Proximal Newton Methods for Large-Scale Machine Learning

Inderjit Dhillon Computer Science & Mathematics UT Austin

ShanghaiTech Symposium on Data Science June 24, 2015

Joint work with C. J. Hsieh, M. Sustik, P. Ravikumar, R. Poldrack, S. Becker, and P. Olsen

伺 ト イ ヨ ト イ ヨ ト

Big Data

Link prediction



LinkedIn.



Spam classification

Gmail -





gene-gene network

Image classification







イロト イポト イヨト イヨト

Modern Machine Learning Problems

Inverse Covariance Estimation

- Learning graph (structure and/or weights) from data.
- Gene networks, fMRI analysis, stock networks, climate analysis.



- 4 同 🕨 - 4 目 🕨 - 4 目

Matrix Completion

- Missing value estimation for recommender systems.
- Assumption: the underlying matrix is low-rank.



-8.72	0.03	-1.03 0.62	-0.07	-0.11	-0.53	-0.46	-0.06	-0.05	-0.53	-0.07	-0.35	-0.19	-0.14
			0.13	-0.42	0.45	0.17	-0.25	-0.17	-0.18	0.27	-0.59	0.05	0.14
			-0.21	-0.43	-0.23	0.16	0.08	0.17	0.57	-0.39	-0.37	-0.08	-0.15

-3.52 3.73 -3.32

-7.78 2.34 2.33

-2.44 -5.29 -3.92

-1.78 1.90 -1.68

Robust PCA

- Image de-noising, face recognition, graphical modeling with latent variables
- Decompose the input matrix as sparse + low-rank matrices.



Multi-task Learning

- Training T classification tasks together.
- Sparse models: Lasso
- Feature sharing: Group Lasso



Inverse Covariance Estimation

Inderjit Dhillon Computer Science & Mathematics UT Austin Proximal Newton Methods for Large-Scale Machine Learning

3

Estimate Relationships between Variables

- Example: FMRI Brain Analysis.
- Goal: Reveal functional connections between regions of the brain.
 (Sun et al, 2009; Smith et al, 2011; Varoquaux et al, 2010; Ng et al, 2011)



Figure from (Varoquaux et al, 2010)

Estimate Relationships between Variables

- Example: Gene regulatory network discovery.
- Goal: Gene network reconstruction using microarray data. (Schafer & Strimmer 2005; Andrei & Kendziorski 2009; Menendez et al, 2010; Yin and Li, 2011)



Figure from (Banerjee et al, 2008)

Other Applications

- Financial Data Analysis:
 - Model dependencies in multivariate time series (Xuan & Murphy, 2007).
 - Sparse high dimensional models in economics (Fan et al, 2011).
- Social Network Analysis / Web data:
 - Model co-authorship networks (Goldenberg & Moore, 2005).
 - Model item-item similarity for recommender system(Agarwal et al, 2011).
- Climate Data Analysis (Chen et al., 2010).
- Signal Processing (Zhang & Fung, 2013).
- Anomaly Detection (Ide et al, 2009).

伺い イラト イラト

Inverse Covariance Estimation

- Given: *n* i.i.d. samples $\{y_1, \ldots, y_n\}$, $y_i \in R^p$, $y_i \sim \mathcal{N}(\mu, \Sigma)$,
- Goal: Estimate relationships between variables.
- The sample mean and covariance are given by:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{y}_i$$
 and $S = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{y}_i - \hat{\mu}) (\mathbf{y}_i - \hat{\mu})^T$.

- Covariance matrix Σ may not directly convey relationships.
- An example Chain graph: $y(j) = 0.5y(j-1) + \mathcal{N}(0,1)$

Inverse Covariance Estimation (cont'd)

• Conditional independence is reflected as zeros in Θ (= Σ^{-1}):

 $\Theta_{ij} = 0 \Leftrightarrow y(i), y(j)$ are conditionally independent given other variables

• Chain graph example: $y(j) = 0.5y(j-1) + \mathcal{N}(0,1)$

$$\Theta = \left(\begin{array}{cc} 1 & -0.5 & 0 \\ -0.5 & 1.25 & -0.5 \\ 0 & -0.5 & 1 \end{array} \right)$$

- In a GMRF G = (V, E), each node corresponds to a variable, and each edge corresponds to a non-zero entry in Θ.
- Goal: Estimate the inverse covariance matrix Θ from the observations.

周 ト イラ ト イラ ト 一 ラ

Maximum Likelihood Estimator: LogDet Optimization

• Given the *n* samples, the likelihood is

$$P(\mathbf{y}_1, \dots, \mathbf{y}_n; \hat{\boldsymbol{\mu}}, \Theta) \propto \prod_{i=1}^n (\det \Theta)^{1/2} \exp\left(-\frac{1}{2}(\mathbf{y}_i - \hat{\boldsymbol{\mu}})^T \Theta(\mathbf{y}_i - \hat{\boldsymbol{\mu}})\right)$$
$$= (\det \Theta)^{n/2} \exp\left(-\frac{1}{2}\sum_{i=1}^n (\mathbf{y}_i - \hat{\boldsymbol{\mu}})^T \Theta(\mathbf{y}_i - \hat{\boldsymbol{\mu}})\right)$$

• The log likelihood can be written as

$$\log(P(\mathbf{y_1},\ldots,\mathbf{y_n};\hat{\boldsymbol{\mu}},\Theta)) = \frac{n}{2}\log(\det\Theta) - \frac{n}{2}\mathrm{tr}(\Theta S) + \mathrm{constant}.$$

• Maximum likelihood estimator:

$$\Theta^* = \arg\min_{X\succ 0} \{-\log \det X + \operatorname{tr}(SX)\}.$$

• Problem: when p > n, sample covariance matrix S is singular.

