

Composite Self-concordant Minimization

Quoc Tran Dinh, Anastasios Kyrillidis and Volkan Cevher

Laboratory for Information and Inference Systems (LIONS)

Abstract

We propose an algorithmic framework for convex minimization problems of composite functions with two terms: a self-concordant part and a possibly nonsmooth regularization part. Our method is a new proximal Newton algorithm with local quadratic convergence rate. As a specific problem instance, we consider sparse precision matrix estimation problems in graph learning. Via a careful dual formulation and a novel analytic step-size selection, we instantiate an algorithm within our framework for graph learning that avoids *Cholesky decompositions* and *matrix inversions*, making it attractive for parallel and distributed implementations.

Composite minimization

$$(P) \quad \min_{\mathbf{x} \in \mathbb{R}^n} \{F(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x})\}$$

f
convex and smooth **g**
convex and possibly nonsmooth

Motivation

Problem (P) covers many practical problems:

- Unconstrained basic LASSO/Basic pursuit denoising problems
- Graphical model selections / latent variable graphical model selection
- Poisson imaging reconstruction / Heteroscedastic LASSO
- Low-rank approximation
- Clustering

Problem (P) has specific separable structure

- splitting methods
- augmented Lagrangian methods

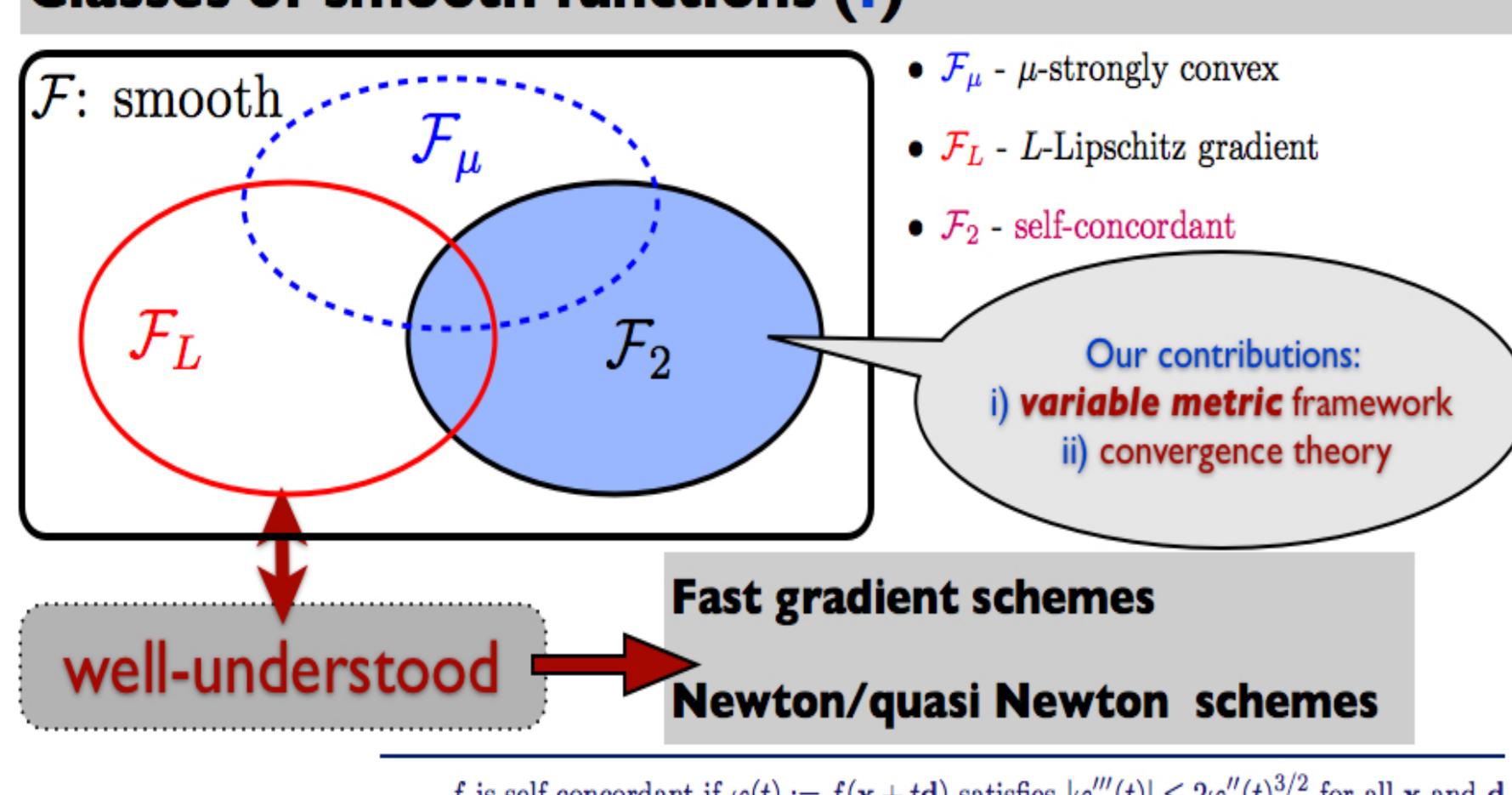
g: ℓ_1 -norm, nuclear norm or indicator functions

