

# Composite Self-concordant Minimization

UT Simons Seminar Series  
October 2nd, 2015

Anastasios Kyrillidis

*Recall*

$$\begin{array}{ll}\text{minimize} & g(\mathbf{X}) \\ \mathbf{X} \in \mathbb{R}^{n \times n} & \\ \text{subject to} & \mathbf{X} \succeq 0\end{array}$$

*Recall*

$$\begin{array}{ll} \text{minimize} & g(\mathbf{X}) \\ \mathbf{X} \in \mathbb{R}^{n \times n} & \end{array}$$

$$\text{subject to } \mathbf{X} \succeq 0$$

**Assumptions:**  $g(\mathbf{X})$  is convex, Lipschitz gradient and strongly convex.

$$\mu \mathbf{I} \preceq \nabla^2 g(\mathbf{X}) \preceq L \mathbf{I}$$

*Recall*

$$\underset{\mathbf{U} \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad g(\mathbf{U}\mathbf{U}^\top)$$

where  $r \leq n$

(since  $\mathbf{U}\mathbf{U}^\top \succeq 0$ )

**Assumptions:**  $g(\mathbf{X})$  is convex, Lipschitz gradient and strongly convex.

$$\mu\mathbf{I} \preceq \nabla^2 g(\mathbf{X}) \preceq L\mathbf{I}$$

**Recall**

$$\underset{\mathbf{U} \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad g(\mathbf{U}\mathbf{U}^\top)$$

- Provable linear convergence rate guarantees, even under non-convexity.
- First-order oracle (i.e., at most  $\nabla g(\cdot)$ ).

*“Dropping convexity for faster semi-definite optimization”, Bhojanapalli, Kyrillidis, Sanghavi, 2015*

**Assumptions:**  $g(\mathbf{X})$  is convex, Lipschitz gradient and strongly convex.

$$\mu\mathbf{I} \preceq \nabla^2 g(\mathbf{X}) \preceq L\mathbf{I}$$

*In this talk*

$$\underset{\mathbf{X} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad f(\mathbf{X}) + \lambda \cdot g(\mathbf{X})$$

## Composite convex optimization

$$\begin{array}{ll}\text{minimize} & f(\mathbf{X}) + \lambda \cdot g(\mathbf{X}) \\ \mathbf{X} \in \mathbb{R}^{n \times n}\end{array}$$

In this talk

## Composite convex optimization

*In this talk*

$$\underset{\mathbf{X} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad f(\mathbf{X}) + \lambda \cdot g(\mathbf{X})$$

- $g(\mathbf{X})$  is a convex function, but possibly non-smooth!
- $g(\mathbf{X})$  ``models`` a-priori knowledge.

## Composite convex optimization

*In this talk*

$$\underset{\mathbf{X} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad f(\mathbf{X}) + \lambda \cdot g(\mathbf{X})$$

- $g(\mathbf{X})$  is a convex function, but possibly non-smooth!
- $g(\mathbf{X})$  ``models`` a-priori knowledge.

Examples:  $g(\mathbf{X}) \equiv \|\mathbf{X}\|_1$  (for sparsity),

$g(\mathbf{X}) \equiv \|\mathbf{X}\|_*$  (for low-rankness),

⋮  
⋮

## Composite convex optimization

*In this talk*

$$\underset{\mathbf{X} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad f(\mathbf{X}) + \lambda \cdot g(\mathbf{X})$$

- $f(\mathbf{X})$  is a **convex function**, but possibly  
non-globally Lipschitz gradient and strongly convex!
- $f(\mathbf{X})$  is usually a data fidelity term.

## Composite convex optimization

*In this talk*

$$\underset{\mathbf{X} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad f(\mathbf{X}) + \lambda \cdot g(\mathbf{X})$$

- $f(\mathbf{X})$  is a convex function, but possibly non-globally Lipschitz gradient and strongly convex!
- $f(\mathbf{X})$  is usually a data fidelity term.

Examples:  $f(\mathbf{X}) \equiv \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X})\|_2^2$  (least squares, MC, PT...),  
 $f(\mathbf{X}) = -\log \det(\mathbf{X}) + \langle \mathbf{X}, \mathbf{C} \rangle$  (GRMF)  
⋮

## Composite convex optimization

*In this talk*

$$\underset{\mathbf{X} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad f(\mathbf{X}) + \lambda \cdot g(\mathbf{X})$$

- $f(\mathbf{X})$  is a convex function, but possibly non-globally Lipschitz gradient and strongly convex!
- $f(\mathbf{X})$  is usually a data fidelity term.

Assumption:  $g(\mathbf{X})$  has a ``tractable`` proximity operator:

$$\text{prox}_{\lambda g}(\mathbf{W}) = \arg \min_{\mathbf{X}} g(\mathbf{X}) + \frac{1}{2\lambda} \|\mathbf{X} - \mathbf{W}\|_F^2$$

## Composite convex optimization

*In this talk*

$$\underset{\mathbf{X} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad f(\mathbf{X}) + \lambda \cdot g(\mathbf{X})$$

- $f(\mathbf{X})$  is a convex function, but possibly non-globally Lipschitz gradient and strongly convex!
- $f(\mathbf{X})$  is usually a data fidelity term.

Assumption:  $g(\mathbf{X})$  has a ``tractable`` proximity operator:

$$\text{prox}_{\lambda \|\cdot\|_1}(\mathbf{W}) = \text{SoftThresh}(\mathbf{W}, \lambda)$$

# Note

$$\underset{\mathbf{X} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad f(\mathbf{X}) + \lambda \cdot g(\mathbf{X})$$

≡

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{x}) + \lambda \cdot g(\mathbf{x})$$

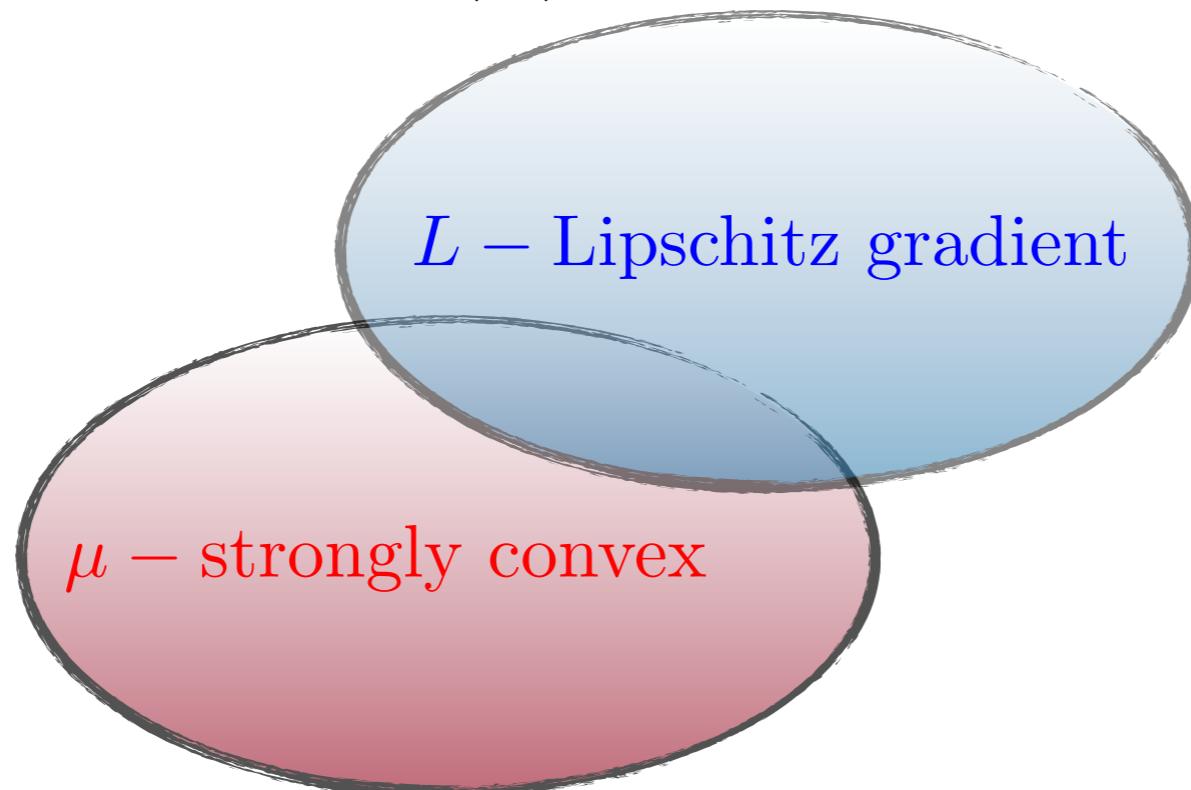
What is known for

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{x}) + \lambda \cdot g(\mathbf{x})$$

- Methods inseparable with Lipschitz gradient and strongly convex assumption:
  - Proximal first order methods [CW05, Nes07, EB92, SRB11]
  - Accelerated first order methods [BT09]
  - Smoothing techniques [Nes03]
  - Proximal Newton methods [BF12, LSS12, TKC15]
  - ADMM / ALM approaches [BPCPE12, GM12]
  - Primal-dual approaches [TC15] (and references therein)

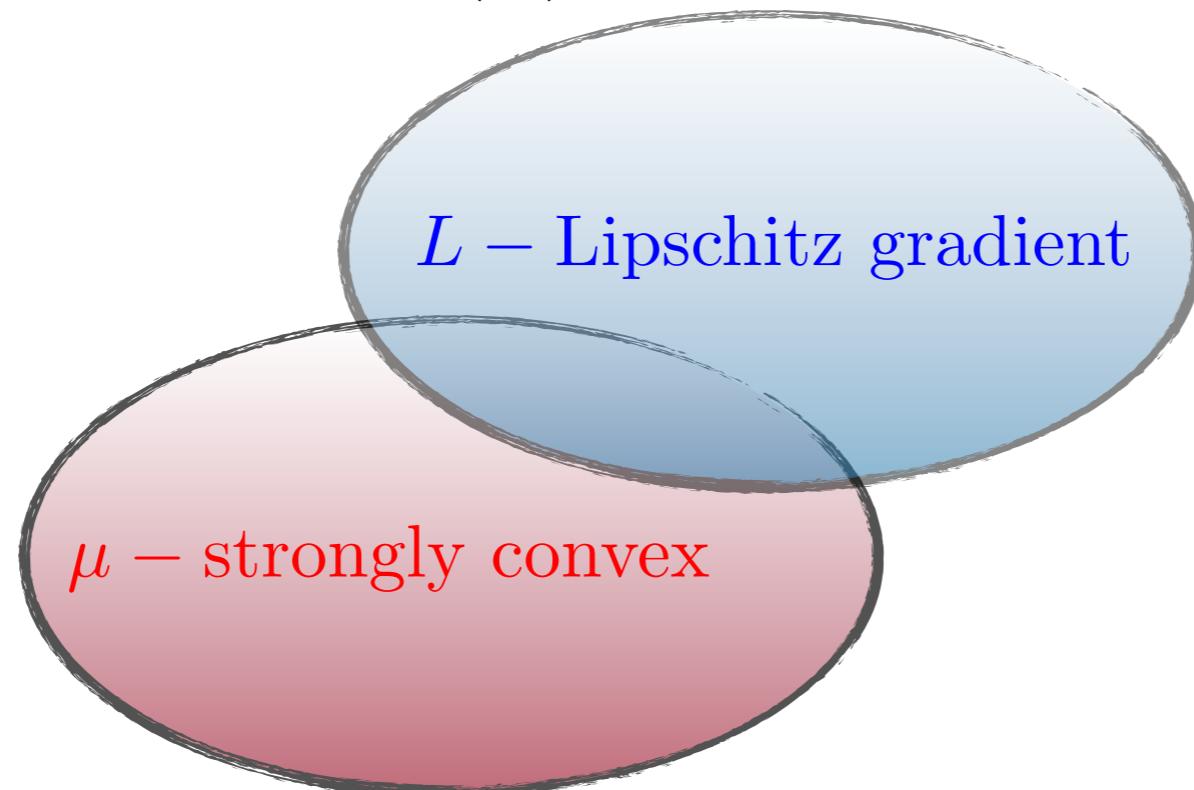
- Methods inseparable with Lipschitz gradient and strongly convex assumption:
  - Proximal first order methods [CW05, Nes07, EB92, SRB11]
  - Accelerated first order methods [BT09]
  - Smoothing techniques [Nes03]
  - Proximal Newton methods [BF12, LSS12, TKC15]
  - ADMM / ALM approaches [BPCPE12, GM12]
  - Primal-dual approaches [TC15] (and references therein)
- Self-concordant functions are not new! [NN94]
  - Known results only for smooth case.

## Differentiable $f(\mathbf{x})$



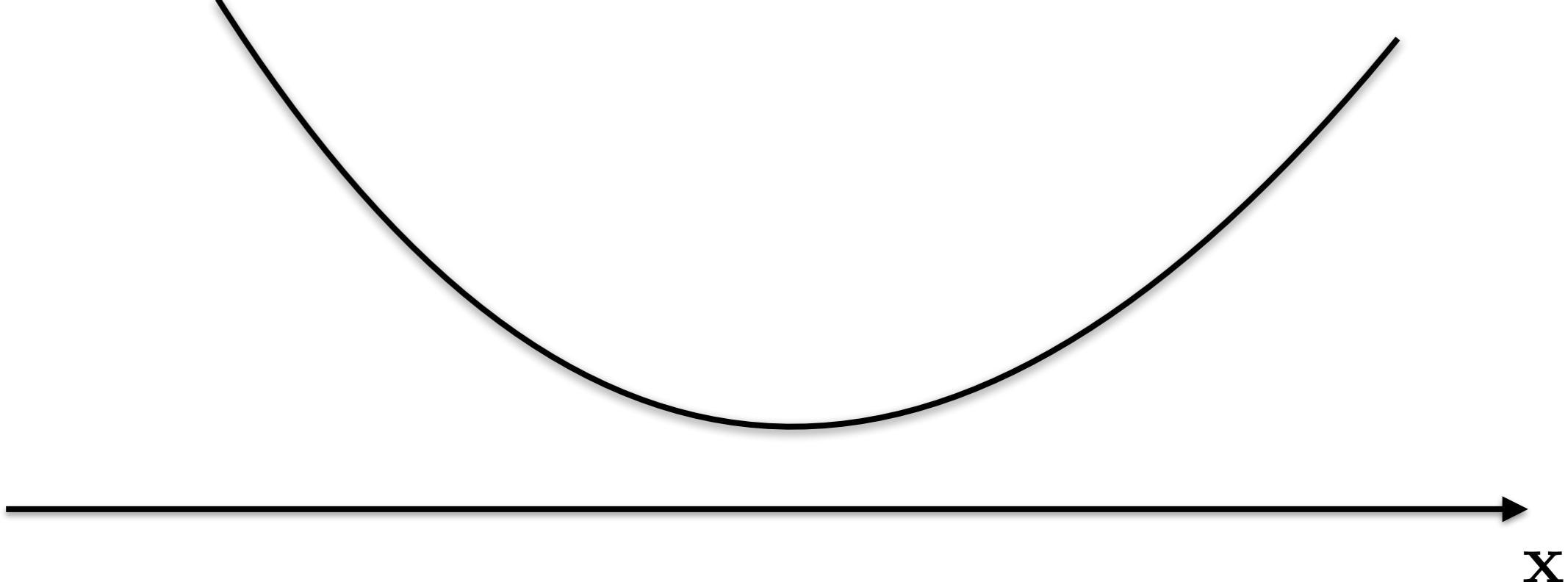
- L-Lipschitz gradient functions:  $\|\nabla f(\mathbf{x}_1) - \nabla f(\mathbf{x}_2)\| \leq L\|\mathbf{x}_1 - \mathbf{x}_2\|$
- μ - strongly convex functions:  $\mu\mathbb{I} \preceq \nabla^2 f(\mathbf{x}) \quad \mu > 0$