L1-regularized inverse covariance selection

- A sparse inverse covariance matrix is preferred add l₁ regularization to promote sparsity.
- The resulting optimization problem:

$$\Theta = \arg\min_{X\succ 0} \left\{ -\log \det X + \operatorname{tr}(SX) + \lambda \|X\|_1 \right\} = \arg\min_{X\succ 0} f(X),$$

where
$$||X||_1 = \sum_{i,j=1}^n |X_{ij}|.$$

- Regularization parameter $\lambda > 0$ controls the sparsity.
- The problem appears hard to solve:
 - Non-smooth log-determinant program.
 - Number of parameters scale quadratically with number of variables.

Optimization in Machine Learning

Linear SVM Kernel SVM Logistic Regression Covariance Selection Metric Learning Matrix Completion Neural Network Link Prediction Clustering



Gradient Descent Stochastic Gradient Coordinate Descent Newton Method Interior Point Method Alternating Linearization Linear Programming

High-dimensional! Large number of samples! Large number of parameters!

Exploiting Structure in Machine Learning Problems

• Regularized Loss Minimization(RLM):

$$\min_{\boldsymbol{\theta}} \{ \underbrace{g(\boldsymbol{\theta}, X)}_{\text{loss}} + \underbrace{h(\boldsymbol{\theta})}_{\text{regularization}} \} \equiv f(\boldsymbol{\theta}),$$

- Exploiting problem structure:
 - Can we have faster computation?
- Exploiting model structure:
 - Can we avoid un-important search space?
- Exploiting Data distribution:
 - Can we speed up optimization algorithms if we roughly estimate the data distribution?

QUIC

QUadratic approximation for Inverse Covariance estimation

Inderjit Dhillon Computer Science & Mathematics UT Austin Proximal Newton Methods for Large-Scale Machine Learning

< ∃ >

Sparse Inverse Covariance Estimation

• p^2 parameters in the non-smooth Log-determinant program:

$$\Theta = \arg\min_{X \succ 0} \left\{ \underbrace{-\log \det X + \operatorname{tr}(SX)}_{\text{negative log likelihood}} + \lambda \|X\|_1 \right\}$$



First Order Methods vs Second Order Methods

- Many algorithms have been proposed.
 - Block coordinate descent (Banerjee et al., 2007), (Friedman et al., 2007),
 - Gradient or Accelerate gradient descent (Lu, 2009), (Duchi et al., 2008),
 - Greedy coordinate descent (Scheinberg & Rich., 2009)
 - ADMM (Scheinberg et al., 2009), (Boyd et al., 2011)
- All of them are first order methods (only gradient information)
- QUIC: the first second order method for this problem (use Hessian information).

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_d} \end{bmatrix}, \quad H = \nabla^2 f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_d} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_d \partial x_1} & \frac{\partial^2 f}{\partial x_d \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_d^2} \end{bmatrix}$$

ヨッ イヨッ イヨッ

First Order Methods vs Second Order Methods



Why do all previous approaches use first order methods?

Second Order Methods for Non-Smooth Functions

• Update rule for Newton method (second order method):

$$\mathbf{x} \leftarrow \mathbf{x} - \eta (\nabla^2 f(\mathbf{x}))^{-1} \nabla f(\mathbf{x})$$

- Two difficulties:
 - Time complexity per iteration may be very high.
 - How to deal with non-smooth functions $(||\mathbf{x}||_1)$?

Second Order Methods for Non-Smooth Functions

• Update rule for Newton method (second order method):

$$\mathbf{x} \leftarrow \mathbf{x} - \eta (\nabla^2 f(\mathbf{x}))^{-1} \nabla f(\mathbf{x})$$

- Two difficulties:
 - Time complexity per iteration may be very high.
 - How to deal with non-smooth functions $(||\mathbf{x}||_1)$?
- Addressed by (Hsieh, Sustik, Dhillon and Ravikumar, 2011).
- Extensions:
 - Theoretical Analysis: (Lee et al., 2012), (Patrinos and Bemporad., 2013), (Tang and Scheinberg, 2014), (Yen, Hsieh, Ravikumar and Dhillon, 2014), ...
 - Optimization Algorithms: (Olsen et al., 2012), (Dinh et al., 2013), (Treister et al., 2014), (Scheinberg and Tang, 2014), (Zhong et al., 2014), ...

伺い イラト イラト

- Newton direction: $H^{-1}g$. (*H*: Hessian, g: gradient)
- Solving the linear system exactly: $O(p^6)$ time.



