

Minimizing composite functions

$$\underset{x}{\text{minimize}} f(x) := g(x) + h(x)$$

- ▶ g and h are closed, convex functions
- ▶ g is continuously differentiable, and its gradient ∇g is Lipschitz continuous
- ▶ h is not necessarily everywhere differentiable, but its *proximal mapping* can be evaluated efficiently

Sparse inverse covariance estimation

- ▶ $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbf{R}^n$ are *i.i.d.* samples from a Gaussian MRF:
 $\Pr(\mathbf{x}; \Theta) \propto \exp(\mathbf{x}^T \Theta \mathbf{x} / 2 - \log \det(\Theta))$
- ▶ We form the sample covariance matrix $\hat{\Sigma} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)} \mathbf{x}^{(i)T}$ and seek a sparse maximum likelihood estimate of Θ :

$$\underset{\Theta \in \mathbf{R}^{n \times n}}{\text{minimize}} -\log \det(\Theta) + \text{tr}(\hat{\Sigma} \Theta) + \lambda \|\text{vec}(\Theta)\|_1$$

Proximal Newton-type methods

Main idea: use a local quadratic model (in lieu of a simple quadratic model) to account for the curvature of g :

$$\Delta \mathbf{x}_k := \underset{d}{\text{arg min}} \nabla g(\mathbf{x}_k)^T d + \frac{1}{2} d^T H_k d + h(\mathbf{x}_k + d).$$

approx. Hessian term

$\Delta \mathbf{x}_k$ can be expressed as

$$H_k \Delta \mathbf{x}_k \in \underbrace{-\nabla g(\mathbf{x}_k)}_{\text{forward/explicit gradient}} - \underbrace{\partial h(\mathbf{x}_k + \Delta \mathbf{x}_k)}_{\text{backward/implicit subgradient}}.$$

A generic proximal Newton-type method

Algorithm 1 A generic proximal Newton-type method

Require: starting point $\mathbf{x}_0 \in \text{dom } f$

- 1: **repeat**
- 2: Choose an approximation to the Hessian H_k .
- 3: Solve the subproblem for a search direction:
 $\Delta \mathbf{x}_k \leftarrow \underset{d}{\text{arg min}} \nabla g(\mathbf{x}_k)^T d + \frac{1}{2} d^T H_k d + h(\mathbf{x}_k + d).$
- 4: Select t_k with a backtracking line search.
- 5: Update: $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + t_k \Delta \mathbf{x}_k.$
- 6: **until** stopping conditions are satisfied.

Convergence of proximal Newton-type methods

Global convergence:

- ▶ smallest eigenvalue of H_k 's bounded away from zero

Local quadratic convergence (prox-Newton method):

- ▶ g is locally strongly convex
- ▶ $\nabla^2 g$ is locally Lipschitz continuous

Local superlinear convergence (prox-quasi-Newton methods):

- ▶ assumptions for quadratic convergence
- ▶ eigenvalues of H_k 's bounded and H_k 's satisfy:

$$\lim_{k \rightarrow \infty} \frac{\| (H_k - \nabla^2 g(\mathbf{x}^*)) (\mathbf{x}_{k+1} - \mathbf{x}_k) \|_2}{\| \mathbf{x}_{k+1} - \mathbf{x}_k \|_2} = 0.$$

Inexact proximal Newton-type methods

Main idea: no need to solve the subproblem exactly only need a good enough search direction.

$$\Delta \mathbf{x}_k \approx \underset{d}{\text{arg min}} \nabla g(\mathbf{x}_k)^T d + \frac{1}{2} d^T H_k d + h(\mathbf{x}_k + d).$$

- ▶ We solve the subproblem approximately with an iterative method, terminating (sometimes very) early
- ▶ number of iterations may increase, but computational expense per iteration is smaller
- ▶ many practical implementations use inexact search directions

Another idea: choose H_k so the subproblem is easy to solve.

Early stopping conditions

Intuition: solve the subproblem almost exactly when

- ▶ \mathbf{x}_k is close to the optimal solution
- ▶ H_{k-1} captures the curvature of g

Typical stopping condition:

$$\| \underbrace{\nabla g(\mathbf{x}_k) + H_k \Delta \mathbf{x}_k + \partial h(\mathbf{x}_k + \Delta \mathbf{x}_k)}_{\text{optimality of subproblem solution}} \| \leq \eta_k \underbrace{\| \mathbf{G}_f(\mathbf{x}_k) \|}_{\text{optimality of } \mathbf{x}_k}$$

choose η_k based on how good the quadratic model is:

$$\eta_k \sim \frac{\| \nabla g_{k-1}(\mathbf{x}_{k-1}) + H_{k-1} \Delta \mathbf{x}_{k-1} - \nabla g(\mathbf{x}_k) \|_2}{\| \nabla g(\mathbf{x}_{k-1}) \|_2}$$

