1 - M, \ldots, k$.

- Then, analogously to the standard L-BFGS method, for any vector $v \in \mathbb{R}^d$, $H_k v$ can be computed efficiently using a two-loop block recursion (in $Mp(4d + 2p) + O(p)$) operations), if all $S_i \in \mathbb{R}^{d \times p}$.

Intuition

- Limited memory - least change aspect of BFGS is important
- Each block update acts like a sketching procedure.

# Choices for the Sketching Matrix $S_k$

We employ one of the following strategies

- Gaussian: $S_k \sim \mathcal{N}(0, I)$ has Gaussian entries sampled i.i.d at each iteration.
- Previous search directions $s_i$ delayed: Store the previous $L$ search directions $S_k = [s_{k+1-L}, \ldots, s_k]$ then update $H_k$ only once every $L$ iterations.
- Self-conditioning: Sample the columns of the Cholesky factors $L_k$ of $H_k$ (i.e., $L_k L_k^T = H_k$) uniformly at random. Fortunately we can maintain and update $L_k$ efficiently with limited memory.

The matrix $S$ is a sketching matrix, in the sense that we are sketching the, possibly very large equation $\nabla^2 f(x) H = I$ to which the solution is the inverse Hessian. Left multiplying by $S^T$ compresses/sketches the equation yielding $S^T \nabla^2 f(x) H = S^T$.

# Stochastic Variance Reduced Gradients

- Stochastic methods converge slowly near the optimum due to the variance of the gradient estimates $\nabla f_{\mathcal{S}}(x)$; hence requiring a decreasing step size.
- We use the control variates approach of Johnson and Zhang (2013) for a SGD method SVRG.
- It uses $\nabla f_{\mathcal{S}}(x_t) - \nabla f_{\mathcal{S}}(w_k) + \nabla f(w_k$, where $w_k$ is a reference point, in place of $\nabla f_{\mathcal{S}}(x_t)$ .
- $w_k$, and the full gradient, are computed after each full pass of the data, hence doubling the work of computing stochastic gradients.
- Other recently proposed SGD variance reduction techniques such as SAG, SAGA, SDCA, and S2GD, can be used in place of SVRG.

# The Basic Algorithm

**Algorithm 0.1:** Stochastic Variable Metric Learning with SVRG

**Input**: $H_{-1} \in \mathbb{R}^{d \times d}$, $w_0 \in \mathbb{R}^d$, $\eta \in \mathbb{R}_+$, $s =$ subsample size, $q =$ sample action size and $m$

1 **for** $k = 0, \ldots, max\_iter$ **do**
2      $\mu = \nabla f(w_k)$
3      $x_0 = w_k$
4      **for** $t = 0, \ldots, m-1$ **do**
5          Sample $\mathcal{S}_t, \mathcal{T}_t \subseteq [n]$ i.i.d from a distribution $\mathcal{S}$
6          Compute the sketching matrix $S_t \in \mathbb{R}^{d \times q}$
7          Compute $\nabla^2 f_{\mathcal{S}}(x_t) S_t$
8          $H_t =$ update_metric$(H_{t-1}, S_t, \nabla^2 f_{\mathcal{T}}(x_t) S_t)$
9          $d_t = -H_t \left( \nabla f_{\mathcal{S}}(x_t) - \nabla f_{\mathcal{S}}(w_k) + \mu \right)$
10         $x_{t+1} = x_t + \eta d_t$
11      **end**
12      **Option I:** $w_{k+1} = x_m$
13      **Option II:** $w_{k+1} = x_i$, $i$ selected uniformly at random from $[m]$;
14 **end**

# Convergence - Assumptions

There exist constants $\lambda, \Lambda \in \mathbb{R}_+$ such that

- $f$ is $\lambda$–strongly convex

$$f(w) \geq f(x) + \nabla f(x)^T (w - x) + \frac{\lambda}{2} \|w - x\|_2^2, \qquad (2)$$

- $f$ is $\Lambda$–smooth

$$f(w) \leq f(x) + \nabla f(x)^T (w - x) + \frac{\Lambda}{2} \|w - x\|_2^2, \qquad (3)$$

- These assumptions imply that

$$\lambda I \preceq \nabla^2 f_{\mathcal{S}}(w) \preceq \Lambda I, \quad \text{for all } x \in \mathbb{R}^d, \mathcal{S} \subseteq [n], \qquad (4)$$

- from which we can prove that there exist constants $\gamma, \Gamma \in \mathbb{R}_+$ such that for all $k$ we have

$$\gamma I \preceq H_k \preceq \Gamma I. \qquad (5)$$

# Linear Convergence

## Theorem

*Suppose that the Assumptions hold. Let $w_*$ be the unique minimizer of $f(w)$. Then in our Algorithm, we have for all $k \geq 0$ that*

$$\mathbb{E} f(w_k) - f(w_*) \leq \rho^k \mathbb{E} f(w_0) - f(w_*),$$

*where the convergence rate is given by*

$$\rho = \frac{1/2m\eta + \eta \Gamma^2 \Lambda (\Lambda - \lambda)}{\gamma \lambda - \eta \Gamma^2 \Lambda^2} < 1,$$