Appears to be prohibitive to apply a second order method.

• Coordinate descent: each time fits one element of g.



• Coordinate descent: each time fits one element of g.



• Use iterative solver—coordinate descent

 $O(p^2)$ per coordinate update, $O(p^4)$ for one sweep. \Rightarrow need 40 days for $p = 10^4$ on a 2.83GHz CPU.



Still prohibitive.

Exploiting Problem Structure Fast Computation of Newton Direction



• Structure of the Hessian:

$$H = \frac{\partial^2}{\partial^2 X} (-\log \det X) = X^{-1} \otimes X^{-1}.$$

• The Kronecker product:

$$A \otimes B = \begin{pmatrix} a_{11}B & \cdots & a_{1p}B \\ \vdots & \ddots & \vdots \\ a_{p1}B & \cdots & a_{pp}B \end{pmatrix}$$

was introduced by G. Zehfuss (1858) and later attributed to L. Kronecker.



The Hessian (p^2 by p^2 matrix) can be represented using p^2 parameters!



- Time complexity: $O(p^2) \Rightarrow O(p)$ for each coordinate update. $O(p^4) \Rightarrow O(p^3)$ per sweep.
- At coordinate descent step to update (i, j) element of descent direction:

$$\begin{pmatrix} H \operatorname{vec}(\Delta) \end{pmatrix}_{ip+j} = \left((X^{-1} \otimes X^{-1}) \operatorname{vec}(\Delta) \right)_{ip+j} = \left(X^{-1} \Delta X^{-1} \right)_{ij}$$
$$= (X^{-1} U)_{ij},$$

where $U = \Delta X^{-1}$ is maintained in memory.

Newton Method for Non-smooth Functions

- Iteration conducts the following updates:
 - Compute the generalized Newton direction

$$D_{t} = \arg\min_{\Delta} \bar{f}_{X_{t}}(\Delta)$$

$$= \arg\min_{\Delta} \langle \nabla g(X_{t}), \operatorname{vec}(\Delta) \rangle + \frac{1}{2} \operatorname{vec}(\Delta)^{T} \nabla^{2} g(X_{t}) \operatorname{vec}(\Delta) + \lambda \|X_{t} + \Delta\|_{1}$$
• $X_{t+1} \leftarrow X_{t} + D_{t}$.
$$\bar{f}_{X_{t}}(\Delta)$$

$$X_{t}$$

 X_{t+1}

Proximal Newton Method

Proximal Newton for sparse Inverse Covariance estimation

For t = 0, 1, ...

- **Compute Proximal Newton direction**: $D_t = \arg \min_{\Delta} \bar{f}_{X_t}(X_t + \Delta)$ (A Lasso problem.)
- **2** Line Search: use an <u>Armijo</u>-rule based step-size selection to get α s.t. $X_{t+1} = X_t + \alpha D_t$ is
 - positive definite,
 - satisfies a sufficient decrease condition f(X_t + αD_t) ≤ f(X_t) + ασΔ_t.

伺 ト イ ヨ ト イ ヨ ト

- Can we scale to $p \approx 20,000$?
 - Coordinate descent updates for computing Newton direction.
 - needs O(p²) storage
 - needs $O(p^4)$ computation per sweep
 - needs $O(p^3)$ computation per sweep (by exploiting problem structure)
 - Line search (compute determinant using Cholesky factorization).
 - needs $O(p^2)$ storage
 - needs $O(p^3)$ computation

Exploiting Model Structure Variable Selection

$$\min_{\theta} \{ \underbrace{-\log \det X + \operatorname{tr}(SX)}_{\text{problem structure}} + \underbrace{h(\theta)}_{\text{regularization}} \} \equiv f(\theta),$$

Sparsity of the Model

ER dataset, p = 692, solutions at iteration 1, 3, 5.



Only 2.7% variables become nonzero.

Inderjit Dhillon Computer Science & Mathematics UT Austin Proximal Newton Methods for Large-Scale Machine Learning

- 4 同 🕨 - 4 目 🕨 - 4 目
- Strategy: before solving the Newton direction, "guess" the variables that should be updated.
- Fixed set: $S_{fixed} = \{(i,j) \mid X_{ij} = 0 \text{ and } |\nabla_{ij}g(X)| \le \lambda.\}.$
- Only update the rest of variables (free set).
- Convergence Guaranteed.

p = 1869, number of variables $= p^2 = 3.49$ million. The size of free set drops to 20,000 very quickly.

Active Variable Selection



Updates on all elements

Updates on Free Set

伺 ト イヨト イヨト

3

Exploiting Model Structure

- Can we scale to $p \approx 20,000$?
 - Coordinate descent updates for computing Newton direction.
 - needs O(p²) storage
 - needs O(p⁴) computation per sweep
 - needs $O(p^3)$ computation per sweep (by exploiting problem structure)
 - needs O(mp) computation per sweep, where $m \approx ||X||_0$ (by exploiting problem and model structure)
 - Line search (compute determinant using Cholesky factorization).
 - needs O(p²) storage
 - needs $O(p^3)$ computation

ヨッ イヨッ イヨッ

$\operatorname{QUIC}:$ QUadratic approximation for sparse Inverse Covariance estimation

Input: Empirical covariance matrix S, scalar λ , initial X_0 . For t = 0, 1, ...