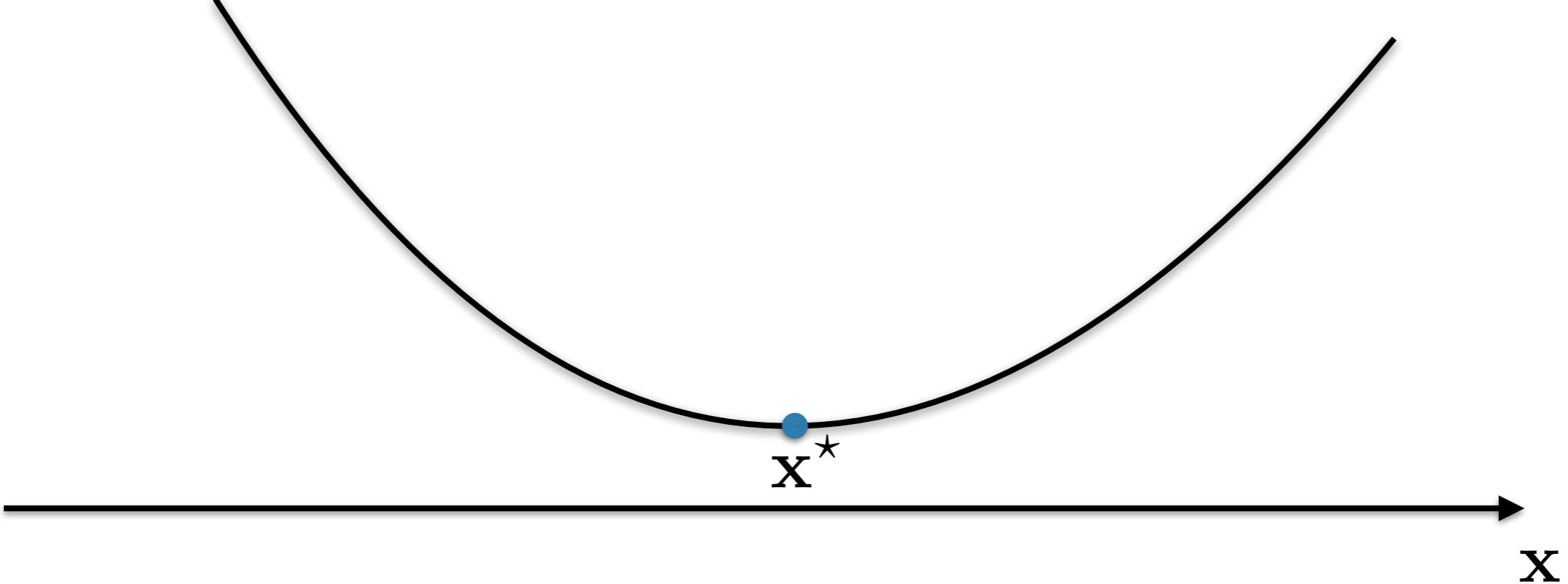
## Differentiable $f(\mathbf{x})$

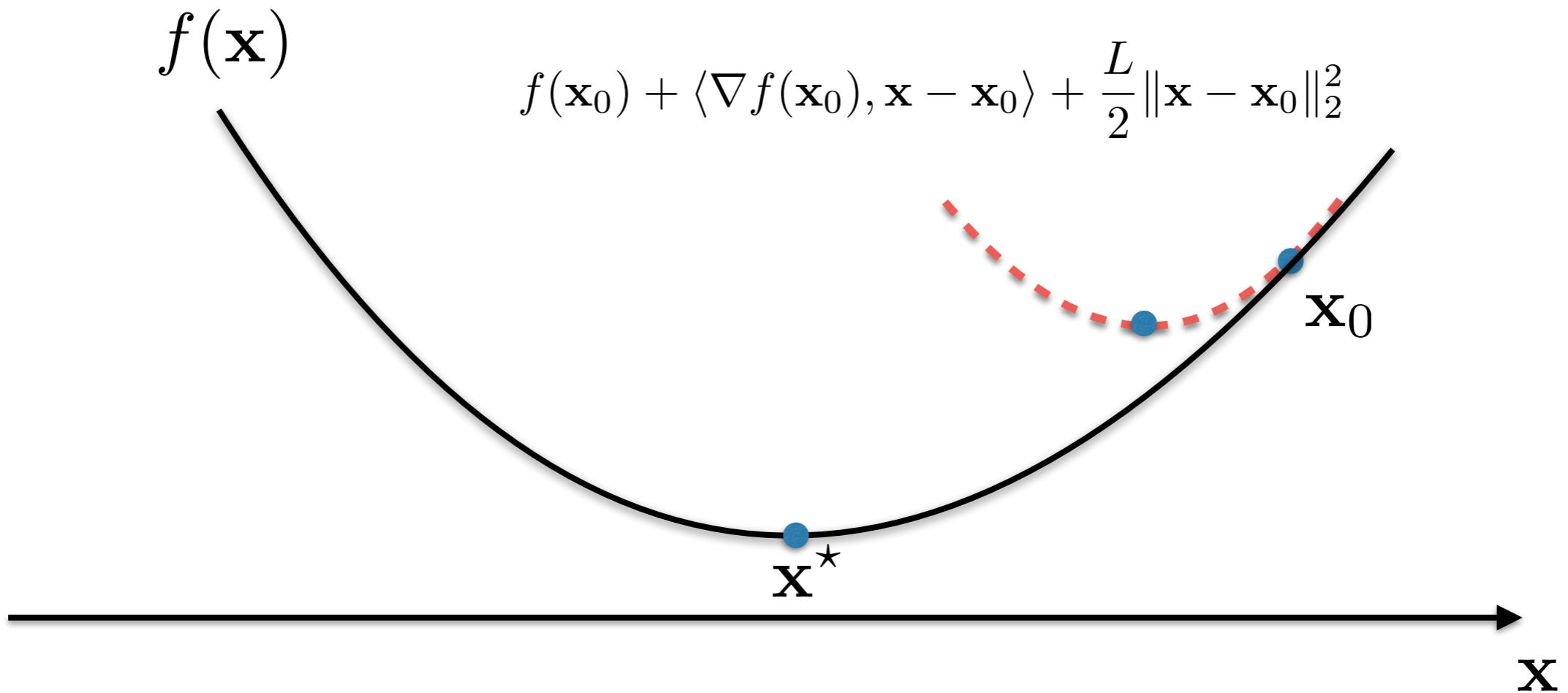


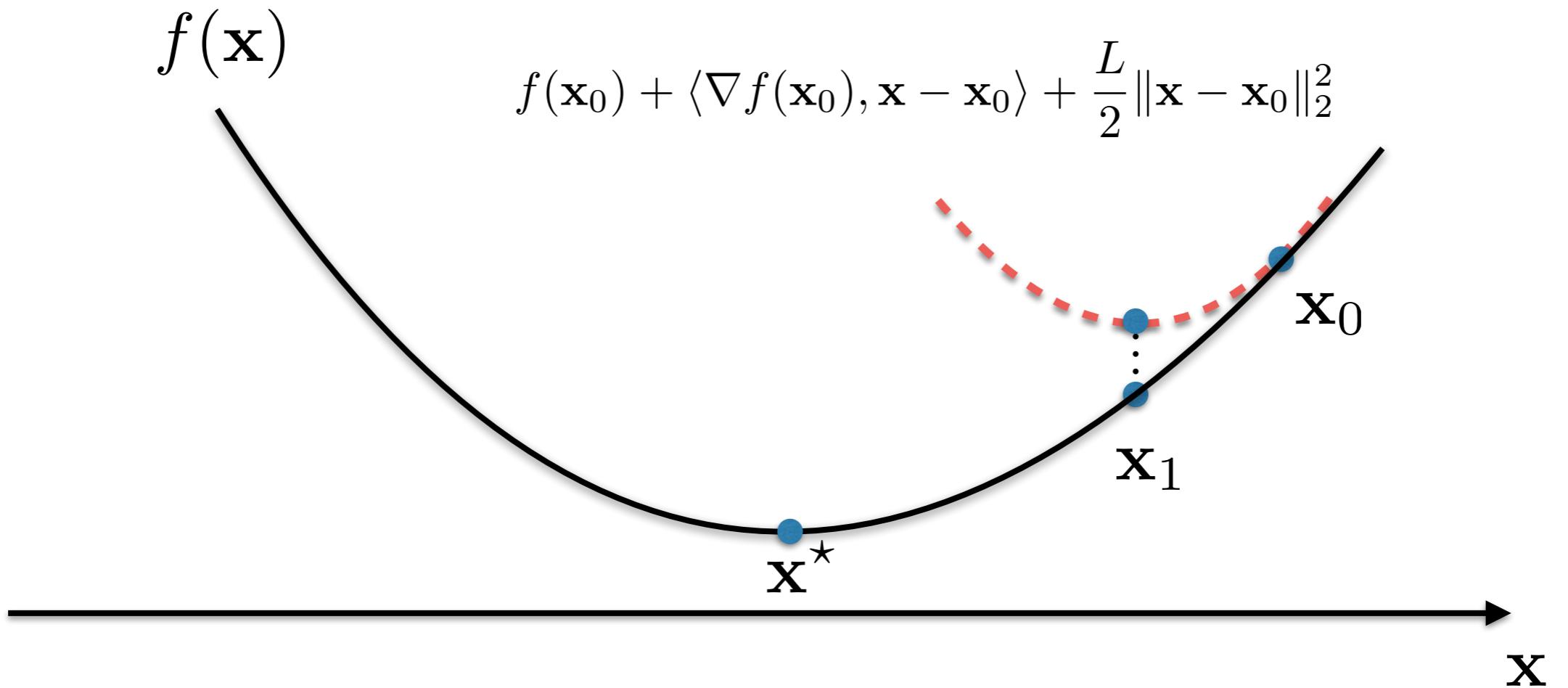
- L-Lipschitz gradient functions:  $\|\nabla f(\mathbf{x}_1) - \nabla f(\mathbf{x}_2)\| \leq L\|\mathbf{x}_1 - \mathbf{x}_2\|$
- $\mu$  - strongly convex functions:  $\mu\mathbb{I} \preceq \nabla^2 f(\mathbf{x}) \quad \mu > 0$
- Such assumptions lead to convergence guarantees, step size selection...

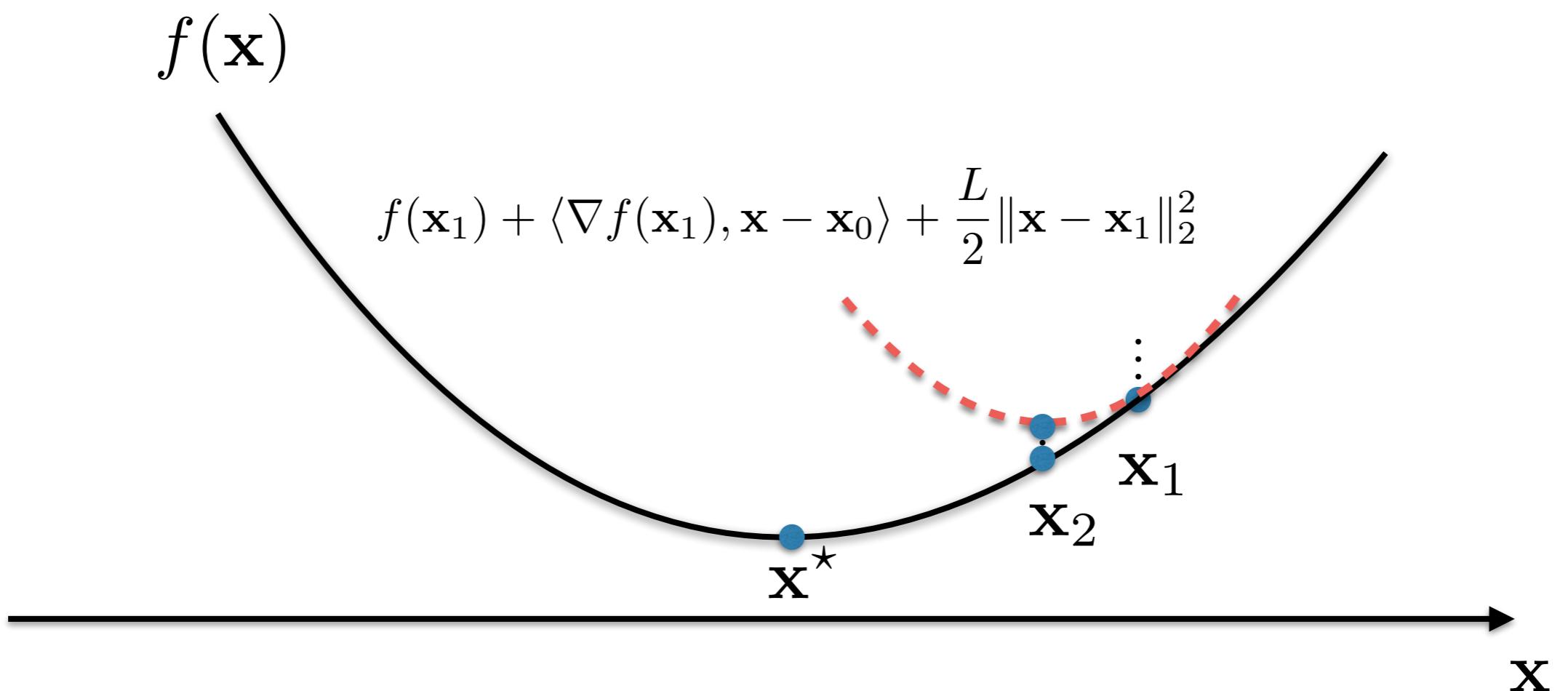


$f(\mathbf{x})$ 

$f(\mathbf{x})$ 



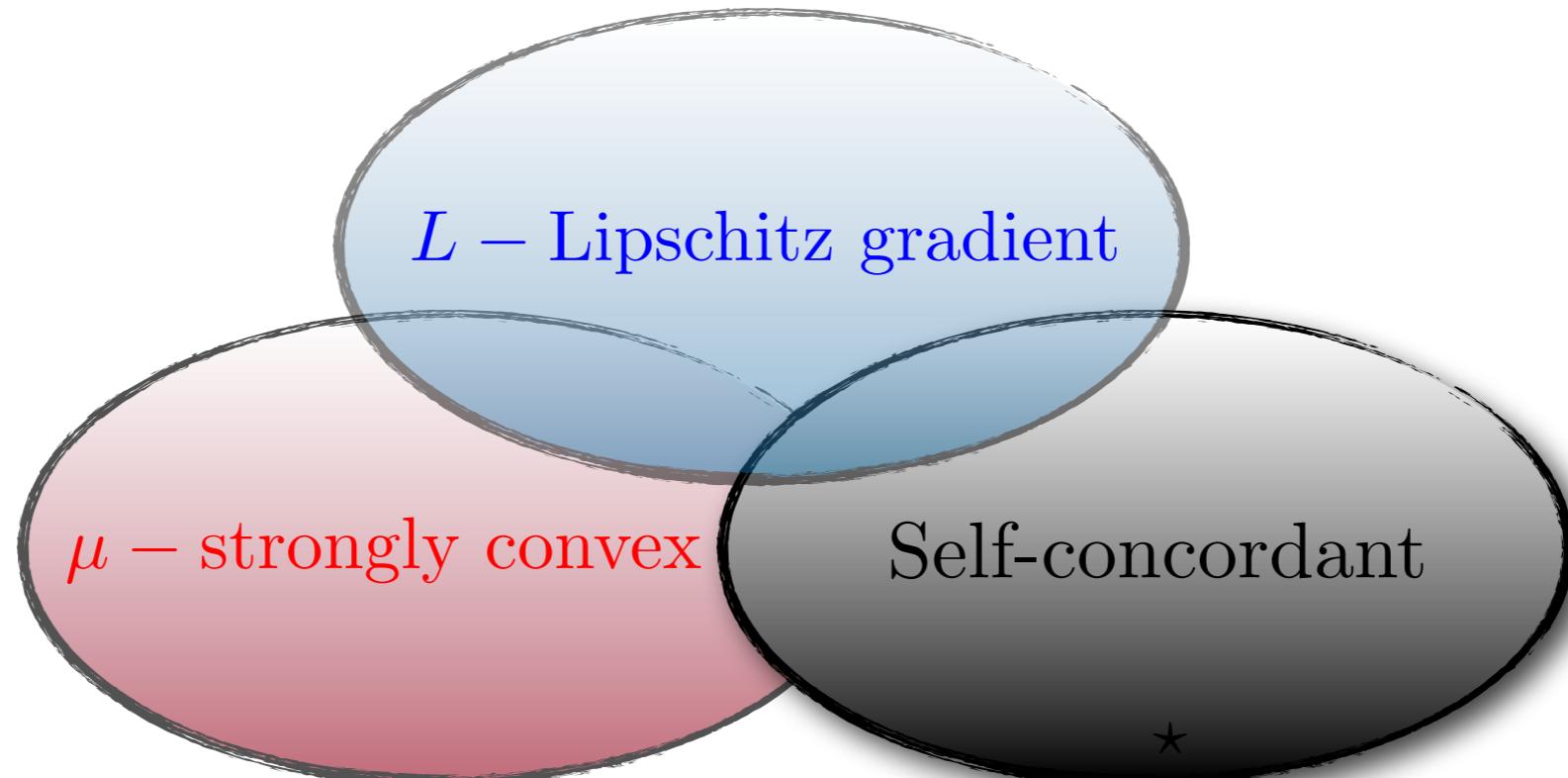




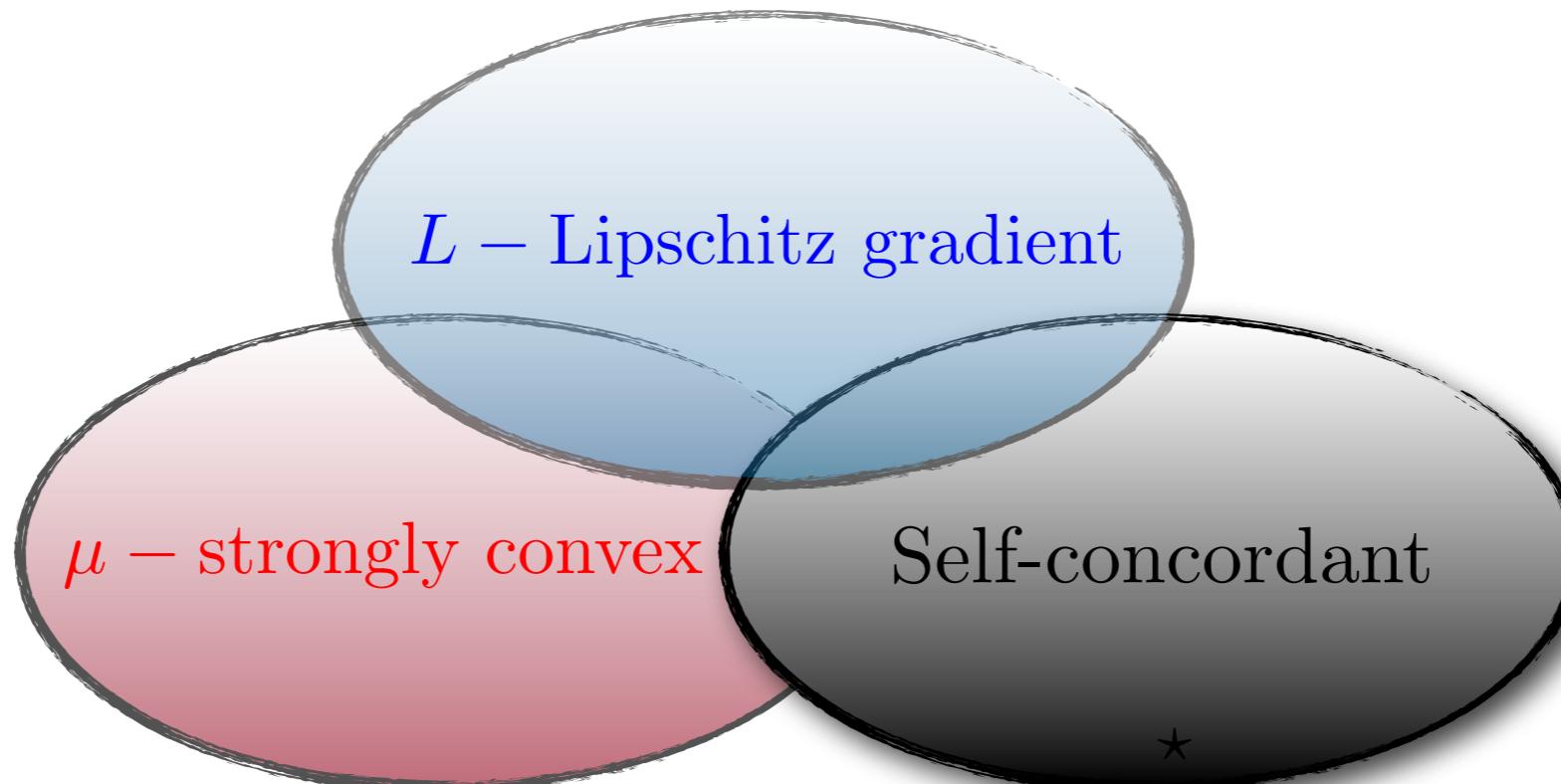
*In this talk*

# Class of self-concordant functions

Differentiable  $f(\mathbf{x})$

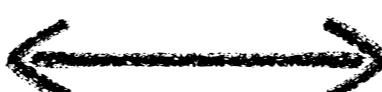


Differentiable  $f(\mathbf{x})$



$f(\mathbf{x})$  is self-concordant<sup>\*</sup>

[NN94]



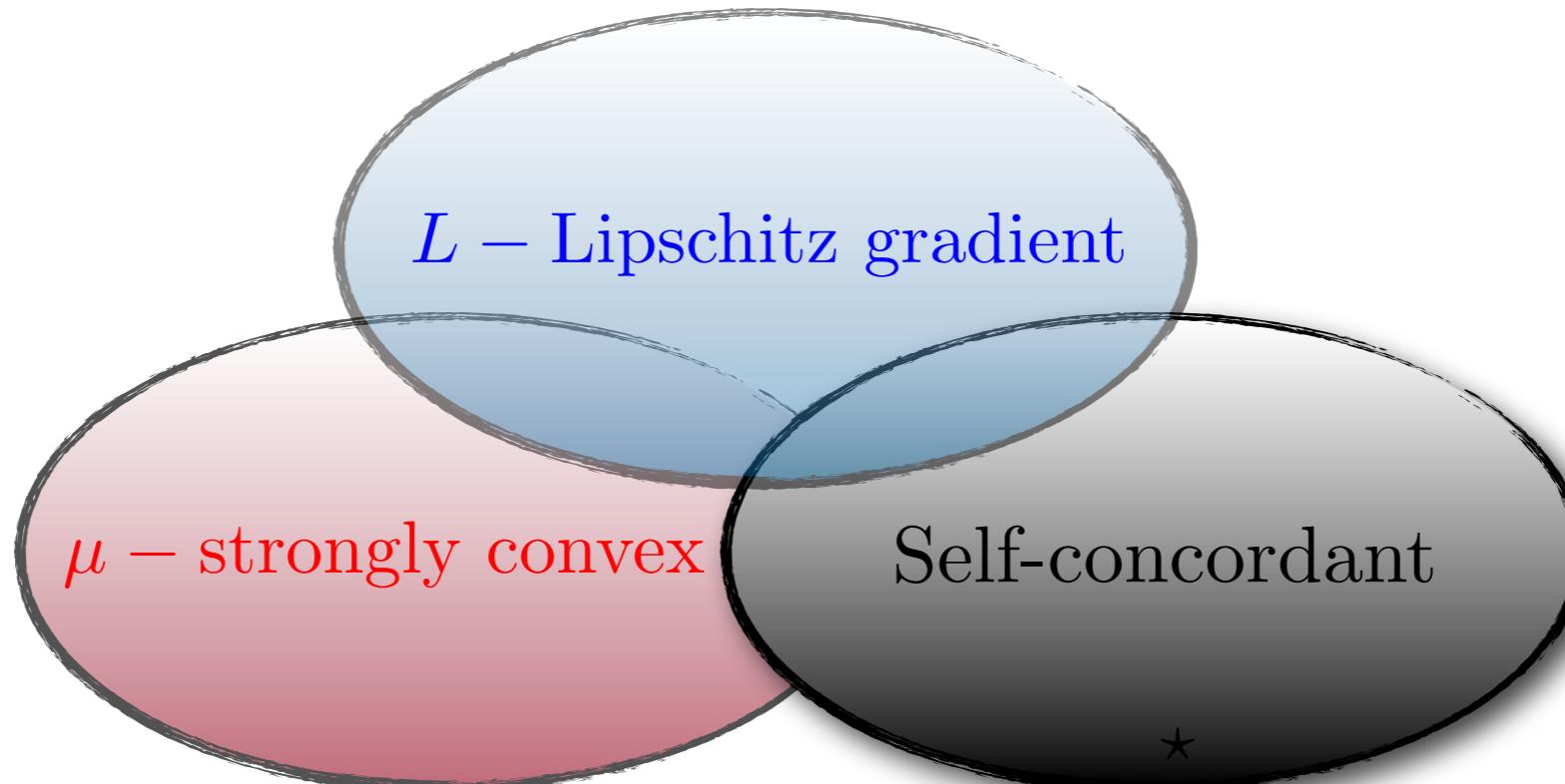
$$|\varphi'''(t)| \leq 2\varphi''(t)^{3/2}$$

$$\varphi(t) := f(\mathbf{x} + t\mathbf{d}), \quad \mathbf{x} \in \text{dom}(f), \mathbf{d} \in \mathbb{R}^n$$

---

\* Related to the affine invariant convergence guarantee of Newton method on self-concordant functions.

