Robust inversion and randomized sampling

Michael P. Friedlander University of British Columbia

International Symposium on Mathematical Programming August 19–24, 2012

Collaborators: Sasha Aravkin, Felix Herrmann, Tristan van Leeuwen, and Mark Schmidt

Model problem

$$\underset{x}{\text{minimize}} \quad f(x) \quad := \frac{1}{m} \sum_{i=1}^{m} f_i(x)$$

Examples:

- least-squares $f_i(x) = (a_i^T x b_i)^2$ $f(x) = ||Ax b||^2$
- log likelihood $f_i(x) = -\log p(b_i; x)$
- sample average $f_i(x) = f(x, \omega_i)$ $f(x) \approx \mathbb{E}_{\omega}[f(x, \omega)]$

Context:

- *m* large
- each $f_i(x)$ and $\nabla f_i(x)$ expensive to evaluate

Computing costs:

- count $f_i / \nabla f_i$ evals, not $f / \nabla f$ evals
- minimize passes through full data set (f_1, \ldots, f_m)

SEISMIC INVERSION

Reflection seismology





Full waveform inversion: velocity models



Full waveform inversion

Each of *m* experiments yields a vector of measurements:

sources: q_1, \ldots, q_m

measurements: d_1, \ldots, d_m

1 source, 1 frequency:

 $\underset{x,u}{\text{minimize }} \|d - Pu\|^2 \quad \text{subj to} \quad H_{\omega}(x)u = q$

All sources, all frequencies:

(eg, 1k sources,
$$\sim$$
 10 freqs)

 \sim

$$\underset{x}{\text{minimize}} \sum_{i}^{m} \sum_{\omega \in \Omega} \|d_{i} - PH_{\omega}(x)^{-1}q_{i}\|^{2}$$

Main cost is solution of Helmholtz equation for each (i, ω) pair:

$$H_{\omega}(x)u=q_i$$

Dimensionality reduction and stochastic optimization

Use all your data:

$$f(x) = \sum_{i}^{m} ||r_i(x)||^2$$
 with $r_i(x) = d_i - F(x)q_i$

Dimensionality reduction: form *s* weighted averages ($s \ll m$)

$$ar{d}_j := \sum_j^m w_{ij} d_i$$
 and $ar{q}_j := \sum_j^m w_{ij} q_i, \qquad j = 1, \dots, s$

Stochastic approximation of the misfit:

$$ar{m{f}}(m{x}) = \sum_j^s \|ar{r}_j(m{x})\|^2$$
 with $ar{r}_j(m{x}) = ar{d}_j - F(m{x})ar{q}_j$

Stochastic optimization interpretation holds if $\mathbb{E}[WW^T] = I$:

$$\mathbb{E}[\overline{f}(x)] = f(x)$$
 and $\mathbb{E}[\nabla \overline{f}(x)] = \nabla f(x)$

Stochastic trace estimation

Least-squares misfit is a matrix trace:

$$f(x) = \sum_{i}^{m} ||r_i||^2 = \operatorname{trace}(R^T R) \quad \text{with} \quad R = [r_1 | \cdots | r_m]$$

Data mixing (dimensionality reduction) equivalent to

$$R \qquad W = \bar{R}$$

gives sample-average function

$$\bar{f}(x) = \sum_{j}^{s} \|\bar{r}_{j}\|^{2} = \operatorname{trace}(\bar{R}^{T}\bar{R}) \quad \text{with} \quad \bar{R} = [\bar{r}_{1} | \cdots | \bar{r}_{s}], \quad s \ll m$$

Connected to stochastic trace estimation.

- Hutchinson ('90) minimize var[trace($\bar{R}^T \bar{R}$)] with $W \sim$ Rademacher
- Avron & Toledo ('11) other optimal choices for \emph{W}
- Haber, Chung, Herrmann ('12) connection to inverse problems

Nonlinear least-squares with corrupted data



good data

4% bad data

collaboration with Total

Robust misfit measures



- Least-squares penalty assumes ε ∼ N(0, 1)
- **Theorem:** Log-concave densities—ie, those induced by convex penalties—are all **exponentially bounded**:

 $\operatorname{prob}(|x| > t + \Delta t \quad \text{given} \quad |x| > t) = \mathcal{O}(e^{-\Delta t})$

Heavy-tailed densities: robust to wrong data & approximate models

Robust fullwaveform inversion results

- Recover velocity on a 2D grid: 201×301
- 151 sources, 6 frequencies: 906 PDE solves per $f/\nabla f$ -eval
- 50% corrupted data