- **(**) Variable selection: select a *free* set of $m \ll p^2$ variables.
- Use coordinate descent to find descent direction:
 - $D_t = \arg \min_{\Delta} \bar{f}_{X_t}(X_t + \Delta)$ over set of free variables, (A Lasso problem.)
- Solution Search: use an Armijo-rule based step-size selection to get α s.t. $X_{t+1} = X_t + \alpha D_t$ is
 - positive definite,
 - satisfies a sufficient decrease condition f(X_t + αD_t) ≤ f(X_t) + ασΔ_t.

- * 同 * * ヨ * * ヨ * - ヨ

(Cholesky factorization of $X_t + \alpha D_t$)

QUIC can solve p = 20,000 in about 2.3 hours.

• Current iterate X_t .



<ロ> (日) (日) (日) (日) (日)

э

• Form the quadratic approximation $\bar{f}_{X_t}(\Delta)$.

$$f(X) = g(X) + \lambda \|X\|_1$$

$$\bar{f}_{X_t}(\Delta) = f(X_t) + \langle \nabla g(X), \Delta \rangle + \frac{1}{2} \operatorname{vec}(\Delta)^T \nabla^2 g(X) \operatorname{vec}(\Delta) + \lambda \| X + \Delta \|_1$$



< ∃ >

э

- Form the quadratic approximation $\bar{f}_{X_t}(\Delta)$.
- Free set selection.



< ∃ →

- Form the quadratic approximation $\bar{f}_{X_t}(\Delta)$.
- Free set selection.
- Coordinate descent to minimize $\bar{f}_{X_t}(\Delta)$



4 3 b

- Form the quadratic approximation $\bar{f}_{X_t}(\Delta)$.
- Free set selection.
- Coordinate descent to minimize $\bar{f}_{X_t}(\Delta)$



< ∃ →

- Form the quadratic approximation $\bar{f}_{X_t}(\Delta)$.
- Free set selection.
- Coordinate descent to minimize $\bar{f}_{X_t}(\Delta)$
- $X_{t+1} \leftarrow X_t + D_t$



/□ ▶ < 글 ▶ < 글

Theorem: Quadratic Convergence Rate for QUIC

Suppose $g(\cdot)$ is strongly convex and the quadratic subproblem at each iteration is solved exactly, QUIC converges to the global optimum with an asymptotic quadratic convergence rate:

$$\lim_{t \to \infty} \frac{\|X_{t+1} - X^*\|_F}{\|X_t - X^*\|_F^2} = \kappa,$$

where κ is a constant.

QUIC Results: Data from Biology Applications



Inderjit Dhillon Computer Science & Mathematics UT Austin Proximal Newton Methods for Large-Scale Machine Learning

Quic & Dirty

Active Subspace Selection for Dirty Statistical Models

Inderjit Dhillon Computer Science & Mathematics UT Austin Proximal Newton Methods for Large-Scale Machine Learning

4 3 b

Exploiting Model Structure

$$\min_{X} \{ \underbrace{\text{Loss}(X; \text{Data})}_{\text{problem structure}} + \underbrace{\text{Regularization}(X)}_{\text{model structure}} \}$$

- If $h(X) = ||X||_1$, the solution is expected to be sparse.
- Active Variable Selection:



Inderjit Dhillon Computer Science & Mathematics UT Austin Proximal Newton Methods for Large-Scale Machine Learning

▲□ → ▲ □ → ▲ □ →

Matrix Completion Problems

$$\min_{X} \{ g(X) + \lambda \|X\|_* \}$$

- $||X||_*$: a nuclear norm regularization to promote the low rank structure.
- Active variable selection ⇒ Active subspace selection (identify the low-dimensional subspace)



Inderjit Dhillon Computer Science & Mathematics UT Austin

Proximal Newton Methods for Large-Scale Machine Learning

Active Subspace Selection

• Given current iterate $X = U\Sigma V^T$, the "active subspace" is:

$$\mathcal{A} = \{ \mathbf{u}\mathbf{v}^{\mathsf{T}} \mid \mathbf{u}^{\mathsf{T}} X \mathbf{v} \neq 0 \text{ or } |\mathbf{u}^{\mathsf{T}} \nabla X \mathbf{v}| > \lambda \}.$$

• Step 1: Active subspace selection:

•
$$U_G \Sigma_G V_G^\top = S_\lambda(X - \nabla g(X)).$$

• $U_A := \operatorname{span}(U, U_G), V_A := \operatorname{span}(V, V_G).$

•
$$\mathcal{A} = \{ \mathbf{u}\mathbf{v}^T \mid \mathbf{u} \in U_A, \mathbf{v} \in \mathcal{V}_A \}.$$

• Step 2: Solve the reduced-sized problem:

$$\operatorname{argmin}_{S \in \mathbb{R}^{k \times k}} \bar{g}_X(U_A S V_A^{\top}) + \lambda \|S\|_*,$$