## Differentiable $f(\mathbf{x})$

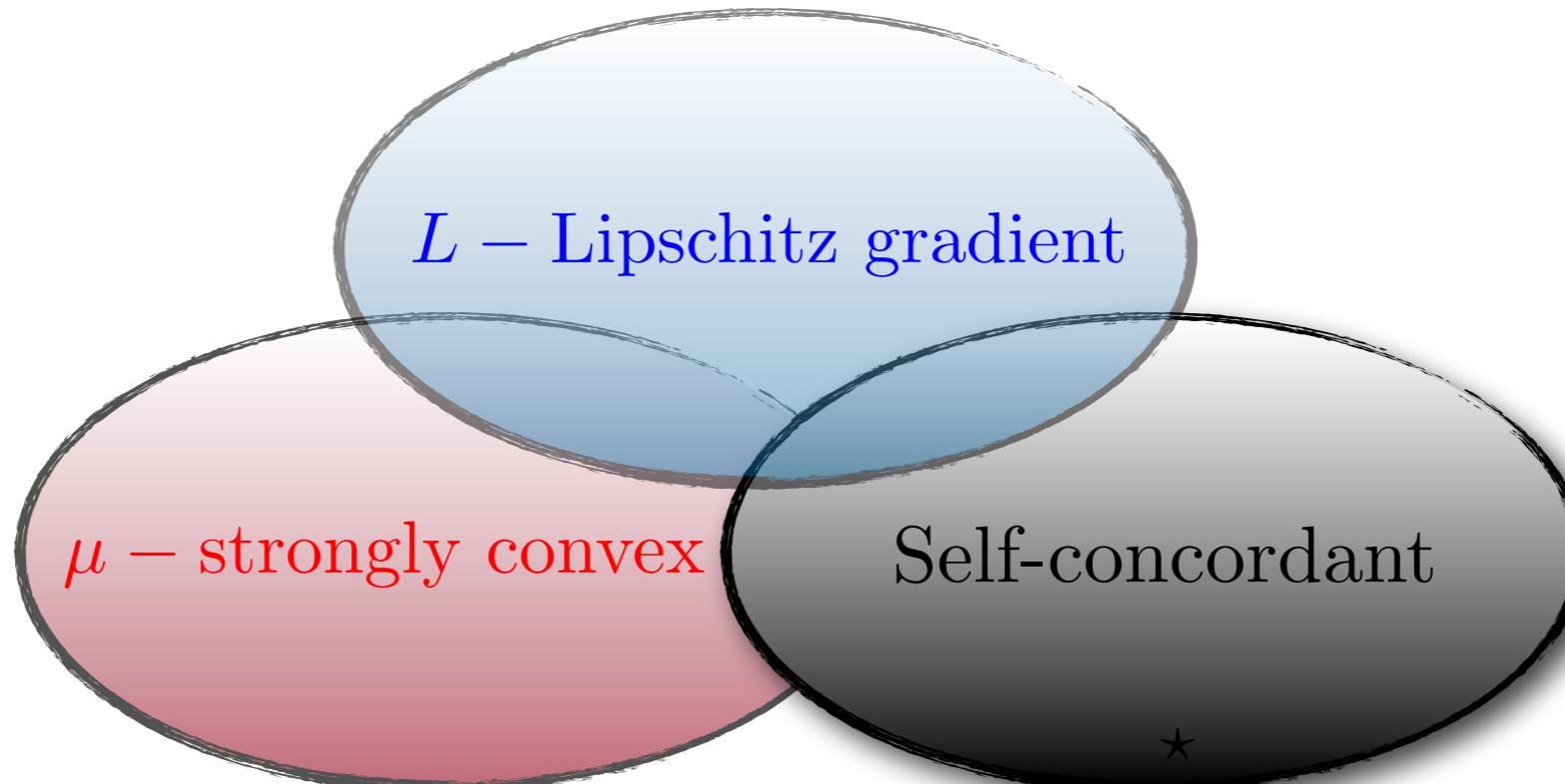


Intuition

$$\|\nabla^2 f(\mathbf{x}_1) - \nabla^2 f(\mathbf{x}_2)\| \leq M \|\mathbf{x}_1 - \mathbf{x}_2\|$$

E.g., Newton  $\|\mathbf{x}_{i+1} - \mathbf{x}^*\| \leq \frac{M \|\mathbf{x}_i - \mathbf{x}^*\|^2}{2(\mu - M \|\mathbf{x}_i - \mathbf{x}^*\|)}$

## Differentiable $f(\mathbf{x})$



Intuition

$$\|\nabla^2 f(\mathbf{x}_1) - \nabla^2 f(\mathbf{x}_2)\| \leq M \|\mathbf{x}_1 - \mathbf{x}_2\|$$

E.g., Newton  $\|\mathbf{x}_{i+1} - \mathbf{x}^*\| \leq \frac{M \|\mathbf{x}_i - \mathbf{x}^*\|^2}{2(\mu - M \|\mathbf{x}_i - \mathbf{x}^*\|)}$

- Constants change with a new basis  $f(\mathbf{Ax})$

**Self-concordance does not lead to global conditions!**

# Self-concordance does not lead to global conditions!

- Main properties of Lipschitz gradient + strongly convex functions

Lower surrogate	$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{\mu}{2}\ \mathbf{y} - \mathbf{x}\ _2^2$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Upper surrogate	$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{L}{2}\ \mathbf{y} - \mathbf{x}\ _2^2$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Hessian surrogates	$\mu\mathbb{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbb{I}$	$\mathbf{x} \in \text{dom}(f)$

# Self-concordance does not lead to global conditions!

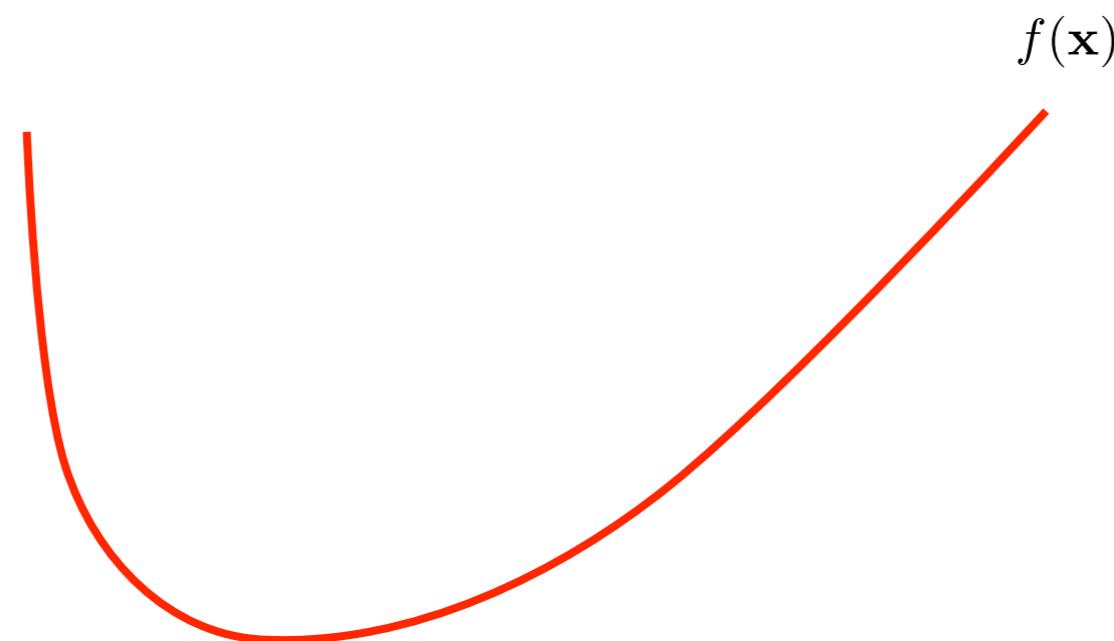
- Main properties of Lipschitz gradient + strongly convex functions

Lower surrogate	$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{\mu}{2}\ \mathbf{y} - \mathbf{x}\ _2^2$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Upper surrogate	$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{L}{2}\ \mathbf{y} - \mathbf{x}\ _2^2$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Hessian surrogates	$\mu\mathbb{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbb{I}$	$\mathbf{x} \in \text{dom}(f)$

# Self-concordance does not lead to global conditions!

- Main properties of Lipschitz gradient + strongly convex functions

Lower surrogate	$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{\mu}{2}\ \mathbf{y} - \mathbf{x}\ _2^2$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Upper surrogate	$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{L}{2}\ \mathbf{y} - \mathbf{x}\ _2^2$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Hessian surrogates	$\mu\mathbb{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbb{I}$	$\mathbf{x} \in \text{dom}(f)$

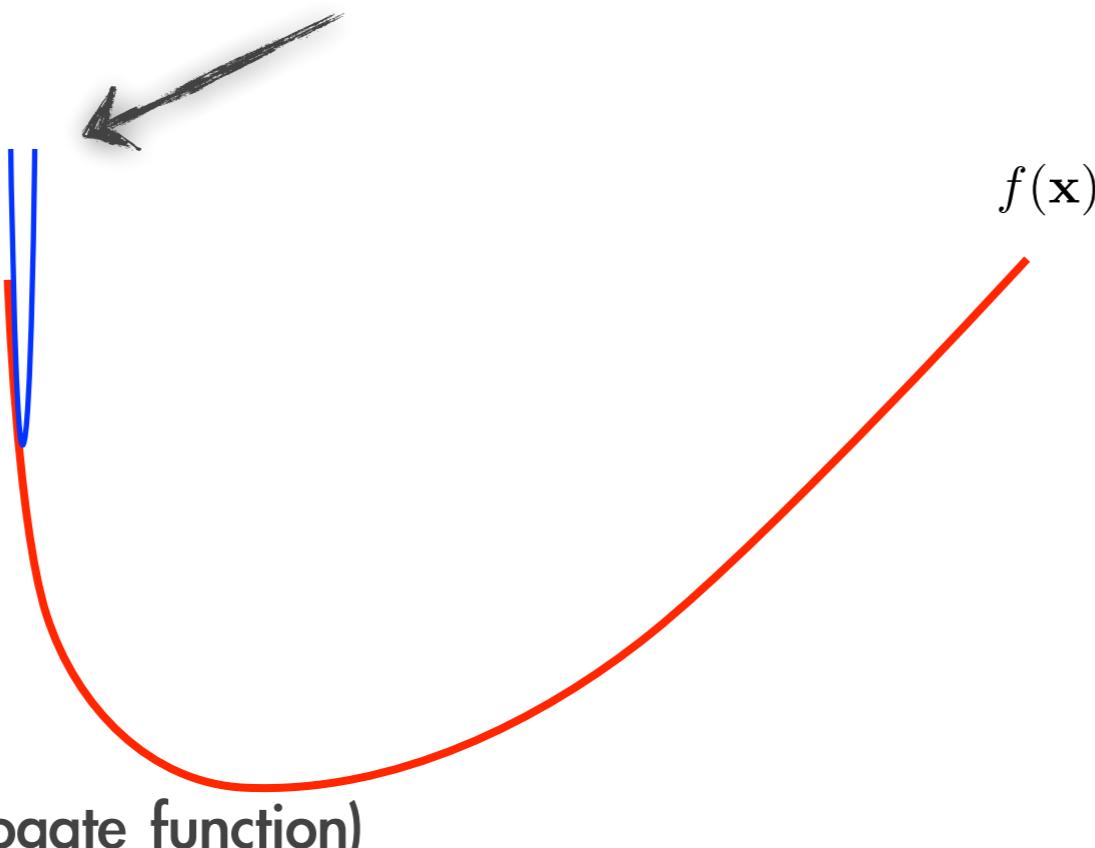


# Self-concordance does not lead to global conditions!

- Main properties of Lipschitz gradient + strongly convex functions

Lower surrogate	$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{\mu}{2}\ \mathbf{y} - \mathbf{x}\ _2^2$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Upper surrogate	$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{L}{2}\ \mathbf{y} - \mathbf{x}\ _2^2$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Hessian surrogates	$\mu\mathbb{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbb{I}$	$\mathbf{x} \in \text{dom}(f)$

$$\mathbf{x}^{k+1} := \arg \min_{\mathbf{x}} \left\{ f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^T(\mathbf{x} - \mathbf{x}^k) + \frac{L}{2}\|\mathbf{x} - \mathbf{x}^k\|_2^2 + g(\mathbf{x}) \right\}$$

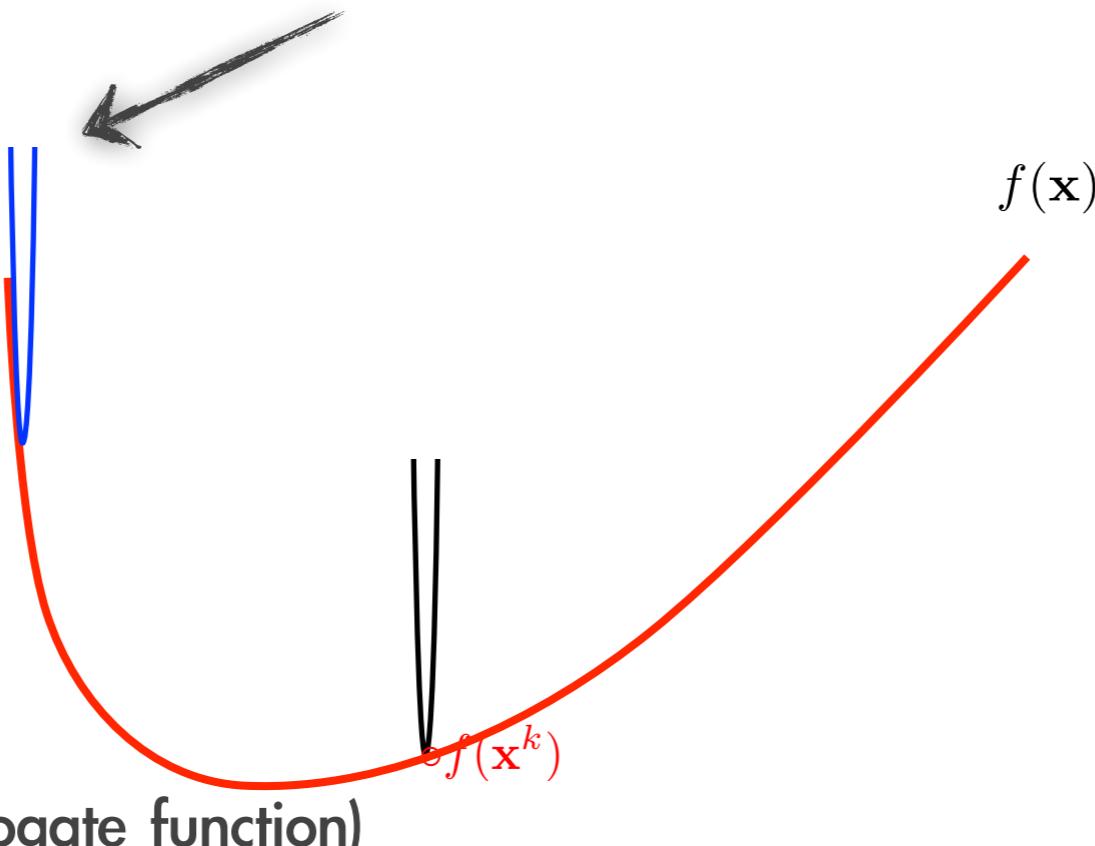


# Self-concordance does not lead to global conditions!

- Main properties of Lipschitz gradient + strongly convex functions

Lower surrogate	$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{\mu}{2}\ \mathbf{y} - \mathbf{x}\ _2^2$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Upper surrogate	$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{L}{2}\ \mathbf{y} - \mathbf{x}\ _2^2$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Hessian surrogates	$\mu\mathbb{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbb{I}$	$\mathbf{x} \in \text{dom}(f)$

$$\mathbf{x}^{k+1} := \arg \min_{\mathbf{x}} \left\{ f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^T(\mathbf{x} - \mathbf{x}^k) + \frac{L}{2}\|\mathbf{x} - \mathbf{x}^k\|_2^2 + g(\mathbf{x}) \right\}$$



# Self-concordance does not lead to global conditions!

- Main properties of Lipschitz gradient + strongly convex functions

Lower surrogate	$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{\mu}{2}\ \mathbf{y} - \mathbf{x}\ _2^2$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Upper surrogate	$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{L}{2}\ \mathbf{y} - \mathbf{x}\ _2^2$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Hessian surrogates	$\mu\mathbb{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbb{I}$	$\mathbf{x} \in \text{dom}(f)$

- Main properties of self-concordant functions

Lower surrogate	$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \omega(\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Upper surrogate	$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \omega_*(\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})$	$\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}} < 1$
Hessian surrogates	$(1 - \ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})^2 \nabla^2 f(\mathbf{x}) \preceq \nabla^2 f(\mathbf{y}) \preceq (1 - \ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})^{-2} \nabla^2 f(\mathbf{x})$	$\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}} < 1$

Local norm:  $\|\mathbf{u}\|_{\mathbf{x}} := [\mathbf{u}^T \nabla^2 f(\mathbf{x}) \mathbf{u}]^{1/2}$

# Self-concordance does not lead to global conditions!

- Main properties of Lipschitz gradient + strongly convex functions

Lower surrogate	$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{\mu}{2}\ \mathbf{y} - \mathbf{x}\ _2^2$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Upper surrogate	$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{L}{2}\ \mathbf{y} - \mathbf{x}\ _2^2$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Hessian surrogates	$\mu\mathbb{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbb{I}$	$\mathbf{x} \in \text{dom}(f)$