Composite self-concordant minimization

$$(P) \quad \min_{\mathbf{x} \in \mathbb{R}^n} \{F(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x})\}$$

f
convex and self-concordant **g**
convex and possibly nonsmooth

Classes of smooth functions (f)

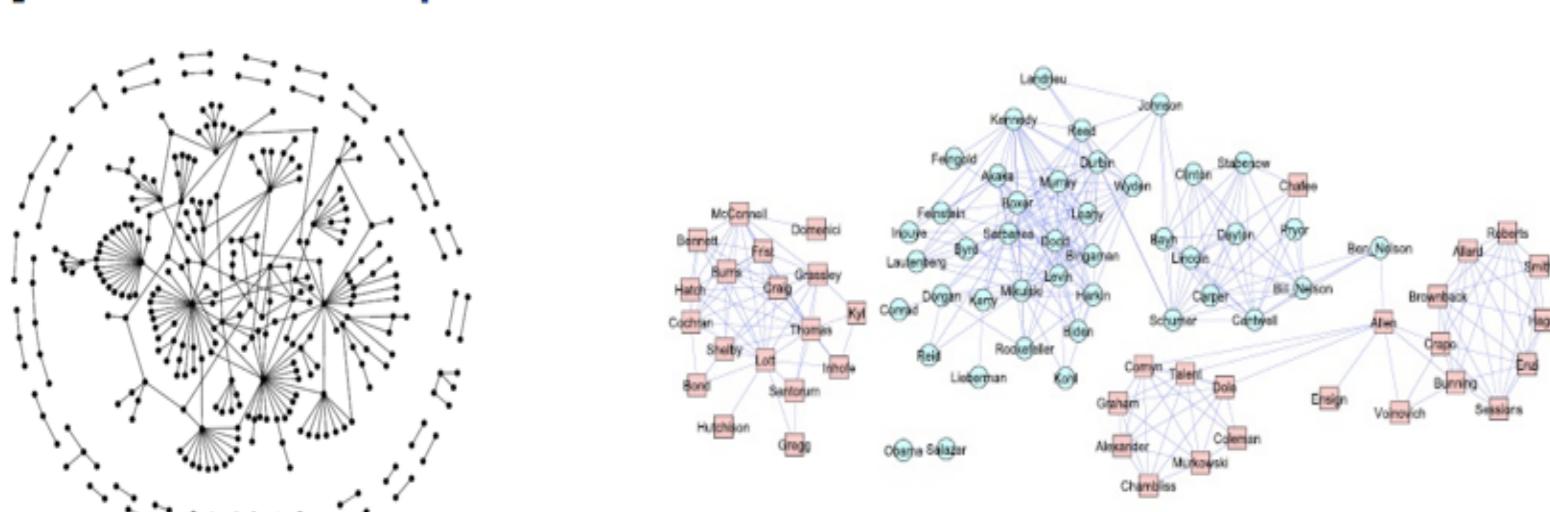


f is self-concordant if $\varphi(t) = f(x+td)$ satisfies $|\varphi'''(t)| \leq 2\varphi''(t)^{3/2}$ for all x and d .

Example: Log-determinant for LMIs

Optimization problem
$\min_{\Theta > 0} \left\{ -\log \det(\Theta) + \text{trace}(\Sigma \Theta) + \rho \ \vec{\Theta}\ _1 \right\}$

Application: Graphical model selection



Log-barrier for linear/quadratic inequalities

Poisson imaging reconstruction

$$x^* \in \arg\min \left\{ \sum_{i=1}^m a_i^T x - \sum_{i=1}^m y_i \log(a_i^T x + b_i) + g(x) \right\}$$

$f(x)$

Basic pursuit denoising problem (BPDN): Barrier formulation

$$x_t^* = \arg\min_x \left\{ -t \log(\sigma^2 - \|Ax - y\|_2^2) + g(x) \right\}$$