 $\bar{g}_X(\cdot)$ is the quadratic approximation of $g(\cdot)$ at current iterate X

- Iterate steps 1 and 2.
- Subspace selection can be generalized to settings where the regularizer is a "decomposable norm" at the solution

Active Subspace Selection for Decomposable Norms

A norm || · || is decomposable at x if there is a subspace T and vector
 e ∈ T such that the sub differential at x is

 $\partial \| \mathbf{x} \|_r = \{ \boldsymbol{\rho} \in \mathbb{R}^n \mid \Pi_{\mathcal{T}}(\boldsymbol{\rho}) = \mathbf{e} \text{ and } \| \Pi_{\mathcal{T}^{\perp}}(\boldsymbol{\rho}) \|_r^* \leq 1 \},$

- Examples: ℓ_1 norm, nuclear norm, group lasso, ...
- Active subspace selection:
 - Divide the space into "fixed" subspace and "free subspace":

$$\mathcal{S}_{\mathsf{free}} := [\mathcal{T}(\boldsymbol{\theta})] \cup [\mathcal{T}(\mathsf{prox}_{\lambda_r}(\mathcal{G}))], \ \ \mathcal{S}_{\mathsf{fixed}} = \mathcal{S}_{\mathsf{free}}^{\perp},$$

where $\mathcal{T}(\cdot)$ is the support of the decomposable norm.

• Solve the reduced sized problem.

Active Subspace Selection for Dirty Models

• Superposition-structured models or "dirty models":

$$\min_{\{\boldsymbol{\theta}^{(r)}\}_{r=1}^{k}} \bigg\{ g\bigg(\sum_{r} \boldsymbol{\theta}^{(r)}\bigg) + \sum_{r} \lambda_{r} h^{(r)}(\boldsymbol{\theta}^{(r)}) \bigg\},\$$

• Examples:

- Sparse + low-rank graphical model learning
- Multi-task learning (sparse + group-sparse or sparse + low-rank).
- Select active subspace for each $\theta^{(r)}$.

Proximal Newton for Dirty Models (QUIC & DIRTY)

For t = 0, 1, ...

- **1** Compute $\bar{\theta} = \sum_{r=1}^{k} \theta^{(r)}$
- 2 Compute the current gradient $G = \nabla g(\bar{\theta})$
- Select free set by

$$\mathcal{S}_{\mathsf{free}}^{(r)} := [\mathcal{T}(\boldsymbol{\theta}^{(r)})] \cup [\mathcal{T}(\mathsf{prox}_{\lambda_r}^{(r)}(\mathcal{G}))], \ \ \mathcal{S}_{\mathsf{fixed}}^{(r)} = \mathcal{S}_{\mathsf{free}}^{(r)^{\perp}},$$

• For $r = 1, \ldots, k$

• Solve the subproblem

$$\Delta^{(r)} = \operatorname*{argmin}_{\boldsymbol{d} \in \mathbb{R}^n} \langle \boldsymbol{d}, \boldsymbol{G} + \sum_{t \neq r} \boldsymbol{H} \Delta^{(t)} \rangle + \frac{1}{2} \boldsymbol{d}^T \boldsymbol{H} \boldsymbol{d} + \lambda_r \| \boldsymbol{\theta}^{(r)} + \boldsymbol{d} \|_r.$$

- 4 同 2 4 日 2 4 日 2 4

So Line Search: use an Armijo-rule based step-size selection to get α s.t. $\theta = \theta + \alpha \Delta$ sufficient decrease the objective function value.

Theorem: Quadratic Convergence Rate for QUIC & DIRTY

Suppose $g(\cdot)$ is strongly convex and the quadratic subproblem at each iteration is solved exactly, QUIC & DIRTY converges to the global optimum with an asymptotic quadratic convergence rate:

$$\lim_{t \to \infty} \frac{\|X_{t+1} - X^*\|_F}{\|X_t - X^*\|_F^2} = \kappa,$$

where κ is a constant.

Results: Data from Biology Applications



 $\min_{S,L:L\succeq 0, S-L\succ 0} -\log \det(S-L) + \langle S-L,\Sigma\rangle + \lambda_S \|S\|_1 + \lambda_L \operatorname{tr}(L).$



Inderjit Dhillon Computer Science & Mathematics UT Austin Proximal Newton Methods for Large-Scale Machine Learning

Application: Multi-task Learning

• Multi-task (multi-class) learning with sparse + group sparse structure.

$$\sum_{r=1}^{k} \ell(y^{(r)}, X^{(r)}(S^{(r)} + B^{(r)})) + \lambda_{S} \|S\|_{1} + \lambda_{B} \|B\|_{1,2},$$

dataset	number of	Dirty Models (sparse + group sparse)			Other Models	
	training data	Quic & Dirty	proximal gradient	ADMM	Lasso	Group Lasso
USPS	100	7.47% / 0.75s	7.49% / 10.8s	7.47% / 4.5s	10.27%	8.36%
	400	2.5% / 1.55s	2.5% / 35.8	2.5% / 11.0s	4.87%	2.93%
RCV1	1000	18.45% / 23.1s	18.49% / 430.8s	18.5% / 259s	22.67%	20.8%
	5000	10.27% / 87s	10.27% / 2254s	10.27% / 1191s	13.67%	12.25%