Convergence of the inexact prox-Newton method

Local linear convergence (inexact prox-Newton method):

- ▶ g is locally strongly convex
- ▶ $\nabla^2 g$ is locally Lipschitz continuous
- ▶ η_k is smaller than some $\bar{\eta}$

Local superlinear convergence (...):

- ▶ assumptions for linear convergence
- ▶ η_k decays to zero (e.g. under our choice of forcing term)

Sparse inverse covariance estimation

- ▶ We have samples $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \in \mathbf{R}^n$ from a Gaussian MRF.
- ▶ We form the sample covariance matrix $\hat{\Sigma} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)} \mathbf{x}^{(i)T}$ and seek a sparse maximum likelihood estimate of Θ :

$$\underset{\Theta \in \mathbf{R}^{n \times n}}{\text{minimize}} -\log \det(\Theta) + \text{tr}(\hat{\Sigma} \Theta) + \lambda \|\text{vec}(\Theta)\|_1$$

Datasets:

- ▶ Estrogen: a gene expression dataset consisting of 682 probe sets collected from 158 patients
- ▶ Leukemia: another gene expression dataset consisting of 1255 genes from 72 patients

Q: How do inexact search directions affect convergence?

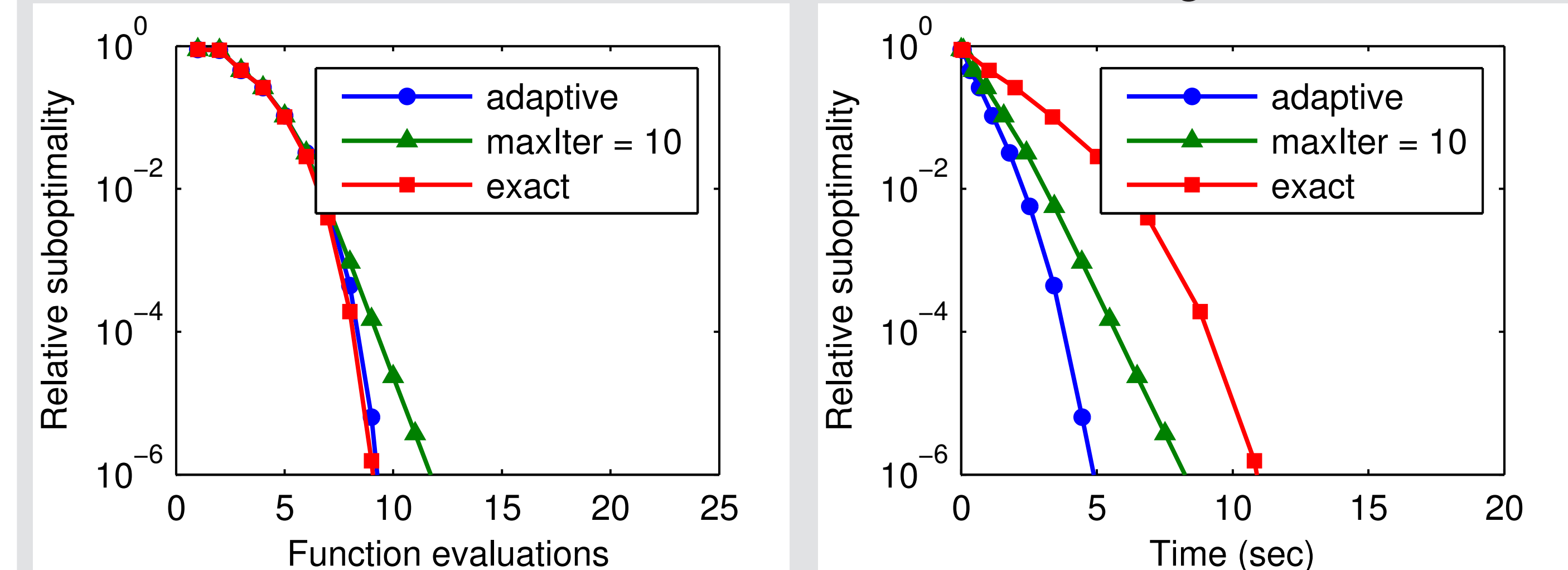


Figure : Estrogen dataset

Summary

Proximal Newton-type methods

- ▶ converge rapidly near the optimal solution, and can produce a solution of high accuracy.
- ▶ are insensitive to the condition number of the sublevel sets of the objective.
- ▶ are suited to problems where g , ∇g is expensive to evaluate compared to h , prox_h .