- Main properties of self-concordant functions

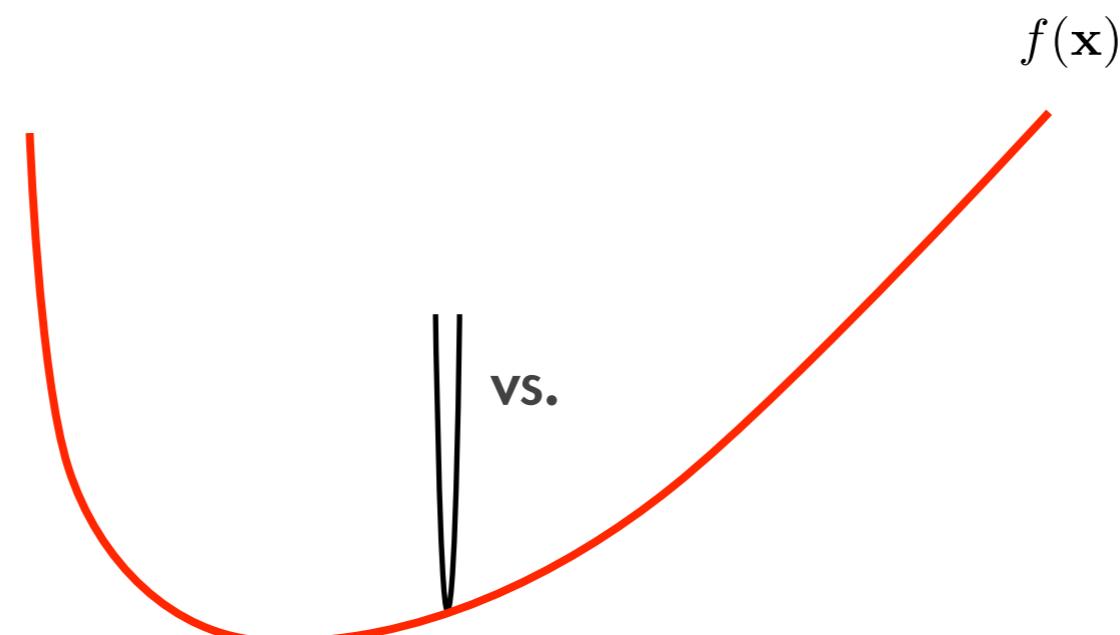
Lower surrogate	$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \omega(\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Upper surrogate	$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \omega_*(\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})$	$\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}} < 1$
Hessian surrogates	$(1 - \ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})^2 \nabla^2 f(\mathbf{x}) \preceq \nabla^2 f(\mathbf{y}) \preceq (1 - \ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})^{-2} \nabla^2 f(\mathbf{x})$	$\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}} < 1$

Local norm:  $\|\mathbf{u}\|_{\mathbf{x}} := [\mathbf{u}^T \nabla^2 f(\mathbf{x}) \mathbf{u}]^{1/2}$

# Self-concordance does not lead to global conditions!

- Main properties of self-concordant functions

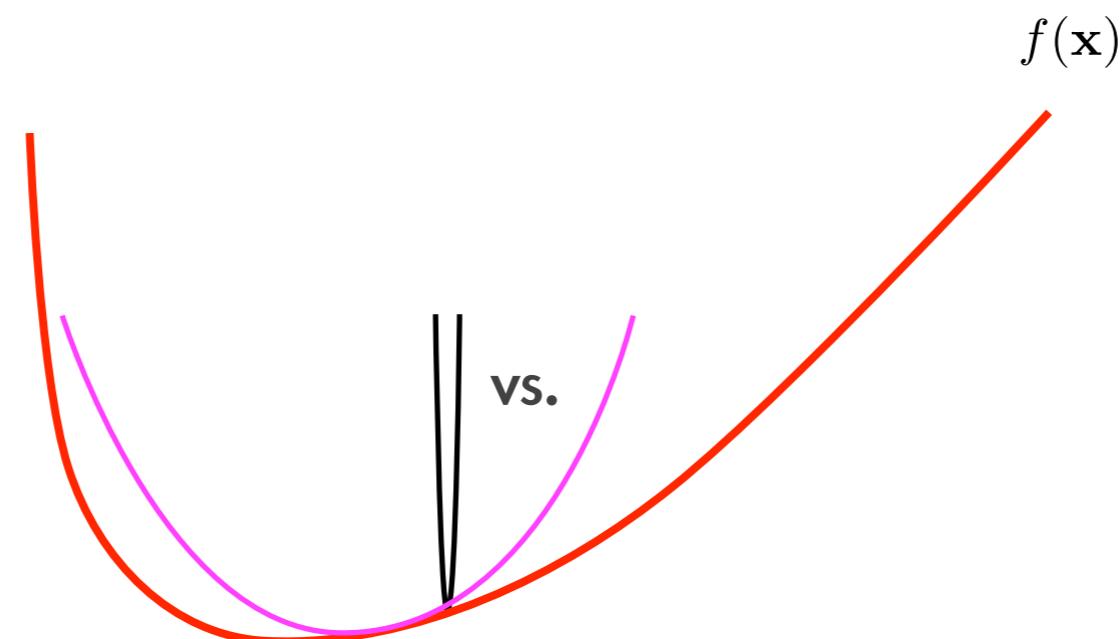
Lower surrogate	$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \omega (\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Upper surrogate	$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \omega_* (\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})$	$\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}} < 1$
Hessian surrogates	$(1 - \ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})^2 \nabla^2 f(\mathbf{x}) \preceq \nabla^2 f(\mathbf{y}) \preceq (1 - \ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})^{-2} \nabla^2 f(\mathbf{x})$	$\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}} < 1$



# Self-concordance does not lead to global conditions!

- Main properties of self-concordant functions

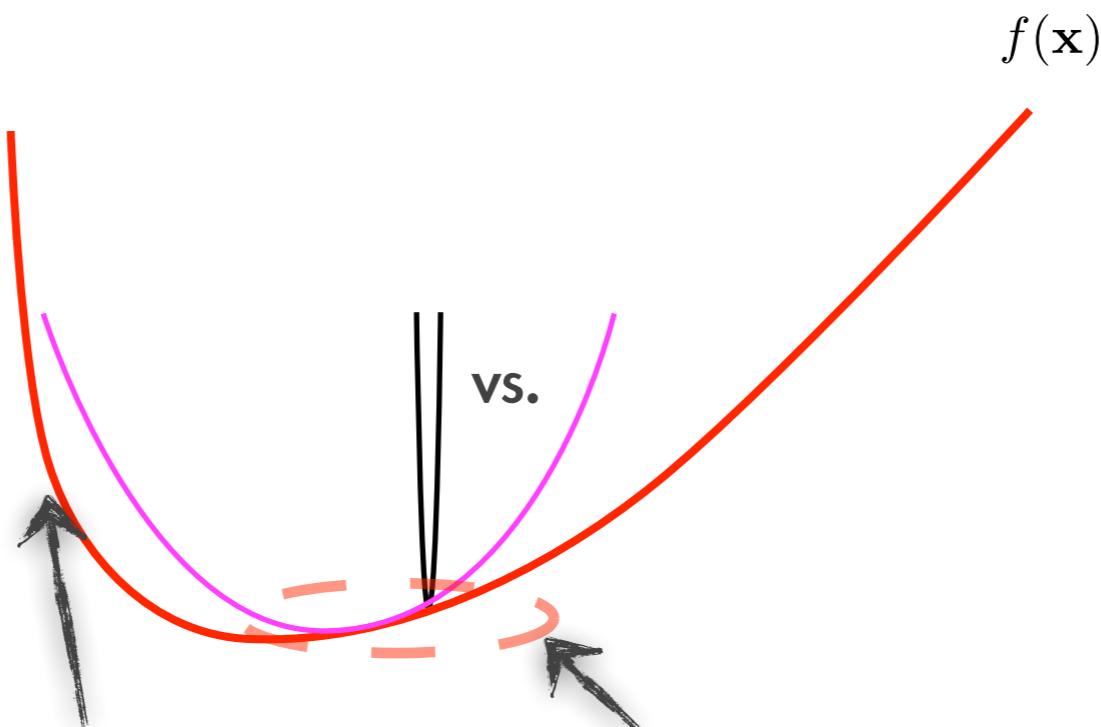
Lower surrogate	$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \omega (\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Upper surrogate	$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \omega_* (\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})$	$\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}} < 1$
Hessian surrogates	$(1 - \ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})^2 \nabla^2 f(\mathbf{x}) \preceq \nabla^2 f(\mathbf{y}) \preceq (1 - \ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})^{-2} \nabla^2 f(\mathbf{x})$	$\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}} < 1$



# Self-concordance does not lead to global conditions!

- Main properties of self-concordant functions

Lower surrogate	$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \omega (\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})$	$\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
Upper surrogate	$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \omega_* (\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})$	$\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}} < 1$
Hessian surrogates	$(1 - \ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})^2 \nabla^2 f(\mathbf{x}) \preceq \nabla^2 f(\mathbf{y}) \preceq (1 - \ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}})^{-2} \nabla^2 f(\mathbf{x})$	$\ \mathbf{y} - \mathbf{x}\ _{\mathbf{x}} < 1$



# Examples

## Examples:

1.  $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$  is L-Lipschitz gradient

## Examples:

1.  $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{Ax}\|_2^2$  is L-Lipschitz gradient  
independent of  $\mathbf{x}_1, \mathbf{x}_2$

$$\|\nabla f(\mathbf{x}_1) - \nabla f(\mathbf{x}_2)\|_2 \leq \|\mathbf{A}^\top \mathbf{A}\|_{\text{sp}} \|\mathbf{x}_1 - \mathbf{x}_2\|_2$$



## Examples:

1.  $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{Ax}\|_2^2$  is L-Lipschitz gradient  
independent of  $\mathbf{x}_1, \mathbf{x}_2$

$$\|\nabla f(\mathbf{x}_1) - \nabla f(\mathbf{x}_2)\|_2 \leq \|\mathbf{A}^\top \mathbf{A}\|_{\text{sp}} \|\mathbf{x}_1 - \mathbf{x}_2\|_2$$

But it is also self-concordant (any linear, quadratic function)

## Examples:

1.  $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{Ax}\|_2^2$  is L-Lipschitz gradient  
independent of  $\mathbf{x}_1, \mathbf{x}_2$

$$\|\nabla f(\mathbf{x}_1) - \nabla f(\mathbf{x}_2)\|_2 \leq \|\mathbf{A}^\top \mathbf{A}\|_{\text{sp}} \|\mathbf{x}_1 - \mathbf{x}_2\|_2$$

But it is also self-concordant (any linear, quadratic function)

2.  $f_{\mathbb{S}_{++}^n}(\mathbf{X}) = -\log \det(\mathbf{X})$  is L-Lipschitz gradient,  
but not globally!

## Examples:

1.  $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{Ax}\|_2^2$  is L-Lipschitz gradient  
independent of  $\mathbf{x}_1, \mathbf{x}_2$

$$\|\nabla f(\mathbf{x}_1) - \nabla f(\mathbf{x}_2)\|_2 \leq \|\mathbf{A}^\top \mathbf{A}\|_{\text{sp}} \|\mathbf{x}_1 - \mathbf{x}_2\|_2$$

But it is also self-concordant (any linear, quadratic function)

2.  $f_{\mathbb{S}_{++}^n}(\mathbf{X}) = -\log \det(\mathbf{X})$  is L-Lipschitz gradient,  
but not globally!

$$\|\nabla f(\mathbf{X}_1) - \nabla f(\mathbf{X}_2)\|_2 \leq \frac{1}{\alpha^2} \|\mathbf{X}_1 - \mathbf{X}_2\|_2$$

where  $\alpha = \min \{\lambda_{\min}(\mathbf{X}_1), \lambda_{\min}(\mathbf{X}_2)\}$  depends on  $\mathbf{X}_1, \mathbf{X}_2$

## Examples:

1.  $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{Ax}\|_2^2$  is L-Lipschitz gradient  
independent of  $\mathbf{x}_1, \mathbf{x}_2$

$$\|\nabla f(\mathbf{x}_1) - \nabla f(\mathbf{x}_2)\|_2 \leq \|\mathbf{A}^\top \mathbf{A}\|_{\text{sp}} \|\mathbf{x}_1 - \mathbf{x}_2\|_2$$

But it is also self-concordant (any linear, quadratic function)

2.  $f_{\mathbb{S}_{++}^n}(\mathbf{X}) = -\log \det(\mathbf{X})$  is L-Lipschitz gradient,  
but not globally!

$$\|\nabla f(\mathbf{X}_1) - \nabla f(\mathbf{X}_2)\|_2 \leq \frac{1}{\alpha^2} \|\mathbf{X}_1 - \mathbf{X}_2\|_2$$

where  $\alpha = \min \{\lambda_{\min}(\mathbf{X}_1), \lambda_{\min}(\mathbf{X}_2)\}$  depends on  $\mathbf{X}_1, \mathbf{X}_2$

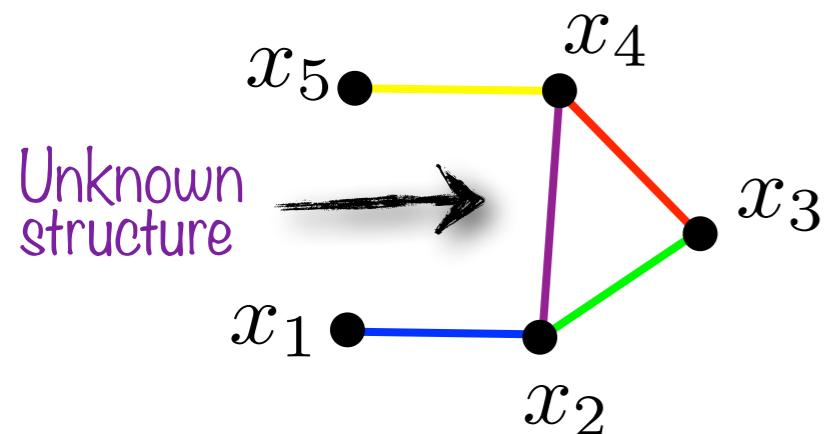
But it is self-concordant!

**...but are all these useful?**

# Application #1: Graphical model selection

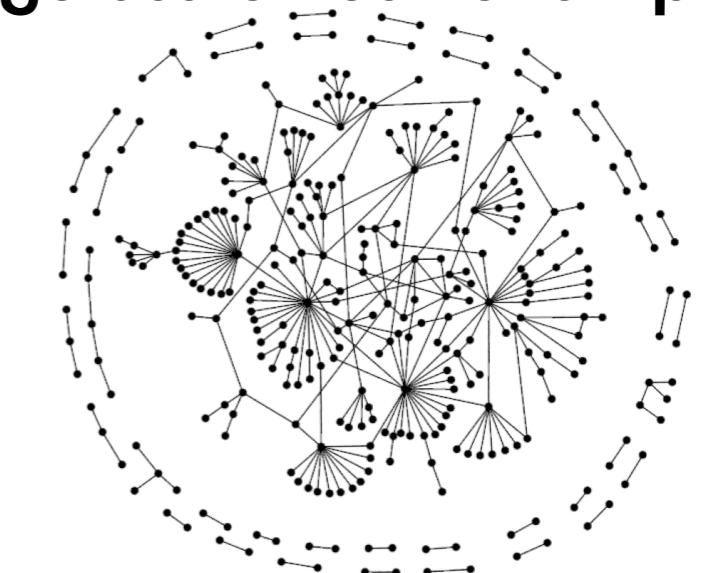
[Dem72, BGA08, HSDR11]

- Toy example



$$\Sigma^{-1} = \begin{matrix} & x_1 & x_2 & x_3 & x_4 & x_5 \\ x_1 & \text{black} & \text{blue} & \text{white} & \text{white} & \text{white} \\ x_2 & \text{blue} & \text{black} & \text{green} & \text{purple} & \text{white} \\ x_3 & \text{white} & \text{green} & \text{black} & \text{red} & \text{white} \\ x_4 & \text{purple} & \text{red} & \text{red} & \text{black} & \text{yellow} \\ x_5 & \text{white} & \text{white} & \text{white} & \text{white} & \text{black} \end{matrix}$$

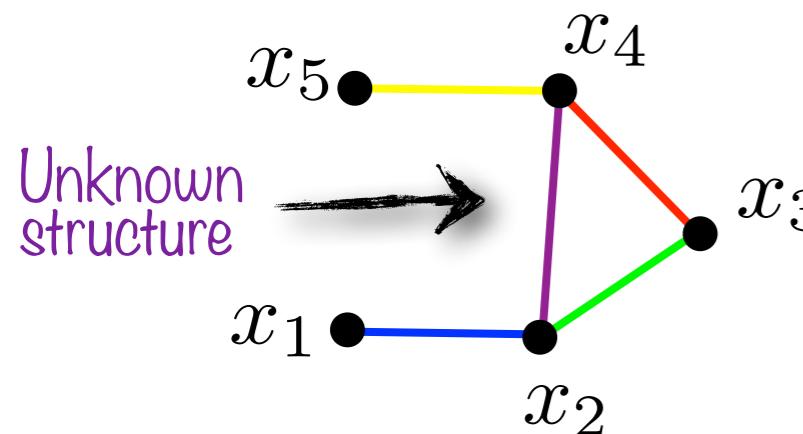
## Large-scale real example



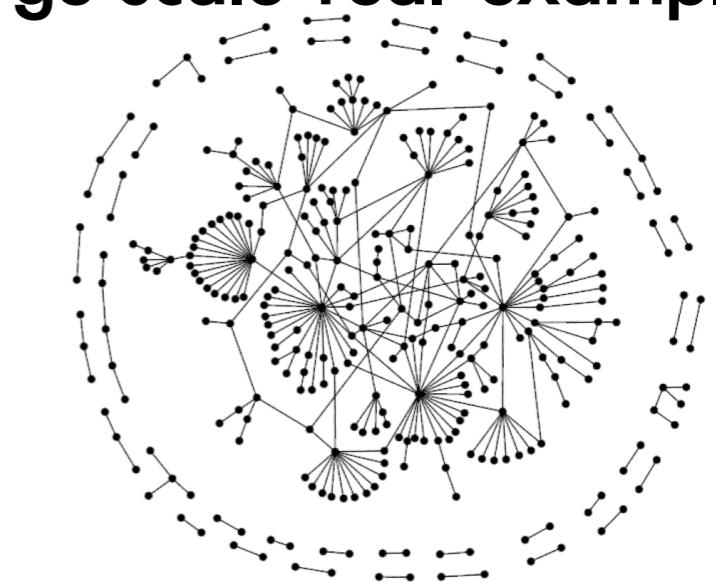
# Application #1: Graphical model selection

[Dem72, BGA08, HSDR11]

## □ Toy example



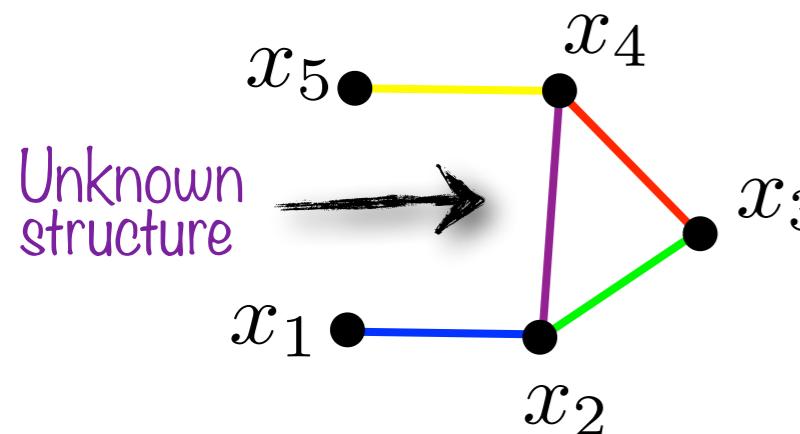
## Large-scale real example



# Application #1: Graphical model selection

[Dem72, BGA08, HSDR11]

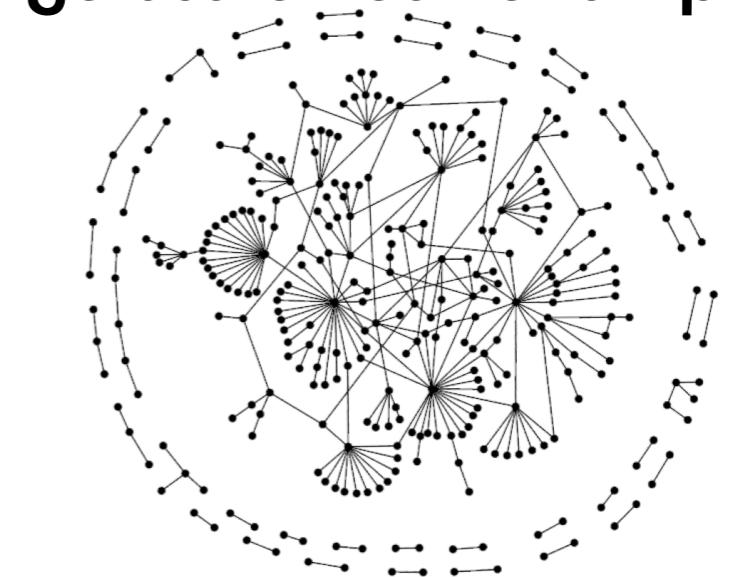
- Toy example



Target

$$\Sigma^{-1} = \begin{matrix} & x_1 & x_2 & x_3 & x_4 & x_5 \\ x_1 & \text{black} & \text{blue} & \text{white} & \text{white} & \text{white} \\ x_2 & \text{blue} & \text{black} & \text{green} & \text{purple} & \text{white} \\ x_3 & \text{white} & \text{green} & \text{black} & \text{red} & \text{white} \\ x_4 & \text{purple} & \text{red} & \text{black} & \text{white} & \text{yellow} \\ x_5 & \text{white} & \text{white} & \text{yellow} & \text{black} & \text{black} \end{matrix}$$