$=: f(x)$

LASSO problem with unknown variance

$$x^* \equiv (\phi^*, \gamma^*) = \arg\min_{\phi, \gamma} \left\{ -\log(\gamma) + \frac{1}{2n} \|\gamma y - X\phi\|_2^2 + \lambda \|\phi\|_1 \right\}$$

$=: f(x)$

Prior state-of-the-art

Proximal point scheme with variable metric [Bonnans, 1993]

Given \mathbf{x}^0 , generate a sequence $\{\mathbf{x}^k\}_{k \geq 0}$ such that

$$\mathbf{x}^{k+1} = \text{prox}_{H_k}(\mathbf{x}^k - H_k^{-1} \nabla f(\mathbf{x}^k))$$

where H_k is symmetric positive definite

Special cases

- If $H_k := \gamma_k I$ then we obtain a proximal-gradient scheme [Nesterov2007, Beck2009]
- If $H_k \approx \nabla^2 f(\mathbf{x}^k)$ then we obtain a proximal-quasi Newton scheme [Becker2012]
- If $H_k = \nabla^2 f(\mathbf{x}^k)$ then we obtain a proximal-Newton scheme [Lee2012]
- If $\text{prox}_{H_k} \equiv \text{Id}$ then we obtain a classical variable metric descent scheme

Connection to splitting

- PPA(H) is a special case of splitting forward-backward methods

Analytic complexity

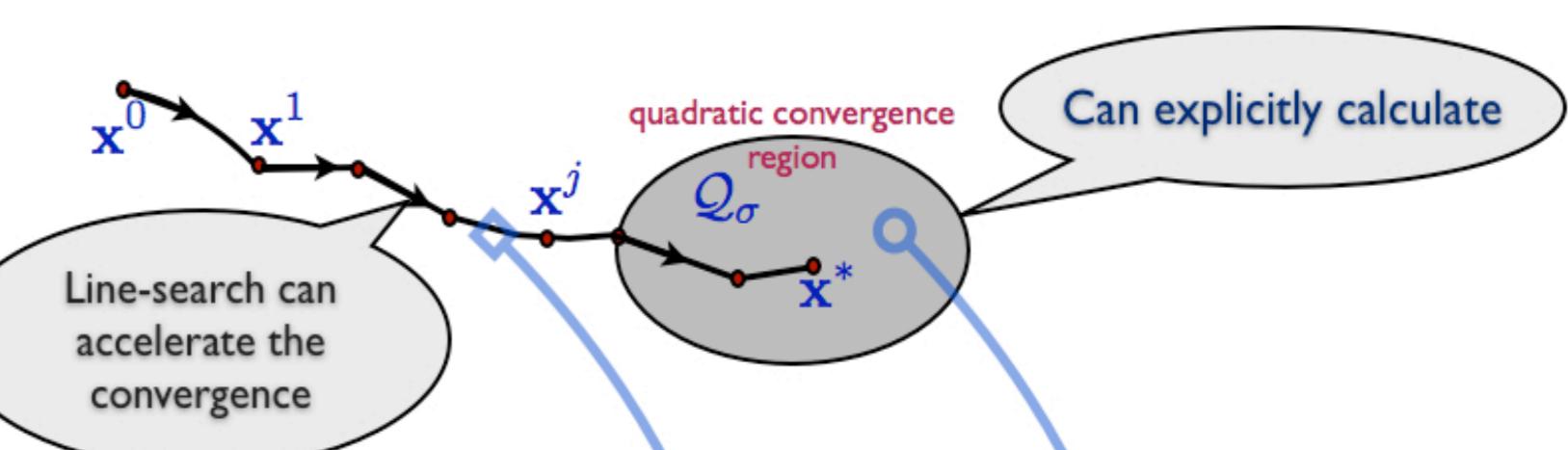
Proximal Newton decrement:

$$\lambda_k := \|\mathbf{d}_k\|_{\mathbf{x}^k}$$

Quadratic convergence region: Let $\sigma := (5 - \sqrt{17})/4 \approx 0.219224$

$$\mathcal{Q}_\sigma := \{\mathbf{x}^k \mid \lambda_k \leq \sigma\}$$

Illustration the convergence behavior



Worst-case complexity

$$\# \text{iterations} = \left\lceil \frac{F(\mathbf{x}^0) - F(\mathbf{x}^*)}{0.021} \right\rceil + O(\ln \ln(\epsilon))$$

Application to graphical model selection

The underlying convex optimization problem

$$\min_{\Theta > 0} \left\{ -\log \det(\Theta) + \text{trace}(\Sigma \Theta) + \rho \|\vec{\Theta}\|_1 \right\}$$

Gradient and Hessian (large-scale, special structure)