*assuming we have chosen $\eta < \gamma \lambda / (2\Gamma^2 \Lambda^2)$ and that we choose $m$ large enough to satisfy*

$$m \geq \frac{1}{2\eta \left(\gamma \lambda - \eta \Gamma^2 \Lambda (2\Lambda - \lambda)\right)},$$

*which is a positive lower bound given our restriction on $\eta$.*

# Upper and lower bounds on eigenvalues of $H_k$

- Under the assumption that

$$\lambda I \preceq \nabla^2 f_{\mathcal{T}}(x) \preceq \Lambda I, \quad \forall x \in \mathbb{R}^d \qquad (6)$$

  there exist constants $\gamma, \Gamma \in \mathbb{R}_+$ such that for all $k$ we have

$$\gamma I \preceq H_k \preceq \Gamma I. \qquad (7)$$

  where

$$\gamma \geq \frac{1}{1 + M\Lambda}, \quad \Gamma \leq (1 + \sqrt{\kappa})^{2M}(1 + \frac{1}{\lambda(2\sqrt{\kappa} + \kappa)}), \quad \kappa \equiv \Lambda/\lambda. \qquad (8)$$

- Previously derived bounds depend on the problem dimension $d$; e.g. $\Gamma \sim ((d + M)\kappa)^{d+M}$

# gisette-scale $d = 5,000, n = 6,000$



Figure: gisette

# covtype-libsvm-binary $d = 54, n = 581,012$



Figure: covtype.libsvm.binary

# Higgs $d = 28, n = 11,000,000$



Figure: HIGGS

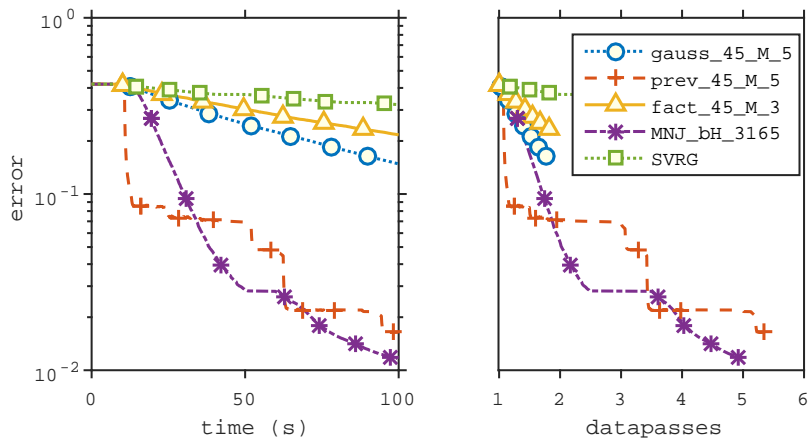Figure: SUSY

# epsilon-normalized $d = 2,000, n = 400,000$



Figure: epsilon-normaliized

Figure: rcv1-train

Figure: url-combined

# Contributions
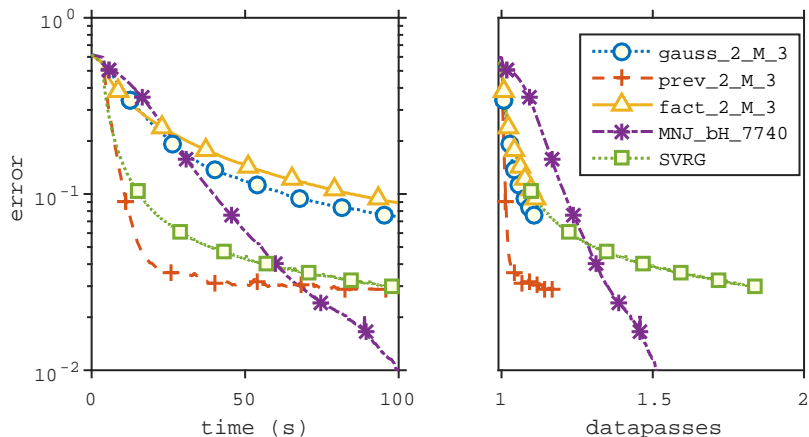
- *New metric learning framework.* A block BFGS framework for gradually learning the metric of the underlying function using a sketched form of the subsampled Hessian matrix
- *New limited-memory block BFGS method.* May also be of interest for non-stochastic optimization
- *New limited-memory factored form block BFGS method.*
- *Several sketching matrix possibilities.*
- *Linear convergence rate* proof for our methods.
- *Tighter upper and lower bounds* on the eigenvalues of the variable metric