## Large-scale real example

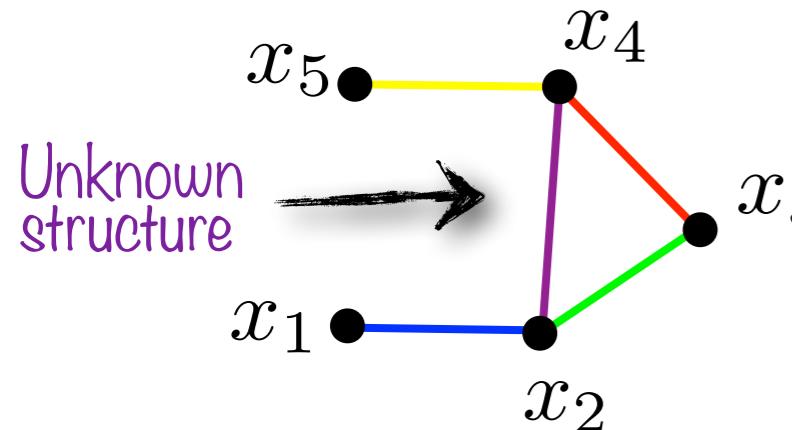


- Given a data set  $\mathcal{D}$ , drawn from a joint pdf with unknown covariance  $\Sigma$ , the aim is to learn a sparse matrix  $\Theta$  that approximates  $\Sigma^{-1}$ .  
Input: sample covariance  $\widehat{\Sigma}$  calculated usually from limited samples

# Application #1: Graphical model selection

[Dem72, BGA08, HSDR11]

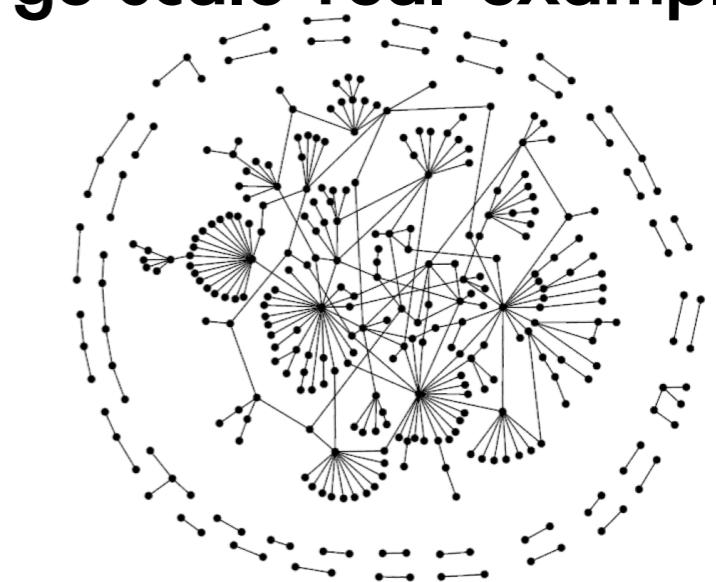
- Toy example



Target

$$\Sigma^{-1} = \begin{matrix} & x_1 & x_2 & x_3 & x_4 & x_5 \\ x_1 & \text{black} & \text{blue} & \text{white} & \text{white} & \text{white} \\ x_2 & \text{blue} & \text{black} & \text{green} & \text{purple} & \text{white} \\ x_3 & \text{white} & \text{green} & \text{black} & \text{red} & \text{white} \\ x_4 & \text{purple} & \text{red} & \text{red} & \text{black} & \text{yellow} \\ x_5 & \text{white} & \text{white} & \text{white} & \text{yellow} & \text{black} \end{matrix}$$

## Large-scale real example



- Given a data set  $\mathcal{D}$ , drawn from a joint pdf with unknown covariance  $\Sigma$ , the aim is to learn a sparse matrix  $\Theta$  that approximates  $\Sigma^{-1}$ .  
Input: sample covariance  $\hat{\Sigma}$  calculated usually from limited samples

Self-concordant function

## Optimization problem

$$\min_{\Theta} \left\{ \underbrace{-\log \det(\Theta) + \text{trace}(\hat{\Sigma}\Theta)}_{f(\cdot)} + \underbrace{\rho \|\Theta\|_1}_{g(\cdot)} \right\}$$

Real-world problems include thousands of variables to optimize

# Theory

*“Composite self-concordant minimization”, Tran-Dinh, Kyrillidis, Cevher, JMLR 2015*

□ Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

□ Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

□ Generic strategy:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i$$

← Direction to move  
↑ Step size

- Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

- Generic strategy:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i$$

← Direction to move  
                  ↑ Step size

- For L-Lipschitz gradient (and strongly convex)  $f(\mathbf{x})$  :

$$\alpha_i \propto (L, \mu)$$

- Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

- Generic strategy:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i$$

← Direction to move  
                  ↑ Step size

- For L-Lipschitz gradient (and strongly convex)  $f(\mathbf{x})$  :

$$\alpha_i \propto (L, \mu)$$

- For direction, proximal first-order schemes:

- $\mathbf{d}_i = \arg \min_{\mathbf{d}} \left\{ f(\mathbf{x}_i) + \langle \nabla f(\mathbf{x}_i), \mathbf{d} \rangle + \frac{L}{2} \|\mathbf{d}\|_2^2 + g(\mathbf{x}_i + \mathbf{d}) \right\}$

- Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

- Generic strategy:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i$$


 Direction to move  
 Step size

- For L-Lipschitz gradient (and strongly convex)  $f(\mathbf{x})$  :

$$\alpha_i \propto (L, \mu)$$

- For direction, proximal first-order schemes:

- $\mathbf{d}_i = \arg \min_{\mathbf{d}} \left\{ f(\mathbf{x}_i) + \langle \nabla f(\mathbf{x}_i), \mathbf{d} \rangle + \frac{L}{2} \|\mathbf{d}\|_2^2 + g(\mathbf{x}_i + \mathbf{d}) \right\}$

- ...and proximal Newton schemes:

- $\mathbf{d}_i = \arg \min_{\mathbf{d}} \left\{ f(\mathbf{x}_i) + \langle \nabla f(\mathbf{x}_i), \mathbf{d} \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x}_i) \mathbf{d}, \mathbf{d} \rangle + g(\mathbf{x}_i + \mathbf{d}) \right\}$

- Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

*In this talk*

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

- Generic strategy:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i \leftarrow \begin{array}{l} \text{Direction to move} \\ \text{Step size} \end{array}$$

- For  $L$ -Lipschitz gradient (and strongly convex)  $f(\mathbf{x})$ :

$$\alpha_i \propto (L, \mu)$$

- For direction, proximal first-order schemes:

- $\mathbf{d}_i = \arg \min_{\mathbf{d}} \left\{ f(\mathbf{x}_i) + \langle \nabla f(\mathbf{x}_i), \mathbf{d} \rangle + \frac{L}{2} \|\mathbf{d}\|_2^2 + g(\mathbf{x}_i + \mathbf{d}) \right\}$

- ...and proximal Newton schemes:

- $\mathbf{d}_i = \arg \min_{\mathbf{d}} \left\{ f(\mathbf{x}_i) + \langle \nabla f(\mathbf{x}_i), \mathbf{d} \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x}_i) \mathbf{d}, \mathbf{d} \rangle + g(\mathbf{x}_i + \mathbf{d}) \right\}$

- Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

*In this talk*

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

- Generic strategy:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i \leftarrow \begin{array}{l} \text{Direction to move} \\ \text{Step size} \end{array}$$

- For  $L$ -Lipschitz gradient (and strongly convex)  $f(\mathbf{x})$ :

$$\alpha_i \propto (L, \mu) \qquad \alpha_i ?$$

- For direction, proximal first-order schemes:

- $\mathbf{d}_i = \arg \min_{\mathbf{d}} \left\{ f(\mathbf{x}_i) + \langle \nabla f(\mathbf{x}_i), \mathbf{d} \rangle + \frac{L}{2} \|\mathbf{d}\|_2^2 + g(\mathbf{x}_i + \mathbf{d}) \right\}$

- ...and proximal Newton schemes:

- $\mathbf{d}_i = \arg \min_{\mathbf{d}} \left\{ f(\mathbf{x}_i) + \langle \nabla f(\mathbf{x}_i), \mathbf{d} \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x}_i) \mathbf{d}, \mathbf{d} \rangle + g(\mathbf{x}_i + \mathbf{d}) \right\}$

- Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

*In this talk*

- Generic strategy:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i \leftarrow \begin{array}{l} \text{Direction to move} \\ \text{Step size} \end{array}$$

**Step size**

- Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

*In this talk*

- Generic strategy:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i \leftarrow \begin{array}{l} \text{Direction to move} \\ \text{Step size} \end{array}$$

**Step size**

- Define Newton decrement:  $\lambda_i = \sqrt{\mathbf{d}_i^T \nabla^2 f(\mathbf{x}_i) \mathbf{d}_i}$

- Step size selection:  $\alpha_i = \frac{1}{\lambda_i + 1} \in (0, 1]$

- Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

*In this talk*

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

- Generic strategy:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i \leftarrow \begin{array}{l} \text{Direction to move} \\ \text{Step size} \end{array}$$

**Step size**

- Define Newton decrement:

$$\lambda_i = \sqrt{\mathbf{d}_i^T \nabla^2 f(\mathbf{x}_i) \mathbf{d}_i}$$

- Step size selection:

$$\alpha_i = \frac{1}{\lambda_i + 1} \in (0, 1]$$



Requires second  
order information

- Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

*In this talk*

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

- Generic strategy:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i$$

Direction to move

Step size

- Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

*In this talk*

- Generic strategy:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i \leftarrow \begin{array}{l} \text{Direction to move} \\ \text{Step size} \end{array}$$

## Algorithm (Newton case)

Damped Newton phase ( $\alpha_i < 1$ )

Newton phase ( $\alpha_i = 1$ )

- Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

*In this talk*

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

- Generic strategy:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i \leftarrow \begin{array}{l} \text{Direction to move} \\ \text{Step size} \end{array}$$

## Algorithm (Newton case)

Damped Newton phase ( $\alpha_i < 1$ )

- Compute  $\mathbf{d}_i$

Newton phase ( $\alpha_i = 1$ )

- Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

*In this talk*

- Generic strategy:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i \leftarrow \begin{array}{l} \text{Direction to move} \\ \text{Step size} \end{array}$$

## Algorithm (Newton case)

### Damped Newton phase ( $\alpha_i < 1$ )

- Compute  $\mathbf{d}_i$
- Compute  $\lambda_i$  and  $\alpha_i$

### Newton phase ( $\alpha_i = 1$ )

- Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

*In this talk*

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

- Generic strategy:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i \leftarrow \begin{array}{l} \text{Direction to move} \\ \text{Step size} \end{array}$$

## Algorithm (Newton case)

### Damped Newton phase ( $\alpha_i < 1$ )

- Compute  $\mathbf{d}_i$
- Compute  $\lambda_i$  and  $\alpha_i$
- Do  $\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i$

### Newton phase ( $\alpha_i = 1$ )

- Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

*In this talk*

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

- Generic strategy:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i$$

*Direction to move*

*Step size*

## Algorithm (Newton case)

### Damped Newton phase ( $\alpha_i < 1$ )

- Compute  $\mathbf{d}_i$
- Compute  $\lambda_i$  and  $\alpha_i$
- Do  $\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i$

### Newton phase ( $\alpha_i = 1$ )

- If  $\lambda_i \leq 0.219$  then,  $\mathbf{x}_{i+1} = \mathbf{x}_i + \mathbf{d}_i$

- Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

*In this talk*

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

- Generic strategy:

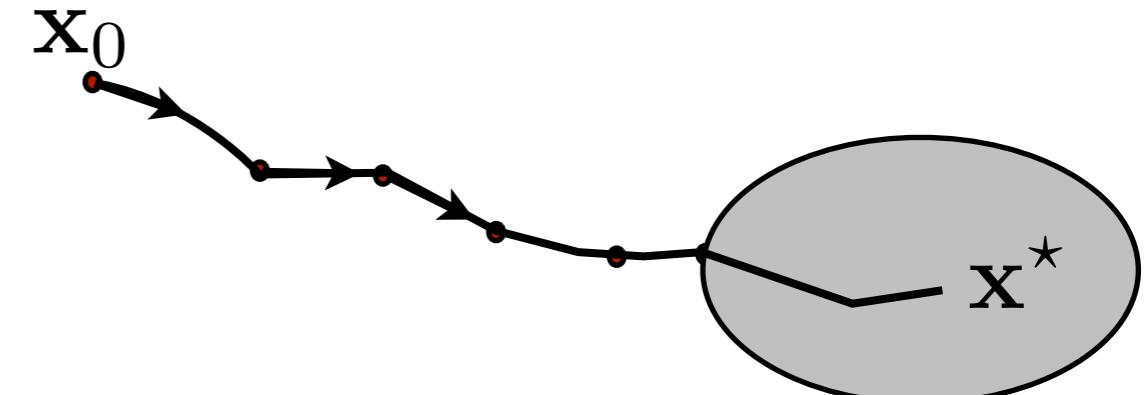
$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i$$

← Direction to move  
↑ Step size

## Algorithm (Newton case)

### Damped Newton phase ( $\alpha_i < 1$ )

- Compute  $\mathbf{d}_i$
- Compute  $\lambda_i$  and  $\alpha_i$
- Do  $\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i$



### Newton phase ( $\alpha_i = 1$ )

- If  $\lambda_i \leq 0.219$  then,  $\mathbf{x}_{i+1} = \mathbf{x}_i + \mathbf{d}_i$

- Recall:

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

*In this talk*

$$\text{prox}_{\lambda g}(\mathbf{w}) = \arg \min_{\mathbf{x}} g(\mathbf{x}) + \frac{1}{2\lambda} \|\mathbf{x} - \mathbf{w}\|_2^2$$

- Generic strategy:

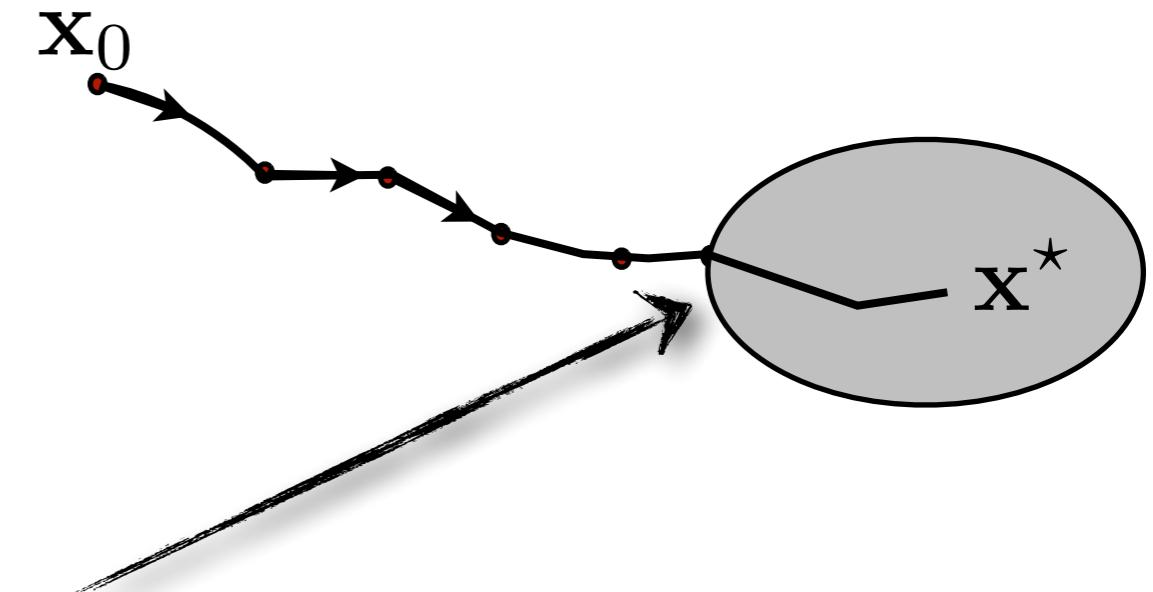
$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i$$

← Direction to move  
↑ Step size

## Algorithm (Newton case)

### Damped Newton phase ( $\alpha_i < 1$ )

- Compute  $\mathbf{d}_i$
- Compute  $\lambda_i$  and  $\alpha_i$
- Do  $\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha_i \mathbf{d}_i$



### Newton phase ( $\alpha_i = 1$ )

- If  $\lambda_i \leq 0.219$  then,  $\mathbf{x}_{i+1} = \mathbf{x}_i + \mathbf{d}_i$

□ Theorem for Newton case:

$$\lambda_i = \sqrt{\mathbf{d}_i^T \nabla^2 f(\mathbf{x}_i) \mathbf{d}_i}$$

If  $\alpha_i = \frac{1}{\lambda_i + 1} \in (0, 1]$ , then  $\mathbf{x}_{i+1}$  satisfies:

$$f(\mathbf{x}_{i+1}) + g(\mathbf{x}_{i+1}) \leq f(\mathbf{x}_i) + g(\mathbf{x}_i) - \omega(\lambda_i), \quad \text{where } \omega(t) = t - \log(1 + t) \geq 0$$

Moreover, step size  $\alpha_i$  is optimal.

□ Theorem for Newton case:

$$\lambda_i = \sqrt{\mathbf{d}_i^T \nabla^2 f(\mathbf{x}_i) \mathbf{d}_i}$$

If  $\alpha_i = \frac{1}{\lambda_i + 1} \in (0, 1]$ , then  $\mathbf{x}_{i+1}$  satisfies:

$$f(\mathbf{x}_{i+1}) + g(\mathbf{x}_{i+1}) \leq f(\mathbf{x}_i) + g(\mathbf{x}_i) - \omega(\lambda_i), \quad \text{where } \omega(t) = t - \log(1 + t) \geq 0$$

Moreover, step size  $\alpha_i$  is optimal.

□ Proof sketch

$$\omega_*(t) = -t - \log(1 - t), \quad t \in [0, 1]$$

□ By convexity of  $f(\cdot)$  and self-concordance:

$$f(\mathbf{x}_{i+1}) \leq f(\mathbf{x}_i) + \alpha_i \nabla f(\mathbf{x}_i)^T \mathbf{d}_i + \omega_* \left( \alpha_i \sqrt{\mathbf{d}_i^T \nabla^2 f(\mathbf{x}_i) \mathbf{d}_i} \right)$$

□ By convexity of  $g(\cdot)$  and optimality conditions:

$$g(\mathbf{x}_{i+1}) \leq g(\mathbf{x}_i) - \alpha_i (\nabla f(\mathbf{x}_i)^T \mathbf{d}_i - \mathbf{d}_i^T \nabla^2 f(\mathbf{x}_i) \mathbf{d}_i)$$

□ Basic algebra:

$$f(\mathbf{x}_{i+1}) + g(\mathbf{x}_{i+1}) \leq f(\mathbf{x}_i) + g(\mathbf{x}_i) - (\alpha_i \lambda_i - \omega_*(\alpha_i \lambda_i))$$

- Theorem for Newton case:

$$\lambda_i = \sqrt{\mathbf{d}_i^T \nabla^2 f(\mathbf{x}_i) \mathbf{d}_i}$$

If  $\alpha_i = \frac{1}{\lambda_i + 1} \in (0, 1]$ , then  $\mathbf{x}_{i+1}$  satisfies:

$$f(\mathbf{x}_{i+1}) + g(\mathbf{x}_{i+1}) \leq f(\mathbf{x}_i) + g(\mathbf{x}_i) - \omega(\lambda_i), \quad \text{where } \omega(t) = t - \log(1 + t) \geq 0$$

Moreover, step size  $\alpha_i$  is optimal.