Sampling strategies for dimensionality reduction

Generic inverse problem (assuming iid observations):

$$\min_{x} f(x) := \frac{1}{m} \sum_{i}^{m} \rho(r_i) \quad \text{with} \quad R(x) = [r_1 | \cdots | r_m]$$

Moment condition $E[WW^T] = I$ is **not generally sufficient** to guarantee

$$\mathbb{E}[\quad \bar{f}(x)] = \quad f(x) \\ \mathbb{E}[\nabla \bar{f}(x)] = \nabla f(x) \quad \text{with} \quad \bar{R}(x) = R(x)W$$

Random subset selection without replacement, ie,

$$\bar{R}(x) = [r_{i(1)} | r_{i(2)} | \cdots | r_{i(s)}]$$

does yield desired "expected objective" property

MODEL PROBLEM

$$\underset{x}{\text{minimize}} \quad f(x) := \frac{1}{m} \sum_{i=1}^{m} f_i(x)$$

Complexity of steepest descent

Baseline Algorithm:

$$x_{k+1} \leftarrow x_k - \alpha_k \nabla f(x_k), \qquad \alpha_k \equiv 1/L$$

Assume: convex f; Lipschitz ∇f with param L

Sublinear rate:

•
$$f(x_k) - f(x_*) = O(1/k)$$

•
$$f(x_k) - f(x_*) = O(1/k^2)$$

(constant stepsize) (optimal rate with extrapolation) [Nesterov '83; Tseng '10]

Linear rate: additionally assume that f is strongly convex w/ param μ

•
$$f(x_k) - f(x_*) = O([1 - \mu/L]^k)$$

Note: if *f* is twice differentiable, $\mu I \preceq \nabla^2 f(x) \preceq LI$

Incremental gradient methods

Algorithm: $x_{k+1} \leftarrow x_k - \alpha \nabla f_i(x_k), i \in \{1, \dots, m\}$ cyclic, randomized

Assume: f strongly convex (μ) and Lipschitz ∇f (L)

Constant stepsize: $\alpha_k \equiv \bar{\alpha}$

• $\|x_k - x_*\|^2 \le \mathcal{O}([1 - \mu/L]^k) + \mathcal{O}(\boldsymbol{m}^2 \bar{\boldsymbol{\alpha}})$ k full cycles

•
$$\mathbb{E} \|x_k - x_*\|^2 \le \mathcal{O}([1 - \mu/L]^k) + \mathcal{O}(\boldsymbol{m}\bar{\alpha})$$
 k iterations

Decreasing stepsize: $\sum_k \alpha_k = \infty$, $\sum_k \alpha_k^2 < \infty$

- $\|x_k x_*\|^2 = \mathcal{O}(1/k)$ k full cycles
- $\mathbb{E} \|x_k x_*\|^2 = \mathcal{O}(1/k)$ k iterations

Many variations:

Luo/Tseng '94; Nedić/Bertsekas '00/'10; Blatt et al '08; Bottou '10

EXAMPLES

Seismic inversion

Recover image of geological structures via nonlinear least squares

$$\underset{x}{\text{minimize}} \sum_{i}^{m} \sum_{\omega \in \Omega} \|d_{i} - PH_{\omega}(x)^{-1}q_{i}\|^{2}$$

Observations: Each of *m* "shots" is an experiment:

sources: q_1, \ldots, q_m , measurements: d_1, \ldots, d_m





0.010.40.822.6471016223039 of 39 passes

Image denoising



- Statistical denoising via conditional random fields
- Dataset of 50 synthetic 64 × 64 images [Kumar/Hebert '04]
- Generalization of logistic model to capture dependencies among labels

$$\underset{x}{\text{maximize}} \sum_{i=1}^{m} \log p(b_i; x)$$

• p is intractable and approximated



full gradient



0.250.500.7512345 of 5 passes

Sampling approach

Increasing batch:

[Bertsekas/Tsitsiklis '96, Shapiro/H-do-M '00]

 $\mathcal{S}_k \subseteq \{1, \dots, m\}, \qquad s_k o m \quad (\text{slowly})$

Sample-average gradient:

$$g_k(x) := \frac{1}{s_k} \sum_{i \in S_k} \nabla f_i(x)$$

Algorithm:

$$x_{k+1} \leftarrow x_k - \alpha_k d$$
 with $H_k d = -g_k(x_k)$

Goal: non-asymptotic analysis based on controlling gradient error

- $g_k(x) = \nabla f(x_k) + e_k$ where $\|e_k\|^2 \le \epsilon_k$ or $\mathbb{E}[\|e_k\|^2] \le \epsilon_k$
- How to control sample size *s_k*?