- Gradient of f : $\nabla f(\mathbf{x}) = \text{vec}(\Sigma - \Theta^{-1})$.
- Hessian of f : $\nabla^2 f(\mathbf{x}) = \Theta^{-1} \otimes \Theta^{-1}$

Dual approach for solving subproblem (SP)

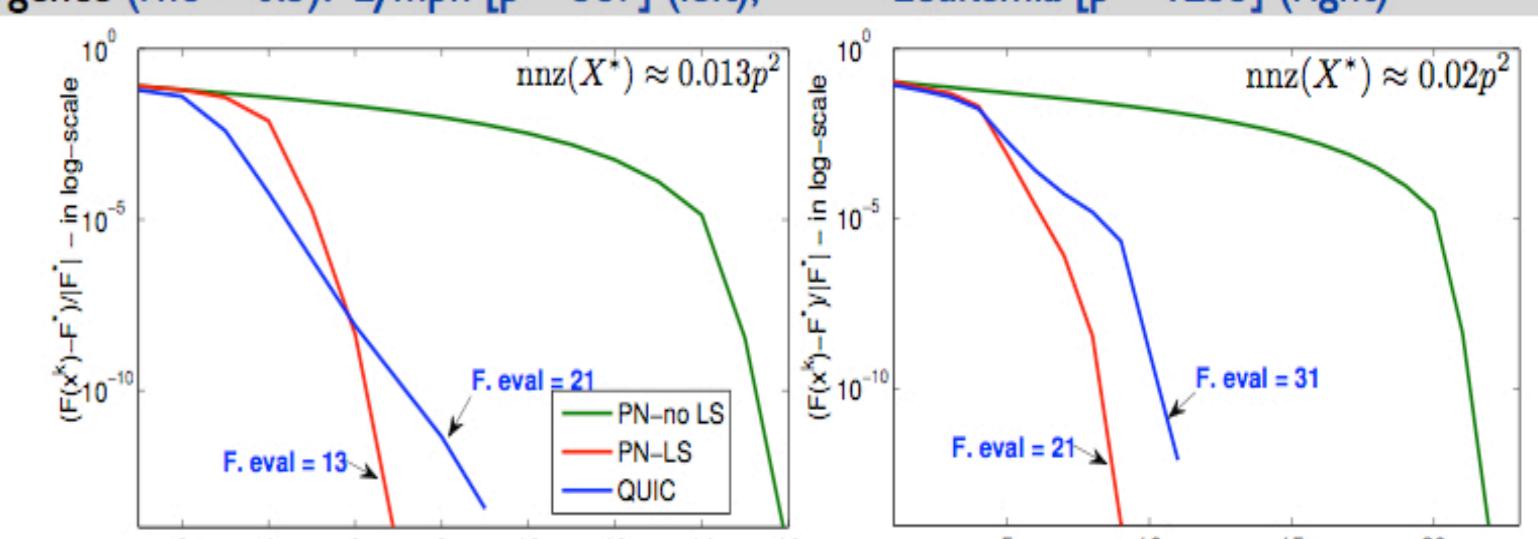
Unconstrained LASSO problems	No Cholesky decomposition and matrix inversion
Primal subproblem	Dual subproblem (SPGL)
$\min_{\Delta} \left\{ \frac{1}{2} \text{trace}((\Theta_\Delta^{-1} \Delta)^2) + \text{trace}(\mathbf{R}_\Delta \Delta) + \rho \ \vec{\Delta}\ _1 \right\}$	$\min_{\ \vec{\Delta}\ _\infty \leq 1} \left\{ \frac{1}{2} \text{trace}((\Theta_\Delta^{-1} \Delta)^2) + \text{trace}(\mathbf{Q}_\Delta \Delta) \right\}$
$\mathbf{R}_\Delta := \Sigma - 2\Theta_\Delta^{-1}$	$\mathbf{Q}_\Delta := \rho^{-1}[\Theta_\Delta \Sigma \Theta_\Delta - 2\Theta_\Delta]$

How to compute proximal Newton decrement $\lambda_i := \|\mathbf{d}_i\|_{\mathbf{x}^i}$?