- Guaranteed objective function decrease per iteration.
- No back-tracking line search = no function evaluations (expensive)
- Guaranteed (local) quadratic convergence rate.
- Exact analytic complexity on the number of iterations.

□ Theorem for Newton case:

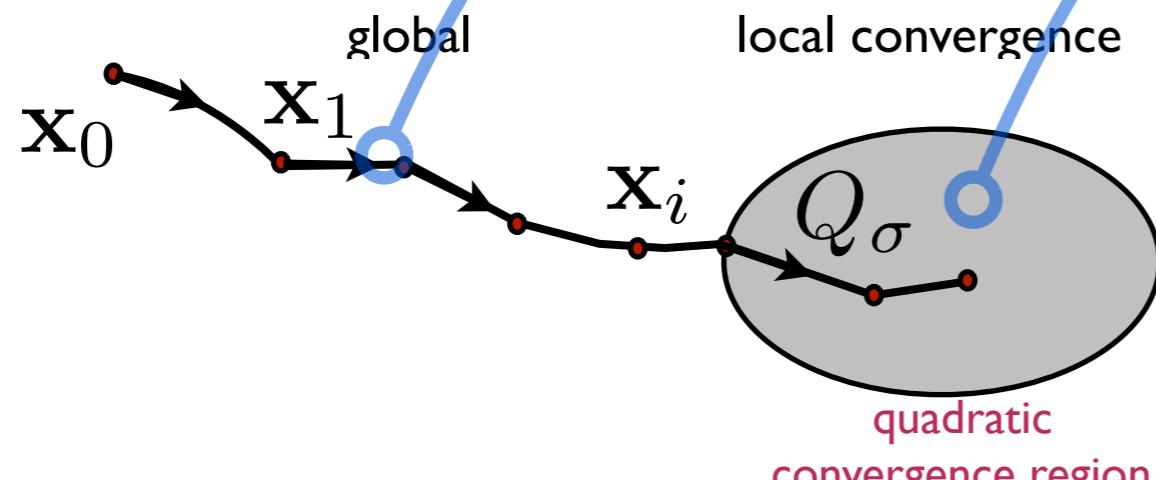
$$\lambda_i = \sqrt{\mathbf{d}_i^T \nabla^2 f(\mathbf{x}_i) \mathbf{d}_i}$$

If  $\alpha_i = \frac{1}{\lambda_i + 1} \in (0, 1]$ , then  $\mathbf{x}_{i+1}$  satisfies:

$$f(\mathbf{x}_{i+1}) + g(\mathbf{x}_{i+1}) \leq f(\mathbf{x}_i) + g(\mathbf{x}_i) - \omega(\lambda_i), \quad \text{where } \omega(t) = t - \log(1+t) \geq 0$$

Moreover, step size  $\alpha_i$  is optimal.

$$\# \text{iterations} = \left\lfloor \frac{\phi(\mathbf{x}_0) - \phi(\mathbf{x}^*)}{0.017} \right\rfloor + O(\log \log (0.3/\epsilon))$$



$$\phi(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x})$$

$$\|\mathbf{x}_i - \mathbf{x}^*\| \leq 2\epsilon$$

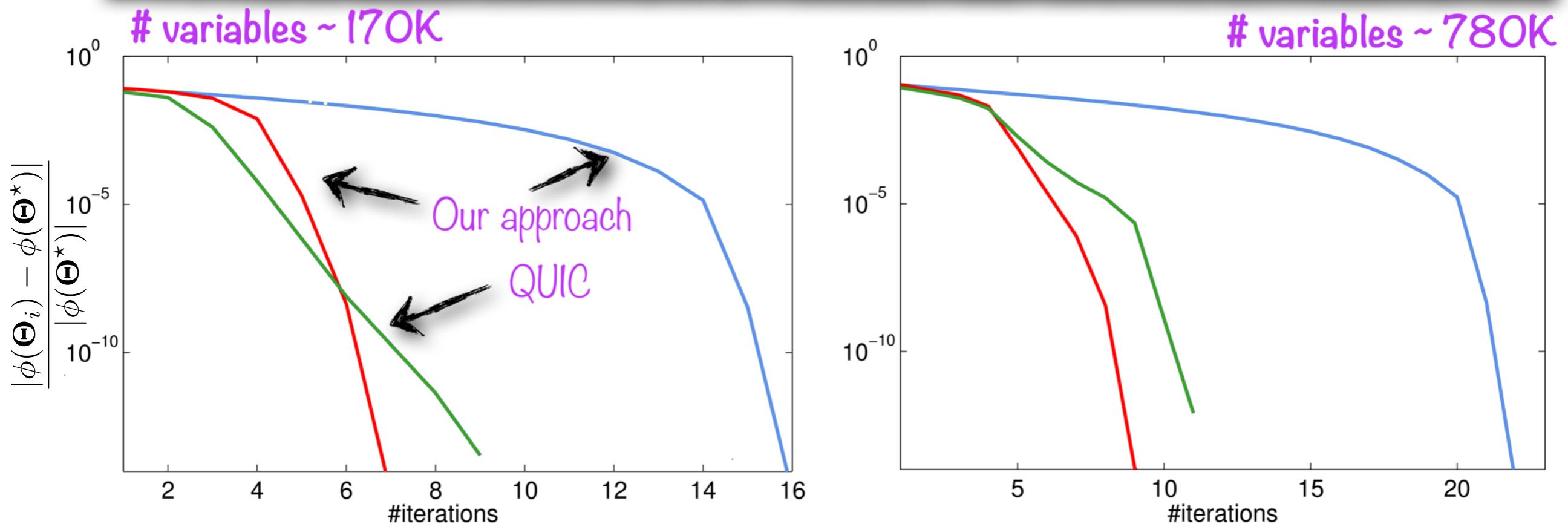
$$Q_\sigma = \{\mathbf{x}_i \mid \lambda_i \leq 0.219\}$$

# Results

$$\phi(\cdot) := f(\cdot) + g(\cdot)$$

$$\min_{\Theta} \left\{ \underbrace{-\log \det(\Theta) + \text{trace}(\widehat{\Sigma}\Theta)}_{f(\cdot)} + \underbrace{\rho \|\Theta\|_1}_{g(\cdot)} \right\}$$

Convergence behavior on gene data [Dem72, BGA08, HSDR11]

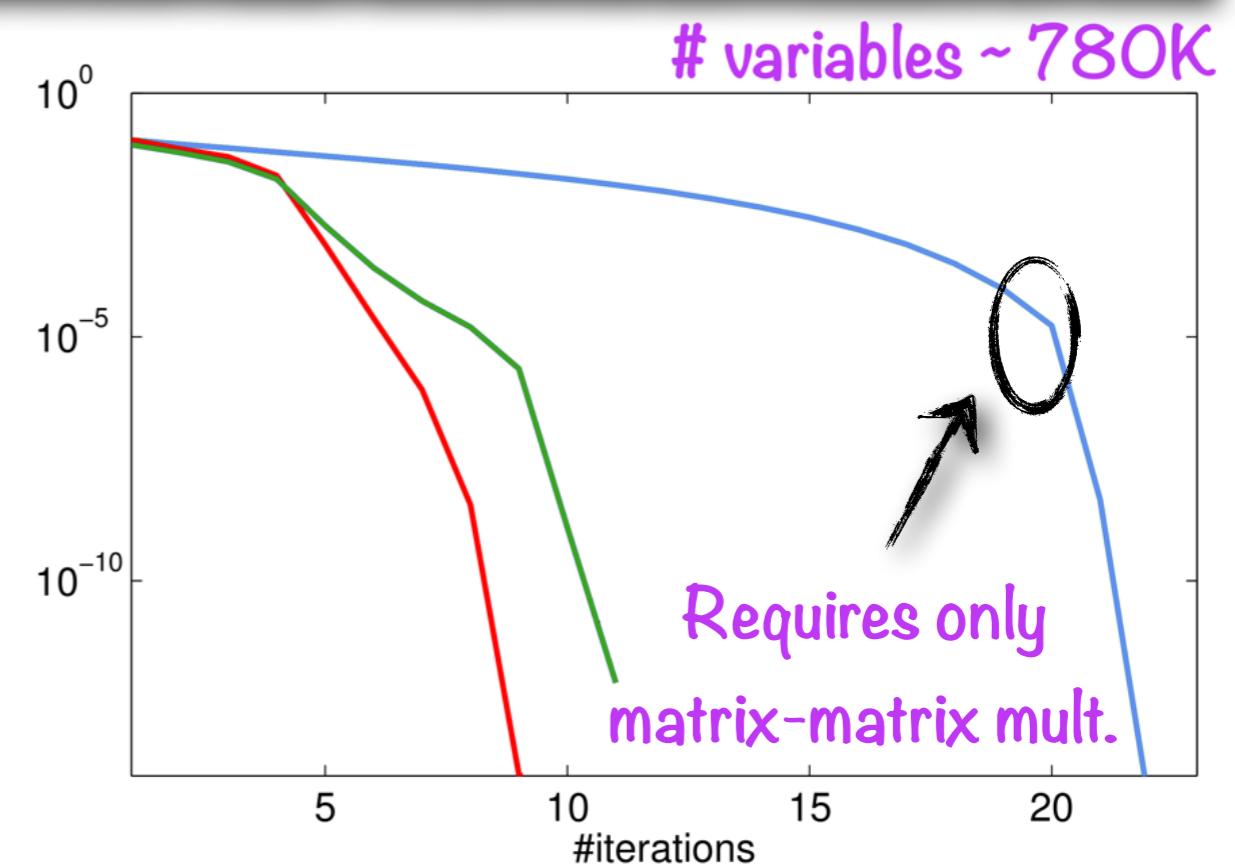
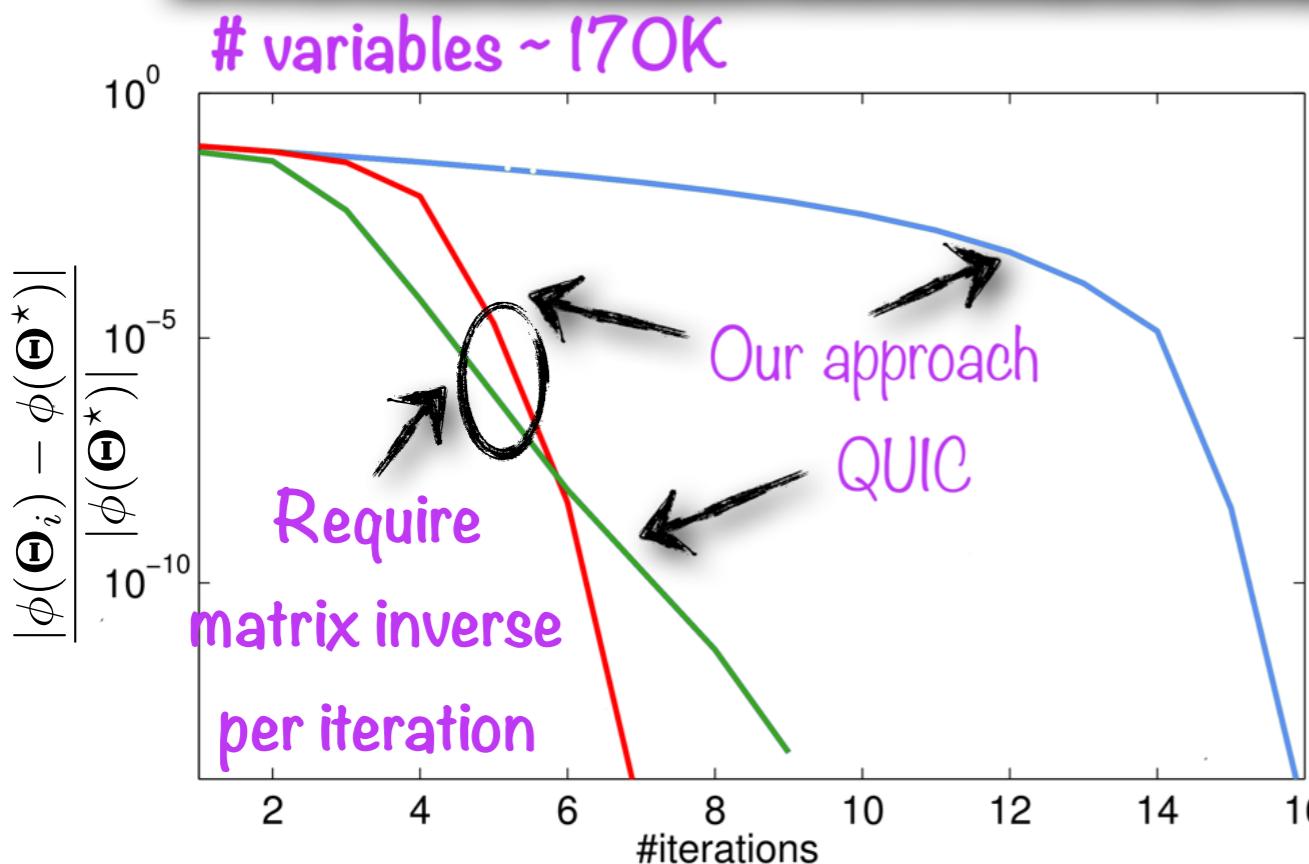


$$\nabla f(\Theta) = \Theta^{-1} + \widehat{\Sigma}$$

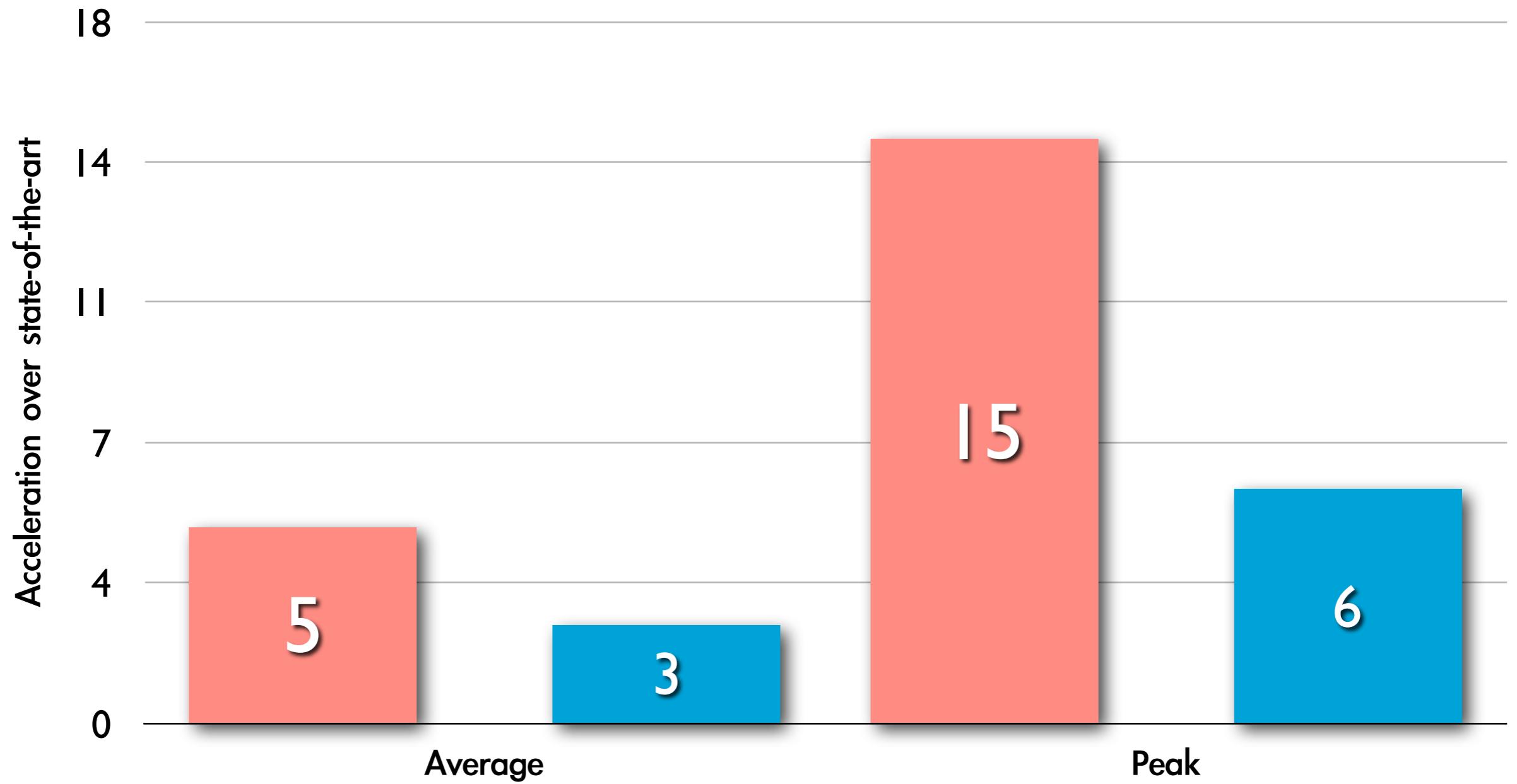
$$\phi(\cdot) := f(\cdot) + g(\cdot)$$

$$\min_{\Theta} \left\{ \underbrace{-\log \det(\Theta) + \text{trace}(\widehat{\Sigma}\Theta)}_{f(\cdot)} + \underbrace{\rho \|\Theta\|_1}_{g(\cdot)} \right\}$$

Convergence behavior on gene data [Dem72, BGA08, HSDR11]

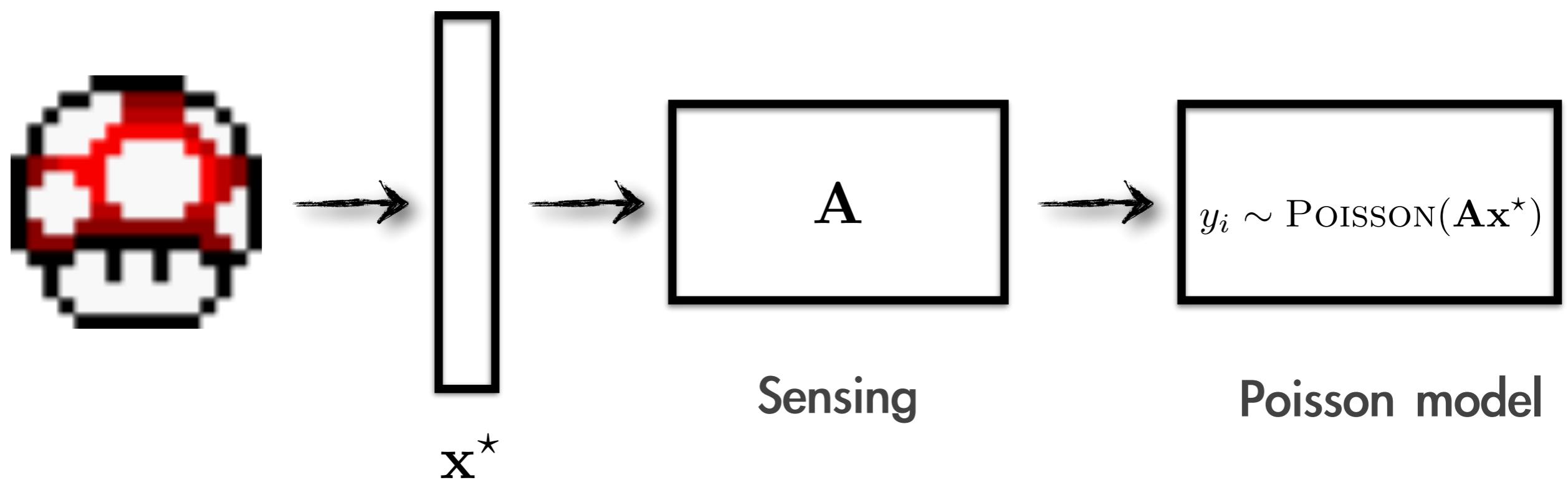


■ Our approach with inversions/Cholesky's dec.      ■ Our approach with matrix mult.