Gradient with generic errors

Prototype algo:

$$x_{k+1} \leftarrow x_k - \alpha g_k$$
, $g_k = \nabla f(x_k) + e_k$, α fixed

Assumptions:

- Lipschitz gradient (L); strong convexity (μ)
- $\|e_k\|^2 \leq \epsilon_k$

Convergence rate: for all k = 1, 2, ... [F. & Schmidt '12]

$$\|x_k - x_*\|^2 \leq \mathcal{O}([1 - \mu/L]^k) + \mathcal{O}(\epsilon_k)$$

Examples:

Growing the sample size

Prototype algo:

$$\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k - \alpha \mathbf{g}_k, \qquad \mathbf{g}_k = \frac{1}{s_k} \sum_{i \in \mathcal{S}_k} \nabla f_i(\mathbf{x}_k), \qquad \mathcal{S}_k \subseteq \{1, \dots, m\}$$

Sampling strategies

- Deterministic: pre-determined sample sequence
- Randomized: uniform sampling

Convergence rates for all k = 1, 2, ... [F. & Schmidt '12]

deterministic: $\|x_{k} - x_{*}\|^{2} = \mathcal{O}([1 - \mu/L]^{k}) + \mathcal{O}\left(\left[\frac{m-s_{k}}{m}\right]^{2}\right)$ sampling w/o replacement $\boldsymbol{E}\|x_{k} - x_{*}\|^{2} = \mathcal{O}([1 - \mu/L]^{k}) + \mathcal{O}\left(\frac{m-s_{k}}{m} \cdot \frac{1}{s_{k}}\right)$ sampling w/ replacement $\boldsymbol{E}\|x_{k} - x_{*}\|^{2} = \mathcal{O}([1 - \mu/L]^{k}) + \mathcal{O}\left(\frac{1}{s_{k}}\right)$

Illustration



In practice

Algorithm:

$$x_{k+1} \leftarrow x_k - \alpha_k d, \qquad \boldsymbol{B}_k d = -\boldsymbol{g}_k(x_k), \qquad \boldsymbol{g}_k(x) = \frac{1}{s_k} \sum_{i \in \mathcal{S}_k} \nabla f_i(x)$$

Hessian approximations **B**_k:

• (limited-memory) quasi-Newton:

[Schraudolph et al '07]

$$s_k := x_{k+1} - x_k, \qquad y_k := g_k(x_{k+1}) - g_k(x_k)$$

sample-average Hessian:

$$\boldsymbol{B}_{\boldsymbol{k}} := rac{1}{h_k} \sum_{i \in \mathcal{H}_k} \nabla^2 f_i(\boldsymbol{x}_k), \quad h_k \ll \boldsymbol{s}_k$$

• Fisher information [for $f(x) = -\sum_{i} \log p_i(\omega_i; x)$]: [Osborne '92]

$$\boldsymbol{B_k} \approx \mathbb{E}_{\omega}[\nabla^2 f(x)] = \mathbb{E}_{\omega}[\nabla f(x)\nabla f(x)^{\mathsf{T}}]$$

APPLICATIONS

Binary logistic regression

$$\max_{x} \sum_{i}^{m} \log p(b_i \mid a_i, x), \quad p(b_i \mid a_i, x) = \frac{1}{1 + \exp(-b_i a_i^T x)}, \quad b_i \in \{-1, 1\}$$

- Email spam classifier (Cormack and Lynam, 2005)
- TREC 2005 dataset: 92,189 email msgs from Enron investigation



Multinomial logistic regression

$$\max_{x} \sum_{i}^{m} \log p(b_{i} = j \mid a_{i}, \{x\}_{j \in \mathcal{C}}), \quad b_{i} \in \mathcal{C}$$

- Digit classification 0 1 2 3 4 5 6 7 8 9
- MNIST dataset: 70,000 handwritten 28×28 images of digits



Chain-structured conditional random fields

$$\max_{x} \sum_{i}^{m} \log p(\{b_{i}^{k} = j_{k}\}_{k \in \Omega} \mid \{a_{i}^{k}\}_{k \in \Omega}, \{x_{j}\}_{j \in \mathcal{C}}), \quad b_{i}^{k} \in \mathcal{C}, \ k \in \Omega$$

- noun-phrase chunking task from natural-language processing
- CoNLL-2000 Shared Task dataset: 211,727 words in 8,936 sentences



General conditional random fields



Seismic inversion

$$\min_{x} \sum_{i}^{m} \sum_{\omega \in \Omega} \|d_i - PH_{\omega}(x)^{-1}q_i\|^2$$

- Recover seismic image via nonlinear least squares
- Marmousi 2D acoustic model; 101 sources/receivers; 8 frequencies