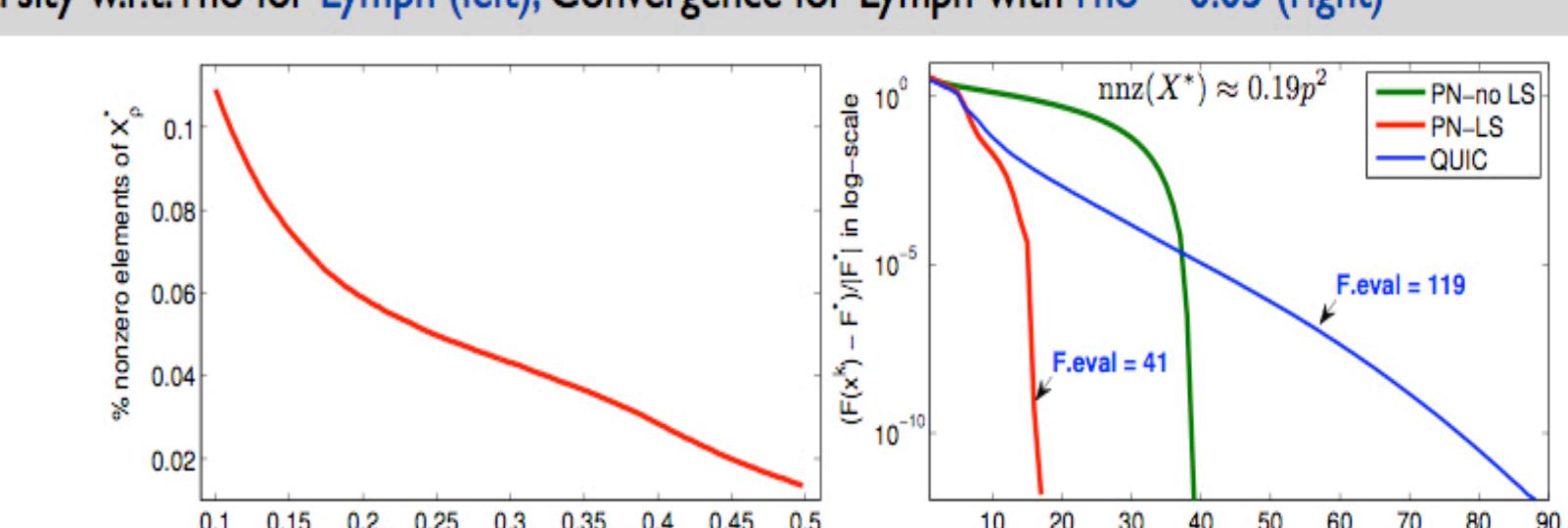
$$\lambda_i := [\rho - 2\text{trace}(\mathbf{W}_i) + \text{trace}(\mathbf{W}_i^2)]^{1/2}, \quad \mathbf{W}_i = \Theta_i(\Sigma - \rho U^*)$$

Numerical experiments on graph learning

Convergence ($\rho = 0.5$): Lymph [$p = 587$] (left), Leukemia [$p = 1255$] (right)



Sparsity w.r.t. rho for Lymph (left), Convergence for Lymph with rho = 0.05 (right)

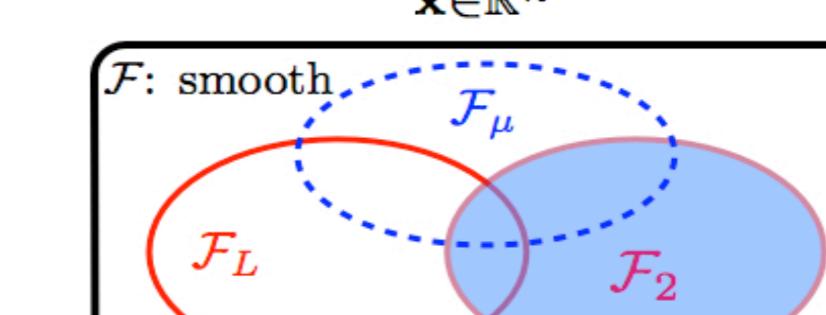


Our method vs QUIC [Hsieh2011]

- QUIC subproblem solver special block-coordinate descent
- Our subproblem solver general proximal algorithms

Conclusions

$$\min_{\mathbf{x} \in \mathbb{R}^n} \{\phi(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x})\}$$



- \mathcal{F}_μ - μ -strongly convex
- \mathcal{F}_L - L -Lipschitz gradient
- \mathcal{F}_2 - self-concordant

Highlights

- **Globalization:** a new strategy for finding step-size explicitly, motivate "forward-looking" line-search strategy
- **Search direction:** efficient (strongly convex program)
- **Local convergence:** quadratic convergence without boundedness of the Hessian analytic quadratic convergence region

- Practical contributions
 - software package has quasi-Newton, pure-Newton, gradient and fast gradient algorithms
 - leverage fast proximal solvers for $g(\mathbf{x})$ (structured norms etc.)
 - robust to subproblem solver accuracy

Software

- A open-source software package called **SCOPT** is available at <http://lions.epfl.ch/software>
- **SCOPT** stands for Self-Concordant OPTimization