# Application #2: Poisson imaging reconstruction

$$\mathbf{x}^* \in \arg \min_{\mathbf{x}} \left\{ \underbrace{\sum_{i=1}^m \mathbf{a}_i^T \mathbf{x} - \sum_{i=1}^m y_i \log (\mathbf{a}_i^T \mathbf{x} + \beta_i)}_{:=\phi(\mathbf{x})} + g(\mathbf{x}) \right\}$$

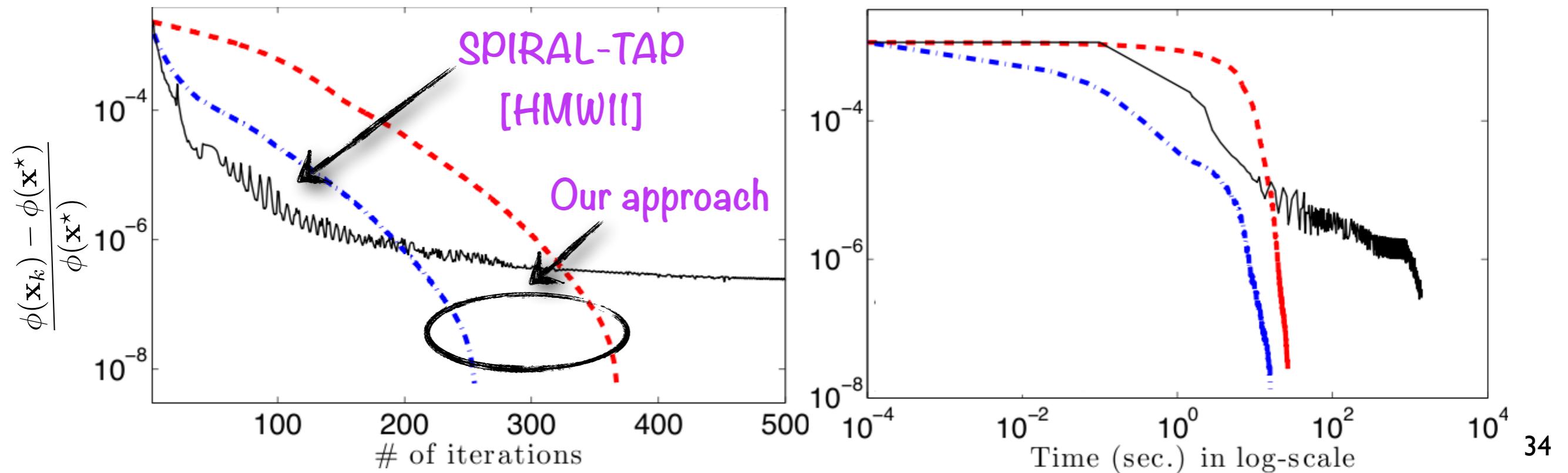


# Application #2: Poisson imaging reconstruction

## Optimization problem

$$\mathbf{x}^* \in \arg \min_{\mathbf{x}} \left\{ \sum_{i=1}^m \mathbf{a}_i^T \mathbf{x} - \sum_{i=1}^m y_i \log (\mathbf{a}_i^T \mathbf{x} + \beta_i) + g(\mathbf{x}) \right\}$$

$\underbrace{\qquad\qquad\qquad}_{:=\phi(\mathbf{x})}$

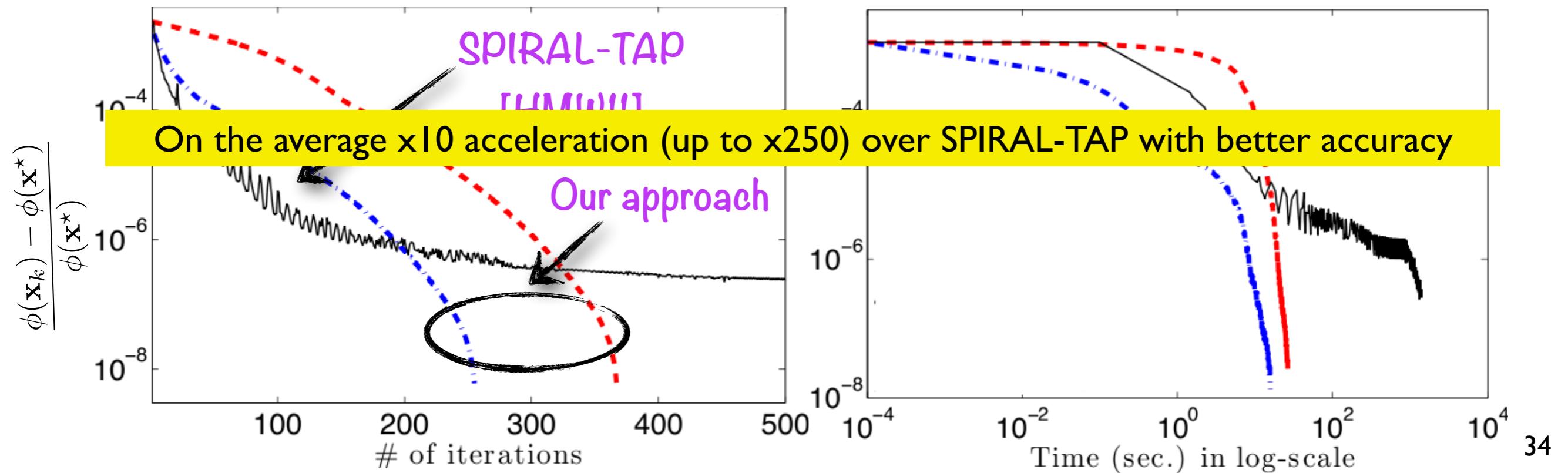


# Application #2: Poisson imaging reconstruction

## Optimization problem

$$\mathbf{x}^* \in \arg \min_{\mathbf{x}} \left\{ \sum_{i=1}^m \mathbf{a}_i^T \mathbf{x} - \sum_{i=1}^m y_i \log (\mathbf{a}_i^T \mathbf{x} + \beta_i) + g(\mathbf{x}) \right\}$$

$\underbrace{\qquad\qquad\qquad}_{:=\phi(\mathbf{x})}$



# Bonus: Path following, composite self-concordant minimization

*“An inexact proximal path-following algorithm for constrained convex minimization”,  
Tran-Dinh, Kyrillidis, Cevher, SIAM Journal on Optimization, 2015*

*Until this point*

$$\underset{\mathbf{X} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad f(\mathbf{X}) + \lambda \cdot g(\mathbf{X})$$

- $f(\mathbf{X})$  appears naturally in the objective
- $f(\mathbf{X})$  is self-concordant

$$\begin{array}{ll}\text{minimize} & g(\mathbf{X}) \\ \mathbf{X} \in \mathbb{R}^{n \times n} & \\ \text{subject to} & \mathbf{X} \in \Omega\end{array}$$

- $\Omega \subseteq \mathbb{R}^{n \times n}$  is endowed with a self-concordant barrier  $f(\mathbf{X})$ .

minimize  $g(\mathbf{X})$   
 $\mathbf{X} \in \mathbb{R}^{n \times n}$

subject to  $\mathbf{X} \in \Omega$

- $\Omega \subseteq \mathbb{R}^{n \times n}$  is endowed with a self-concordant barrier  $f(\mathbf{X})$ .

■ Examples:

$$\begin{aligned}\Omega : \mathbf{X} \succeq 0 &\Rightarrow f_\Omega(\mathbf{X}) = -\log \det(\mathbf{X}) \text{ (solves SDP problem)} \\ \Omega : \mathbf{a}^T \mathbf{x} \geq 0 &\Rightarrow f_\Omega(\mathbf{x}) = -\log(\mathbf{a}^T \mathbf{x}) \\ \Omega : \|\mathbf{A}\mathbf{x} - \mathbf{b}\| \leq \sigma &\Rightarrow f_\Omega(\mathbf{x}) = -\log(\sigma^2 - \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2)\end{aligned}$$

$$\begin{array}{ll}\text{minimize} & g(\mathbf{X}) \\ \mathbf{X} \in \mathbb{R}^{n \times n} &\end{array}$$

subject to  $\mathbf{X} \in \Omega$

Approach

- Solve a sequence of composite self-concordant problems

$$\begin{array}{ll}\text{minimize} & g(\mathbf{X}) + t \cdot f(\mathbf{X}) \\ \mathbf{X} \in \text{int}(\Omega) &\end{array}$$

as  $t$  decreases to zero

## Contributions in a nutshell:

$$\underset{\mathbf{X} \in \text{int}(\Omega)}{\text{minimize}} \quad g(\mathbf{X}) + t \cdot f(\mathbf{X})$$

- First (?) path following scheme in non-smooth, composite form.

## Contributions in a nutshell:

$$\underset{\mathbf{X} \in \text{int}(\Omega)}{\text{minimize}} \quad g(\mathbf{X}) + t \cdot f(\mathbf{X})$$

- First (?) path following scheme in non-smooth, composite form.
- Adaptive selection of  $t_k$  values such that we approximately track the solution trajectory.

## Contributions in a nutshell:

$$\underset{\mathbf{X} \in \text{int}(\Omega)}{\text{minimize}} \quad g(\mathbf{X}) + t \cdot f(\mathbf{X})$$

- First (?) path following scheme in non-smooth, composite form.
- Adaptive selection of  $t_k$  values such that we approximately track the solution trajectory.
- For each  $t_k$ , new estimate is computed from a one-iteration proximal Newton scheme!

## Contributions in a nutshell:

$$\underset{\mathbf{X} \in \text{int}(\Omega)}{\text{minimize}} \quad g(\mathbf{X}) + t \cdot f(\mathbf{X})$$

- First (?) path following scheme in non-smooth, composite form.
- Adaptive selection of  $t_k$  values such that we approximately track the solution trajectory.
- For each  $t_k$ , new estimate is computed from a one-iteration proximal Newton scheme!
- Rigorous convergence guarantees (not presented here).

## Contributions in a nutshell:

$$\underset{\mathbf{X} \in \text{int}(\Omega)}{\text{minimize}} \quad g(\mathbf{X}) + t \cdot f(\mathbf{X})$$

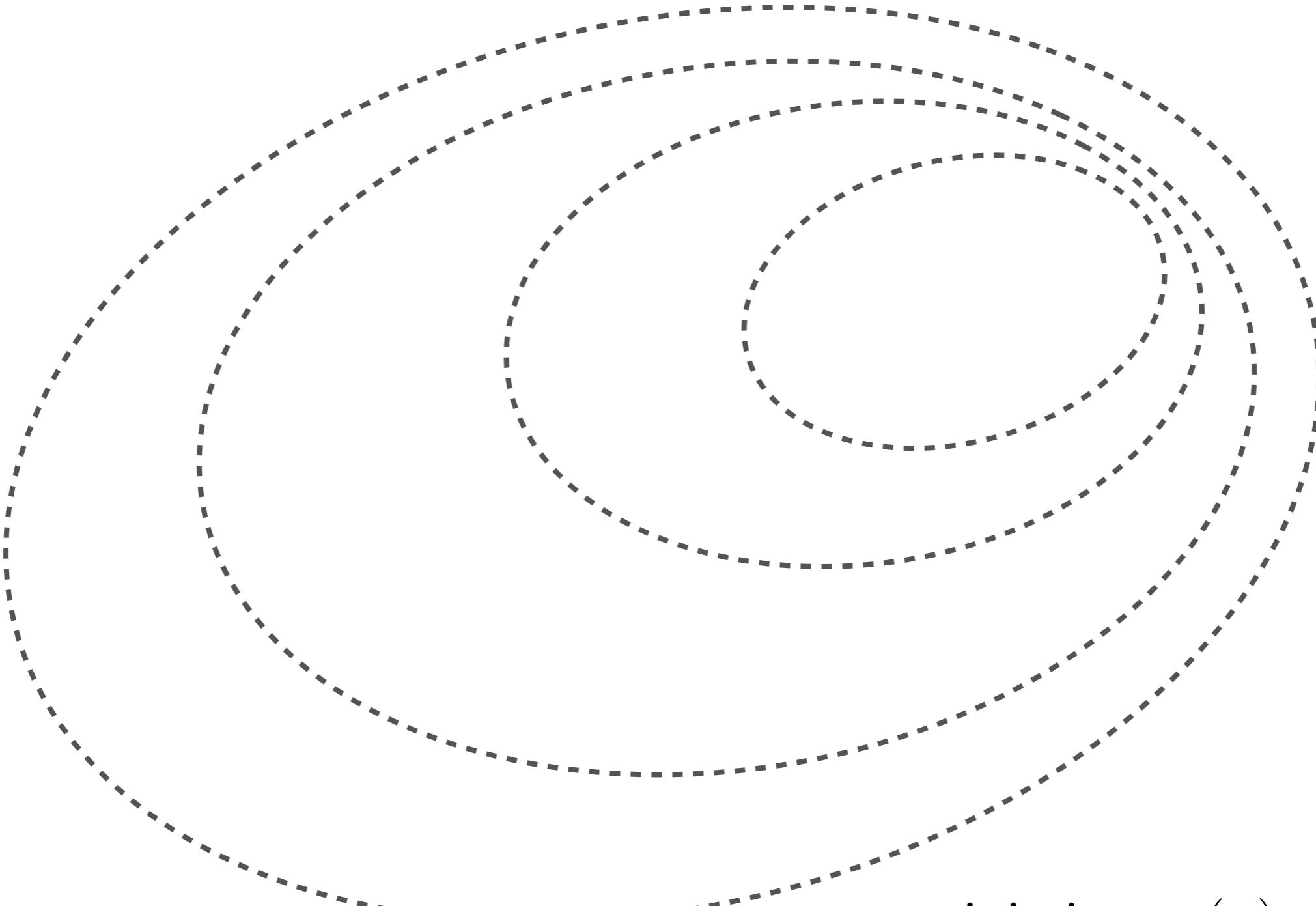
- First (?) path following scheme in non-smooth, composite form.
- Adaptive selection of  $t_k$  values such that we approximately track the solution trajectory.
- For each  $t_k$ , new estimate is computed from a one-iteration proximal Newton scheme!
- Rigorous convergence guarantees (not presented here).

## Algorithm in a nutshell:

- Adaptively set  $t_k$  (closed form solution).
- Solve  $\arg \min \{g(\mathbf{x}) + t_k \cdot f(\mathbf{x})\}$  up to  $\delta$  accuracy with one-iteration proximal Newton scheme.

- Points on central path
  - Central path trajectory
  - ▲ Inexact solutions
- 

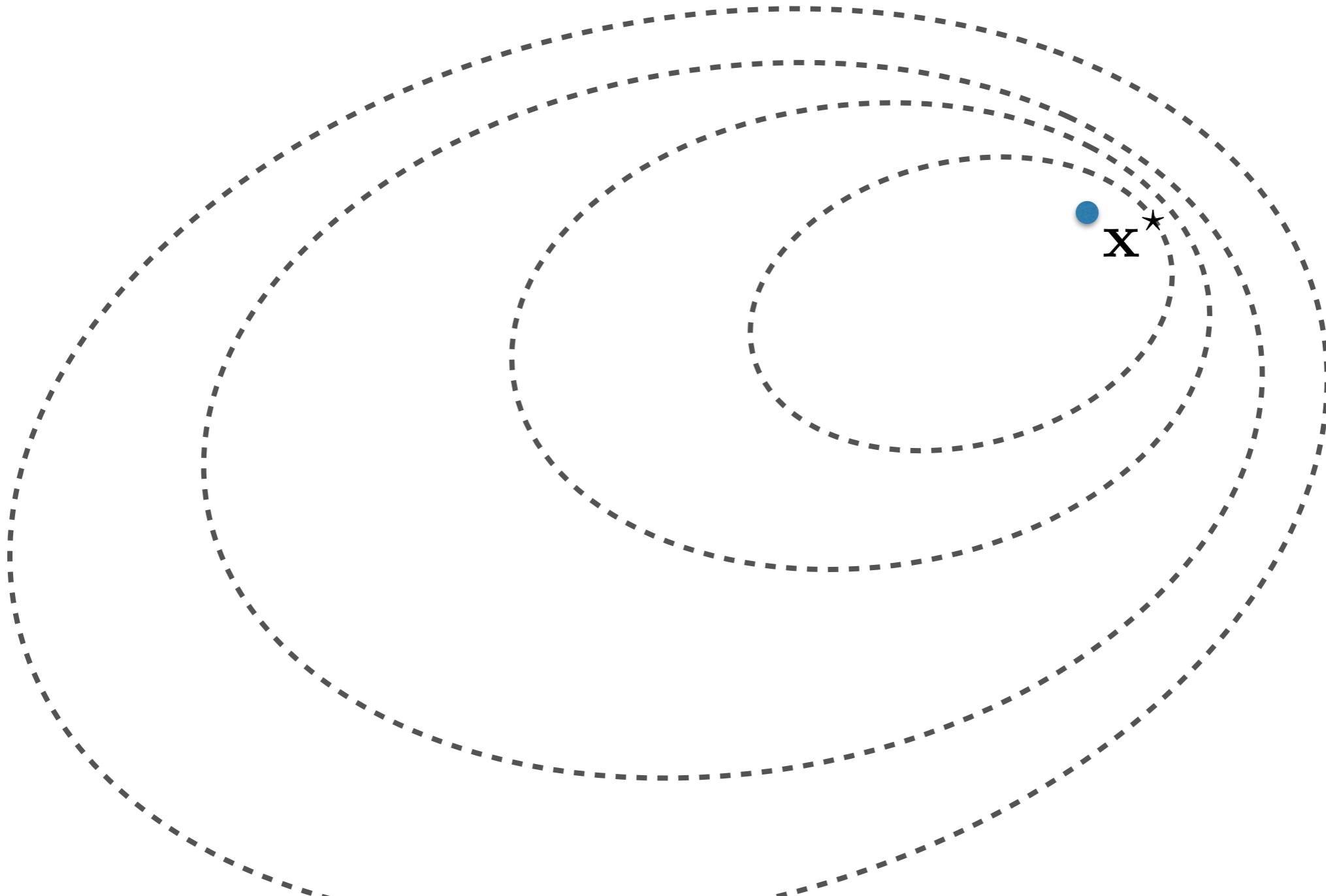
- Points on central path
- Central path trajectory
- ▲ Inexact solutions



minimize  $\mathbf{x} \in \text{int}(\Omega) g(\mathbf{x}) + t \cdot f(\mathbf{x})$