A B > A B >

BIG & QUIC

Sparse Inverse Covariance Estimation with 1 Million Random Variables

Inderjit Dhillon Computer Science & Mathematics UT Austin Proximal Newton Methods for Large-Scale Machine Learning



Inderjit Dhillon Computer Science & Mathematics UT Austin Proximal Newton Methods for Large-Scale Machine Learning

Difficulties in Scaling QUIC

- Coordinate descent requires X_t and X_t^{-1} ,
 - needs $O(p^2)$ storage
 - needs O(mp) computation per sweep, where $m = \|X_t\|_0$
- Line search (compute determinant of a big sparse matrix).
 - needs $O(p^2)$ storage
 - needs $O(p^3)$ computation

ヨッ イヨッ イヨッ

Line Search

- Given sparse matrix $A = X_t + \alpha D$, we need to
 - Check its positive definiteness.
 - Ompute log det(A).
- Cholesky factorization in QUIC requires $O(p^3)$ computation.

• Time complexity: $T_{CG} = O(mT)$, where T is number of CG iterations.

伺 と く ヨ と く ヨ と …

Difficulties in Scaling QUIC

- Coordinate descent requires X_t and X_t^{-1} ,
 - needs $O(p^2)$ storage
 - needs O(mp) computation per sweep, where $m = \|X_t\|_0$
- Line search (compute determinant of a big sparse matrix).
 - needs O(p) storage
 - needs O(mp) computation

ヨッ イヨッ イヨッ

Back-of-the-envelope calculations

- Co-ordinate descent requires *m* computations of $\boldsymbol{w}_i^T D \boldsymbol{w}_j$.
- To solve a problem of size $p = 10^4$ (avg degree=10), QUIC takes 12 minutes
- Extrapolate: to solve a problem with $p = 10^6$, QUIC would take $12 \text{min} \times 10^4 = 83.33$ days (2.6 days using 32 cores).
- p^2 memory means 8 Terabyte main memory for $p = 10^6$.
- Naive solution that has O(m) memory requirement:
 - Compute *w_i*, *w_j* by solving two linear systems by CG,
 - O(T_{CG}) computation per coordinate update
 - $O(mT_{CG})$ computation per sweep.
 - Would need 8,333 days to solve the problem for $p = 10^6$ (260 days using 32 cores).

・ 同 ト ・ ヨ ト ・ ヨ ト

Back-of-the-envelope calculations

- Co-ordinate descent requires *m* computations of $\boldsymbol{w}_i^T D \boldsymbol{w}_j$.
- To solve a problem of size $p = 10^4$ (avg degree=10), QUIC takes 12 minutes
- Extrapolate: to solve a problem with $p = 10^6$, QUIC would take $12 \text{min} \times 10^4 = 83.33$ days (2.6 days using 32 cores).
- p^2 memory means 8 Terabyte main memory for $p = 10^6$.
- Naive solution that has O(m) memory requirement:
 - Compute *w_i*, *w_j* by solving two linear systems by CG,
 - $O(T_{CG})$ computation per coordinate update
 - $O(mT_{CG})$ computation per sweep.
 - Would need 8,333 days to solve the problem for $p = 10^6$ (260 days using 32 cores).

イロト 不得 トイヨト イヨト 二日

• BIGQUIC solves problems of size $p = 10^6$ in one day on 32 cores.

- Assume we can store M columns of W in memory.
- Coordinate descent update (i, j): compute $\boldsymbol{w}_i^T D \boldsymbol{w}_j$.
- If w_i, w_j are not in memory: recompute by CG:

 $X \boldsymbol{w}_i = \boldsymbol{e}_i$: $O(T_{CG})$ time.

- Assume we can store M columns of W in memory.
- Coordinate descent update (i, j): compute $\boldsymbol{w}_i^T D \boldsymbol{w}_j$.
- If w_i, w_j are not in memory: recompute by CG:





4 3 5 4 3 5

- Assume we can store M columns of W in memory.
- Coordinate descent update (i, j): compute $\boldsymbol{w}_i^T D \boldsymbol{w}_j$.
- If $\boldsymbol{w}_i, \boldsymbol{w}_j$ are not in memory: recompute by CG:

 $X \boldsymbol{w}_i = \boldsymbol{e}_i$: $O(mT_{CG})$ time.





- Assume we can store M columns of W in memory.
- Coordinate descent update (i, j): compute $\boldsymbol{w}_i^T D \boldsymbol{w}_j$.
- If w_i, w_j are not in memory: recompute by CG:

 $X \boldsymbol{w}_i = \boldsymbol{e}_i$: $O(mT_{CG})$ time.