No strong convexity, no linear convergence

Steepest descent

$$x_{k+1} \leftarrow x_k - \alpha \nabla f(x_k)$$

is sublinear:

$$f(x_k) - f(x_*) = \mathcal{O}(1/k)$$

Approximate steepest descent

$$x_{k+1} \leftarrow x_k - lpha d_k$$
 with $d_k =
abla f(x_k) + e_k$

Summable errors, ie, $\sum_{k}^{\infty} \|e_k\| < \infty$, gives sublinear iterate average:

$$f(ar{x}_k)-f(x_*)=\mathcal{O}(1/k), \quad ar{x}_k:=rac{1}{k}\sum_{i=1}^k x_i$$

References I

- H. Avron and S. Toledo. Randomized algorithms for estimating the trace of an implicit symmetric positive semi-definite matrix. J. ACM, 58:8:1–8:34, April 2011.
- D. Bertsekas. Incremental Gradient, Subgradient, and Proximal Methods for Convex Optimization: A Survey. In S. Sra, S. Nowozin, and S. J. Wright, editors, *Optimization for* machine learning, chapter 4, pages 85–115. MIT Press, 2012.
- D. Bertsekas and J. Tsitsiklis. Neuro-dynamic programming. Athena Scientific, 1996.
- D. Blatt, A. Hero, and H. Gauchman. A convergent incremental gradient method with a constant step size. SIAM J. Optim., 18(1):29–51, 2008.
- L. Bottou. Large-scale machine learning with stochastic gradient descent. In Y. Lechevallier and G. Saporta, editors, *Proceedings of the 19th International Conference on Computational Statistics (COMPSTAT'2010)*, pages 177–187, Paris, France, August 2010. Springer.
- R. H. Byrd, G. Chin, J. Nocedal, and Y. Wu. Sample size selection in optimization methods for machine learning. *Math. Program.*, 134:127–155, 2012.
- R. H. Byrd, G. M. Chin, W. Neveitt, and J. Nocedal. On the use of stochastic hessian information in optimization methods for machine learning. *SIAM J. Optim.*, 21(3):977–995, 2011.
- L. Grippo. A class of unconstrained minimization methods for neural network training. Optim. Methods Softw., 4(2):135–150, 1994.
- E. Haber, M. Chung, and F. J. Herrmann. An effective method for parameter estimation with PDE constraints with multiple right-hand sides. *SIAM J. Optim.*, 22(3):739–757, 2012.

References II

- M. Hutchinson. A stochastic estimator of the trace of the influence matrix for laplacian smoothing splines. *Comm. Statist. Simulation Comput.*, 19(2):433–450, 1990.
- A. Kleywegt, A. Shapiro, and T. Homem-de Mello. The sample average approximation method for stochastic discrete optimization. SIAM J. Optim., 12(2):479–502, 2002.
- Z. Luo and P. Tseng. Error bounds and convergence analysis of feasible descent methods: A general approach. Ann. Oper. Res., 46(1):157–178, 1993.
- A. Nedic and D. Bertsekas. Convergence rate of incremental subgradient algorithms. Stochastic Optimization: Algorithms and Applications, pages 263–304, 2000.
- Y. Nesterov. A method for unconstrained convex minimization problem with the rate of convergence O(1/k²). Soviet Math. Dokl., 269, 1983.
- M. Osborne. Fisher's method of scoring. Intern. Stat. Rev., pages 99-117, 1992.
- N. N. Schraudolph, J. Yu, and S. Günter. A Stochastic Quasi-Newton Method for Online Convex Optimization. In M. Meila and X. Shen, editors, Proc. 11th Intl. Conf. Artificial Intelligence and Statistics (AIStats), volume 2 of Workshop Conf. Proc., pages 436–443, San Juan, Puerto Rico, 2007. J. Machine Learning Res.
- P. Tseng. Approximation accuracy, gradient methods, and error bound for structured convex optimization. *Math. Program.*, 125:263–295, 2010.
- K. van den Doel and U. M. Ascher. Adaptive and stochastic algorithms for electrical impedance tomography and dc resistivity problems with piecewise constant solutions and many measurements. *SIAM J. Comput.*, 34(1), 2012.

Thanks!

Read:

- M. P. Friedlander and M. Schmidt. Hybrid deterministic-stochastic methods for data fitting. *SIAM J. Scientific Computing*, 34(3), April 2012.
- A. Aravkin, M. P. Friedlander, F. Herrmann, and T. van Leeuwen. Robust inversion, dimensionality reduction, and randomized sampling. *Mathematical Programming*, 134(1):101–125, 2012.

Email:

mpf@cs.ubc.ca

Surf:

http://www.cs.ubc.ca/~mpf