- 4 E 6 4 E 6

- Assume we can store M columns of W in memory.
- Coordinate descent update (i, j): compute $\boldsymbol{w}_i^T D \boldsymbol{w}_j$.
- If *w_i*, *w_j* are not in memory: recompute by CG:
 X w_i = *e_i*: *O(mT_{CG})* time.





Coordinate Updates - ideal case

- Want to find update sequence that minimizes number of cache misses: probably NP Hard.
- Our strategy: update variables block by block (assume k = p/M blocks).



Ideal Case

O(p) column computations



Random Case

O(kp) column computations

Inderjit Dhillon Computer Science & Mathematics UT Austin

Proximal Newton Methods for Large-Scale Machine Learning

General case: block diagonal + sparse

• If the block partition is not perfect:

extra column computations can be characterized by boundary nodes.

• Given a partition $\{S_1, \ldots, S_k\}$, we define boundary nodes as

$$B(S_q) \equiv \{j \mid j \in S_q \text{ and } \exists i \in S_z, z
eq q \text{ s.t. } F_{ij} = 1\},$$

where F is adjacency matrix of the free set.


Graph Clustering Algorithm

• The number of columns to be computed in one sweep is

$$p+\sum_{q}|B(S_{q})|.$$

• Can be upper bounded by

$$|p + \sum_{q} |B(S_q)| \leq p + \sum_{z \neq q} \sum_{i \in S_z, j \in S_q} F_{ij}.$$

• Minimize the right hand side \rightarrow minimize off-diagonal entries. Use Graph Clustering (METIS or Graclus) to find the partition.

Size of boundary nodes in real datasets

ER dataset with p = 692.



Partition by clustering compute 775 columns

< ∃ > < ∃ >

э



 $O(kpT_{CG}) \rightarrow O(pT_{CG})$ flops per sweep.

BIGQUIC

- Block co-ordinate descent with clustering,
 - needs $O(m + p^2/k)$ storage
 - needs O(mp) computation per sweep, where $m = \|X_t\|_0$
- Line search (compute determinant of a big sparse matrix).
 - needs O(p) storage
 - needs O(mp) computation

Theorem (Hsieh, Sustik, Dhillon & Ravikumar, NIPS, 2011)

QUIC converges quadratically to the global optimum, that is for some constant 0 < κ < 1:

$$\lim_{t\to\infty}\frac{\|X_{t+1}-X^*\|_F}{\|X_t-X^*\|_F^2}=\kappa.$$

- Recall $W = X^{-1}$.
- When each \boldsymbol{w}_i is computed by CG $(X \boldsymbol{w}_i = \boldsymbol{e}_i)$:
 - The gradient $\nabla_{ij}g(X) = S_{ij} W_{ij}$ on free set can be computed once and stored in memory.
 - Hessian ($w_i^T D w_j$ in coordinate updates) needs to be repeatedly computed.
- To reduce the time overhead, Hessian should be computed approximately.

直 と く ヨ と く ヨ と

- Recall $W = X^{-1}$.
- When each \boldsymbol{w}_i is computed by CG $(X \boldsymbol{w}_i = \boldsymbol{e}_i)$:
 - The gradient $\nabla_{ij}g(X) = S_{ij} W_{ij}$ on free set can be computed once and stored in memory.
 - Hessian ($w_i^T D w_j$ in coordinate updates) needs to be repeatedly computed.
- To reduce the time overhead, Hessian should be computed approximately.

Theorem (Hsieh, Sustik, Dhillon, Ravikumar & Poldrack, 2013)

If $\nabla g(X)$ is computed exactly and $\bar{\eta}I \succeq \hat{H}_t \succeq \eta I$ for some constant $\bar{\eta}, \eta > 0$ at every Newton iteration, then BIGQUIC converges to the global optimum.

(4月) (日) (日) 日

BIGQUIC Convergence Analysis

- To have super-linear convergence rate, we require the Hessian computation to be more and more accurate as X_t approaches X*.
- Measure the accuracy of Hessian computation by R = [r₁..., r_p], where r_i = X w_i e_i.
- The optimality of X_t can be measured by

$$\nabla_{ij}^{S} f(X) = \begin{cases} \nabla_{ij} g(X) + \operatorname{sign}(X_{ij}) \lambda & \text{if } X_{ij} \neq 0, \\ \operatorname{sign}(\nabla_{ij} g(X)) \max(|\nabla_{ij} g(X)| - \lambda, 0) & \text{if } X_{ij} = 0. \end{cases}$$

Theorem (Hsieh, Sustik, Dhillon, Ravikumar & Poldrack, 2013)

In BIGQUIC, if residual matrix $||R_t|| = O(||\nabla^S f(X_t)||^\ell)$ for some $0 < \ell \le 1$ as $t \to \infty$, then $||X_{t+1} - X^*|| = O(||X_t - X^*||^{1+\ell})$ as $t \to \infty$.

(4月) (4日) (4日) 日

BIGQUIC Convergence Analysis

- To have super-linear convergence rate, we require the Hessian computation to be more and more accurate as X_t approaches X^{*}.
- Measure the accuracy of Hessian computation by $R = [\mathbf{r}_1 \dots, \mathbf{r}_p]$, where $\mathbf{r}_i = X \mathbf{w}_i \mathbf{e}_i$.
- The optimality of X_t can be measured by

$$\nabla_{ij}^{S} f(X) = \begin{cases} \nabla_{ij} g(X) + \operatorname{sign}(X_{ij}) \lambda & \text{if } X_{ij} \neq 0, \\ \operatorname{sign}(\nabla_{ij} g(X)) \max(|\nabla_{ij} g(X)| - \lambda, 0) & \text{if } X_{ij} = 0. \end{cases}$$

Theorem (Hsieh, Sustik, Dhillon, Ravikumar & Poldrack, 2013)

In BIGQUIC, if residual matrix $||R_t|| = O(||\nabla^S f(X_t)||^\ell)$ for some $0 < \ell \le 1$ as $t \to \infty$, then $||X_{t+1} - X^*|| = O(||X_t - X^*||^{1+\ell})$ as $t \to \infty$.

When ||*R_t*|| = *O*(||∇^S*f*(*X_t*)||), BIGQUIC has asymptotic quadratic convergence rate.

Inderjit Dhillon Computer Science & Mathematics UT Austin Proximal Newton Methods for Large-Scale Machine Learning

Experimental results (scalability)



Figure: BIGQUIC can solve one million dimensional problems.

Inderjit Dhillon Computer Science & Mathematics UT Austin Proximal Newton Methods for Large-Scale Machine Learning

同 ト イ ヨ ト イ ヨ ト

Experimental results

 $\operatorname{BIGQUIC}$ is faster even for medium size problems.



Figure: Comparison on FMRI data with a p = 20000 subset (maximum dimension that previous methods can handle).

Results on FMRI dataset

- 228,483 voxels, 518 time points.
- $\lambda = 0.6 \Longrightarrow$ average degree 8, BIGQUIC took 5 hours.
 - $\lambda = 0.5 \Longrightarrow$ average degree 38, $\operatorname{BIGQUIC}$ took 21 hours.
- Findings:
 - Voxels with large degree were generally found in the gray matter.
 - Can detect meaningful brain modules by modularity clustering.



(b) (4) (3) (4)

Conclusions

- **Proximal Newton method**: a second order method for optimizing a smooth convex function plus a non-smooth regularizer.
- How to develop an efficient proximal Newton method?

Exploit problem structure and geometry of problem

- For sparse inverse covariance estimation problem, our proposed algorithm (BIGQUIC) can solve a 1-million dimensional problem in 1-day on a single 32-core machine.
- Main ingredients for BIGQUIC:
 - Fast coordinate descent solver for computing Newton direction.
 - Avoid un-important updates by active subspace selection.
 - Memory efficient coordinate descent updates.
- QUIC & DIRTY: extension to "Dirty" Statistical Models.

伺 と く ヨ と く ヨ と

Future Work

- "Effective" order of Newton-like methods for specific problems?
- Multi-core Computation? "Wild" Co-ordinate Descent?
- Prioritized Scheduling (Co-ordinate Descent)?
- Suitability for Distributed for Out-of-core Computing?
- Develop scalable & reliable software
- No dearth of interesting applications!

References

[1] C. J. Hsieh, I. S. Dhillon, P. Ravikumar, S. Becker, and P. Olsen. *QUIC & DIRTY: A Quadratic Approximation Approach for Dirty Statistical Models*. NIPS, 2014.

[2] C. J. Hsieh, M. Sustik, I. S. Dhillon, P. Ravikumar, and R. Poldrack. *BIG & QUIC: Sparse inverse covariance estimation for a million variables*. NIPS (oral presentation), 2013.

[3] C. J. Hsieh, M. Sustik, I. S. Dhillon, and P. Ravikumar. *Sparse Inverse Covariance Matrix Estimation using Quadratic Approximation*. NIPS, 2011.

[4] C. J. Hsieh, I. S. Dhillon, P. Ravikumar, A. Banerjee. *A Divide-and-Conquer Procedure for Sparse Inverse Covariance Estimation*. NIPS, 2012.

[5] P. A. Olsen, F. Oztoprak, J. Nocedal, and S. J. Rennie. *Newton-Like Methods for Sparse Inverse Covariance Estimation*. Optimization Online, 2012.

[6] O. Banerjee, L. El Ghaoui, and A. d'Aspremont *Model Selection Through Sparse Maximum Likelihood Estimation for Multivariate Gaussian or Binary Data*. JMLR, 2008.

[7] J. Friedman, T. Hastie, and R. Tibshirani. *Sparse inverse covariance estimation with the graphical lasso*. Biostatistics, 2008.

[8] L. Li and K.-C. Toh. *An inexact interior point method for I1-reguarlized sparse covariance selection*. Mathematical Programming Computation, 2010.

[9] K. Scheinberg, S. Ma, and D. Glodfarb. *Sparse inverse covariance selection via alternating linearization methods.* NIPS, 2010.

[10] K. Scheinberg and I. Rish. *Learning sparse Gaussian Markov networks using a greedy coordinate ascent approach*. Machine Learning and Knowledge Discovery in Databases, 2010.

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ● ● ● ● ● ●