- Points on central path
- Central path trajectory
- ▲ Inexact solutions

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \Omega} g(\mathbf{x})$$

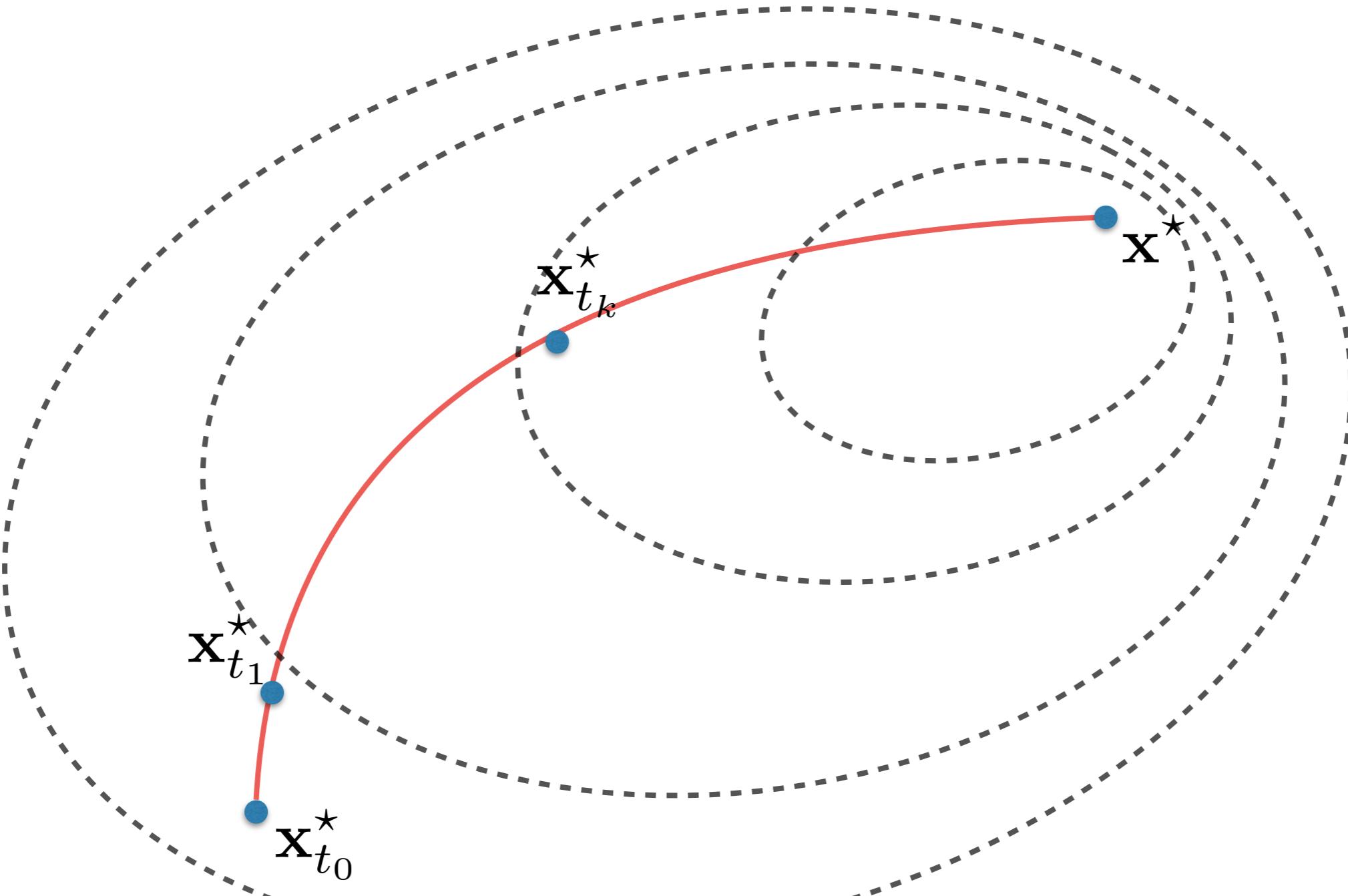


$$\underset{\mathbf{x} \in \text{int}(\Omega)}{\text{minimize}} \quad g(\mathbf{x}) + t \cdot f(\mathbf{x})$$

- Points on central path
- Central path trajectory
- ▲ Inexact solutions

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \Omega} g(\mathbf{x})$$

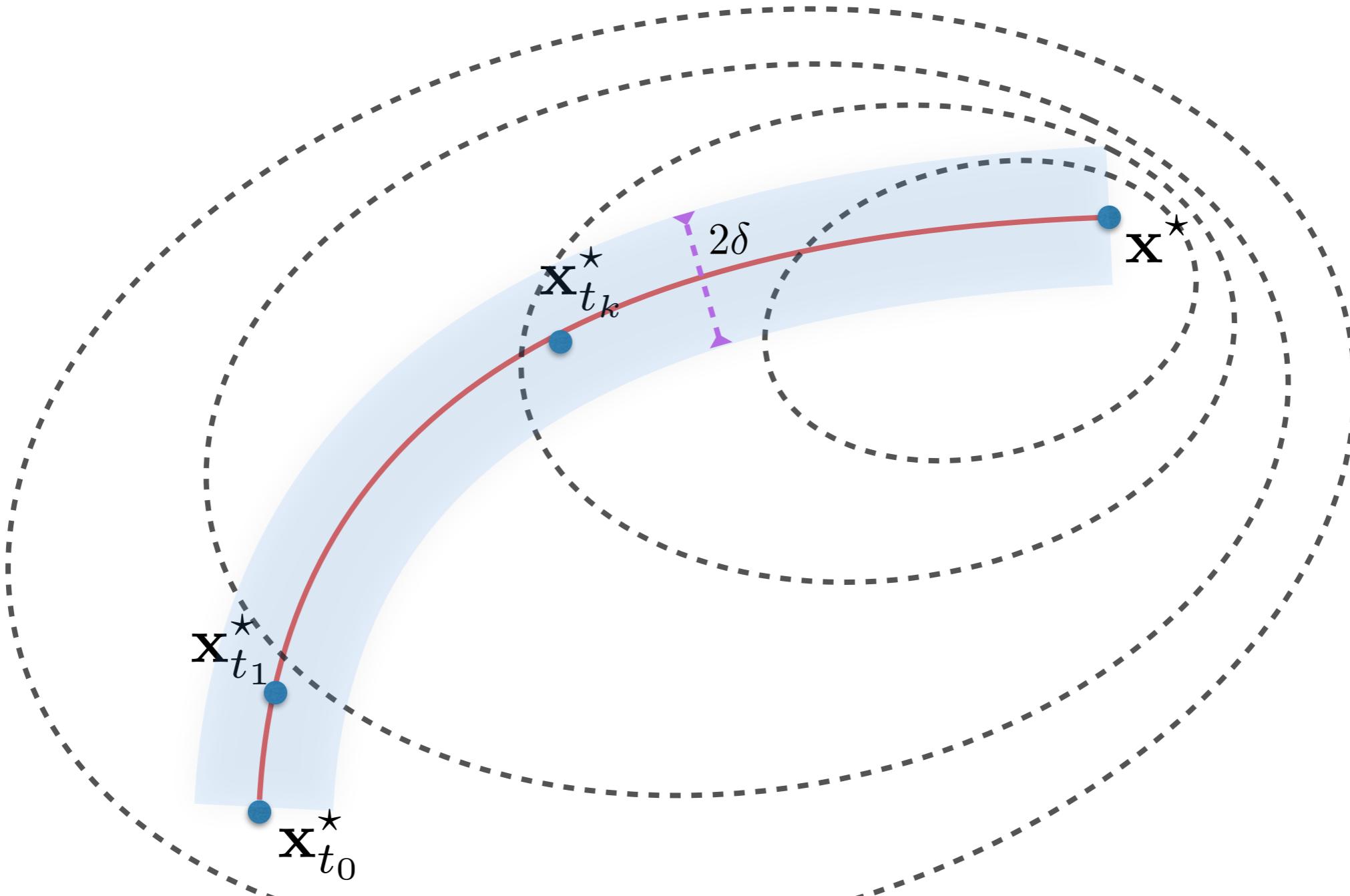
$$\mathbf{x}_{t_k}^* = \arg \min \{g(\mathbf{x}) + t_k \cdot f(\mathbf{x})\}$$



minimize  $\mathbf{x} \in \text{int}(\Omega)$   $g(\mathbf{x}) + t \cdot f(\mathbf{x})$

- Points on central path
- Central path trajectory
- ▲ Inexact solutions

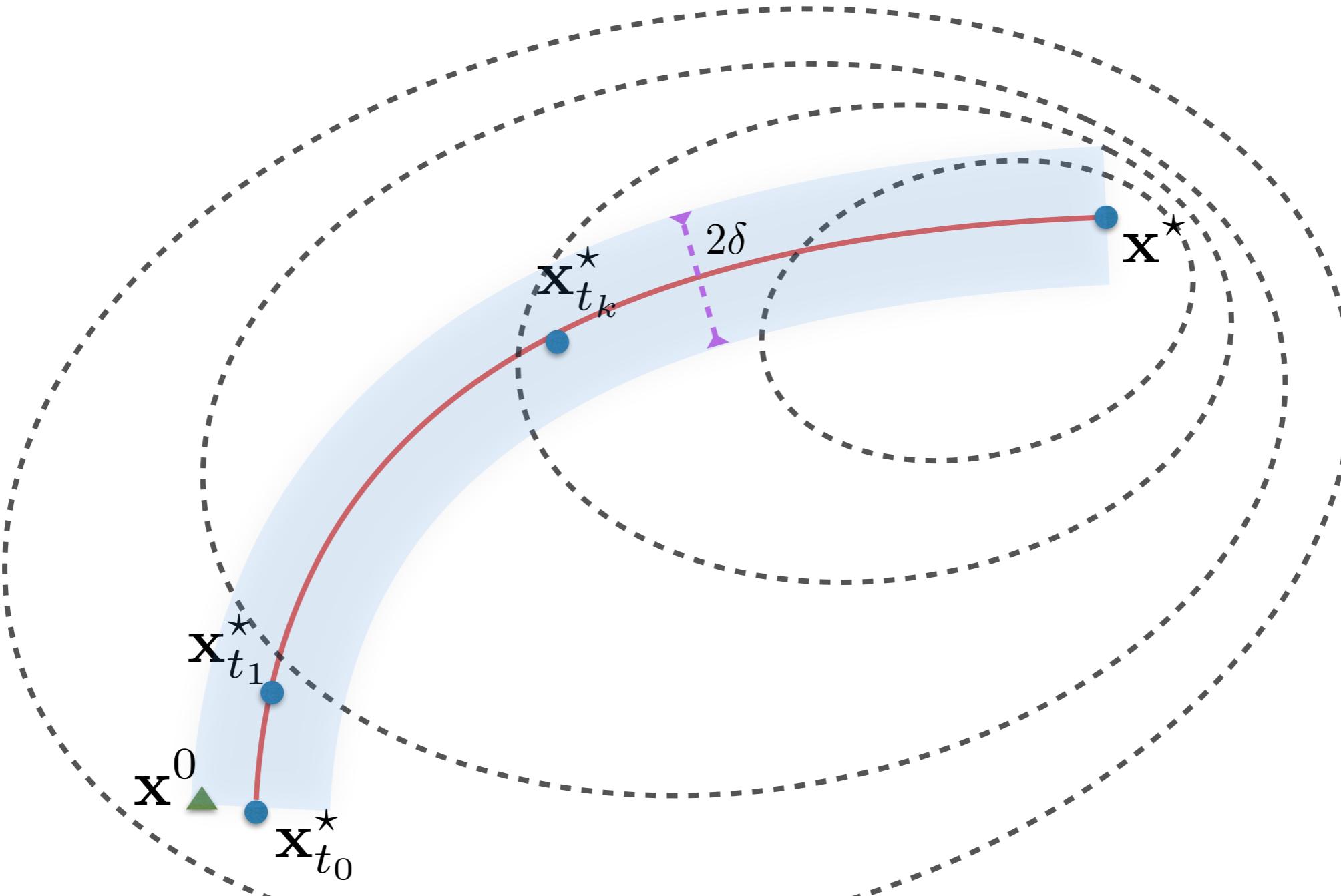
$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \Omega} g(\mathbf{x})$   
 $\mathbf{x}_{t_k}^* = \arg \min \{g(\mathbf{x}) + t_k \cdot f(\mathbf{x})\}$   
 $\mathbf{x}^k$ :  $\delta$ -approximate solution with  
 one proximal Newton step!



$$\underset{\mathbf{x} \in \text{int}(\Omega)}{\text{minimize}} \quad g(\mathbf{x}) + t \cdot f(\mathbf{x})$$

- Points on central path
- Central path trajectory
- ▲ Inexact solutions

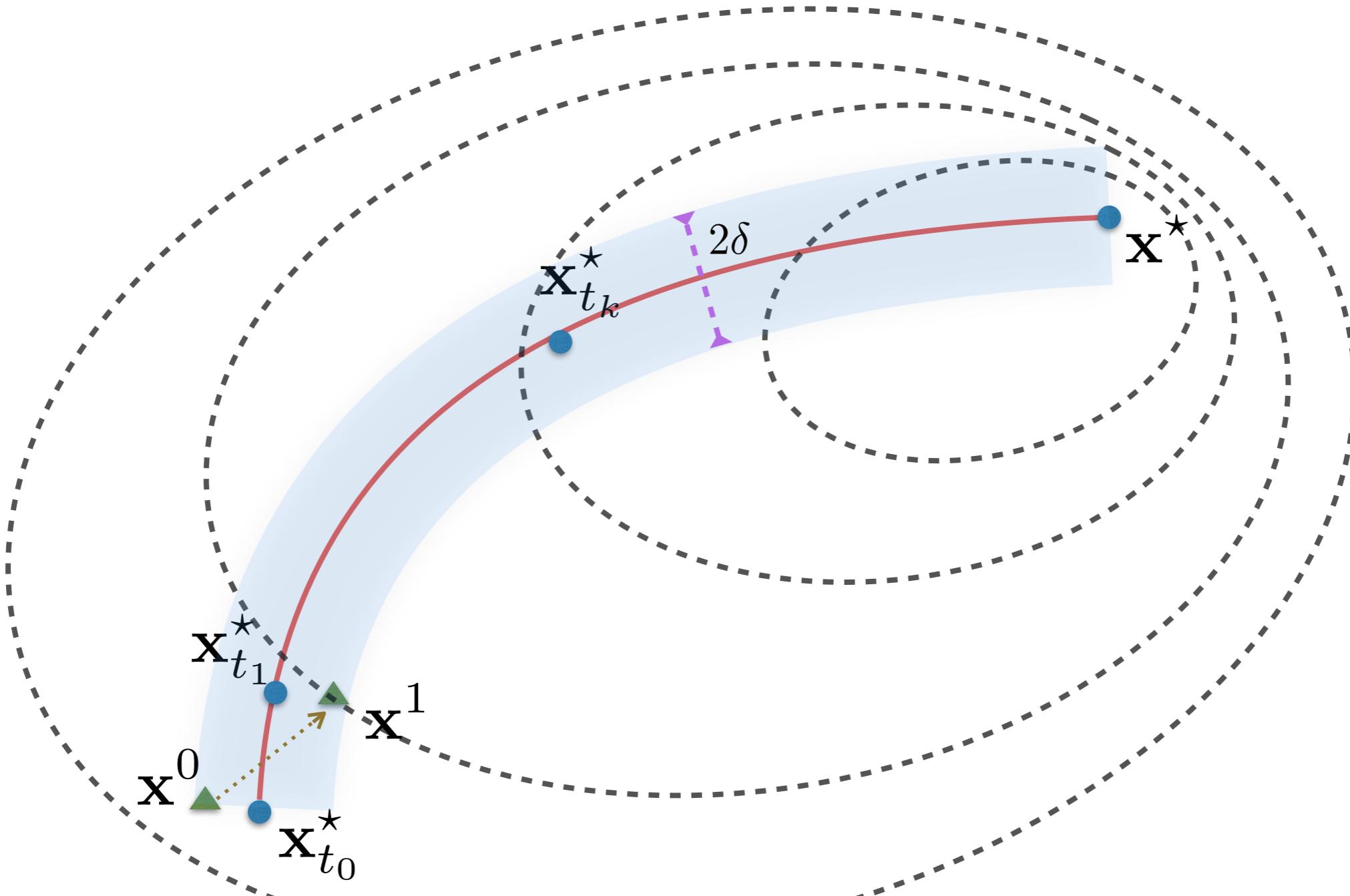
$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \Omega} g(\mathbf{x})$   
 $\mathbf{x}_{t_k}^* = \arg \min \{g(\mathbf{x}) + t_k \cdot f(\mathbf{x})\}$   
 $\mathbf{x}^k$ :  $\delta$ -approximate solution with  
 one proximal Newton step!



$$\underset{\mathbf{x} \in \text{int}(\Omega)}{\text{minimize}} \quad g(\mathbf{x}) + t \cdot f(\mathbf{x})$$

- Points on central path
- Central path trajectory
- ▲ Inexact solutions

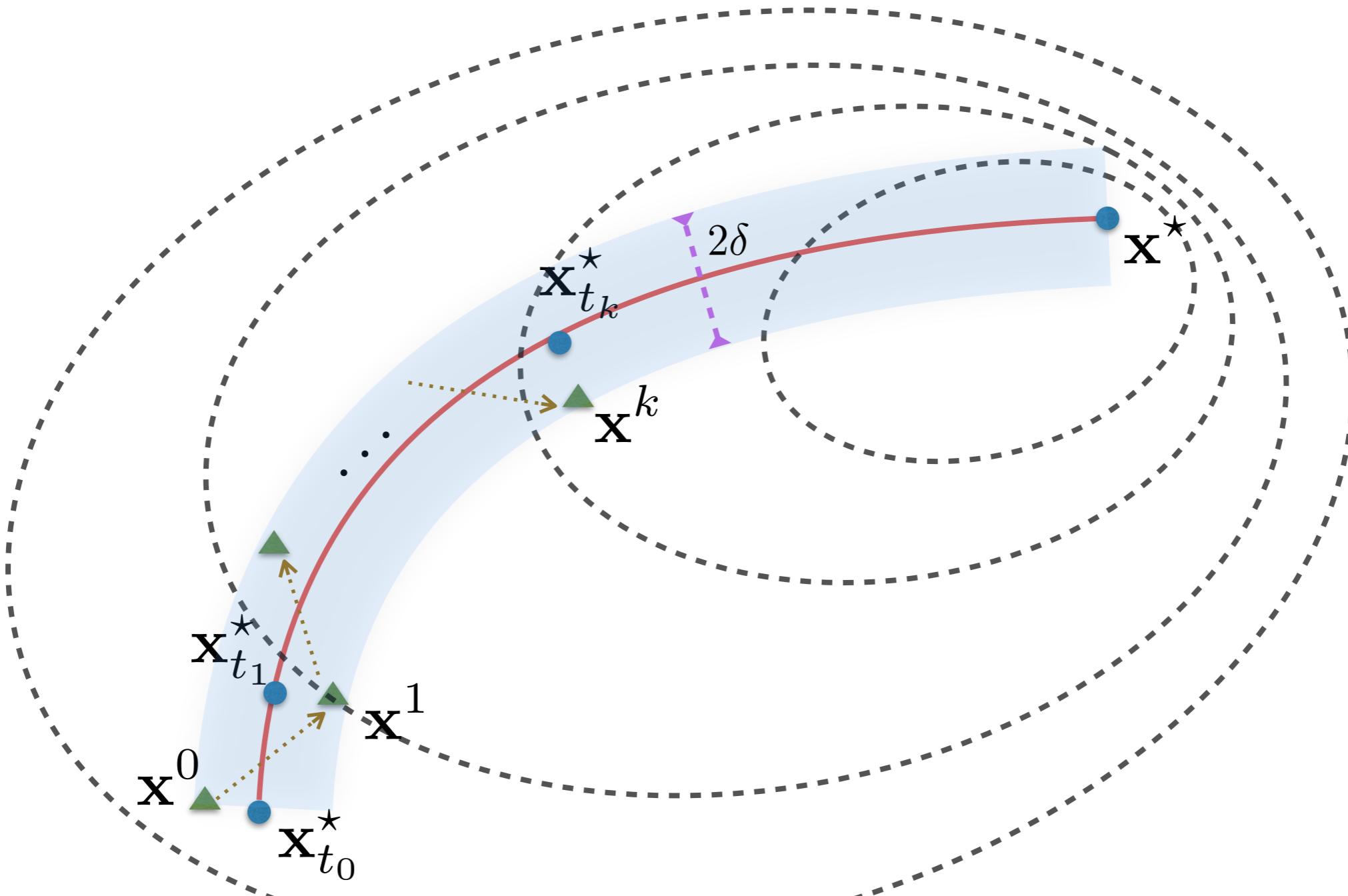
$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \Omega} g(\mathbf{x})$   
 $\mathbf{x}_{t_k}^* = \arg \min \{g(\mathbf{x}) + t_k \cdot f(\mathbf{x})\}$   
 $\mathbf{x}^k$ :  $\delta$ -approximate solution with  
 one proximal Newton step!



$$\underset{\mathbf{x} \in \text{int}(\Omega)}{\text{minimize}} \quad g(\mathbf{x}) + t \cdot f(\mathbf{x})$$

- Points on central path
- Central path trajectory
- ▲ Inexact solutions

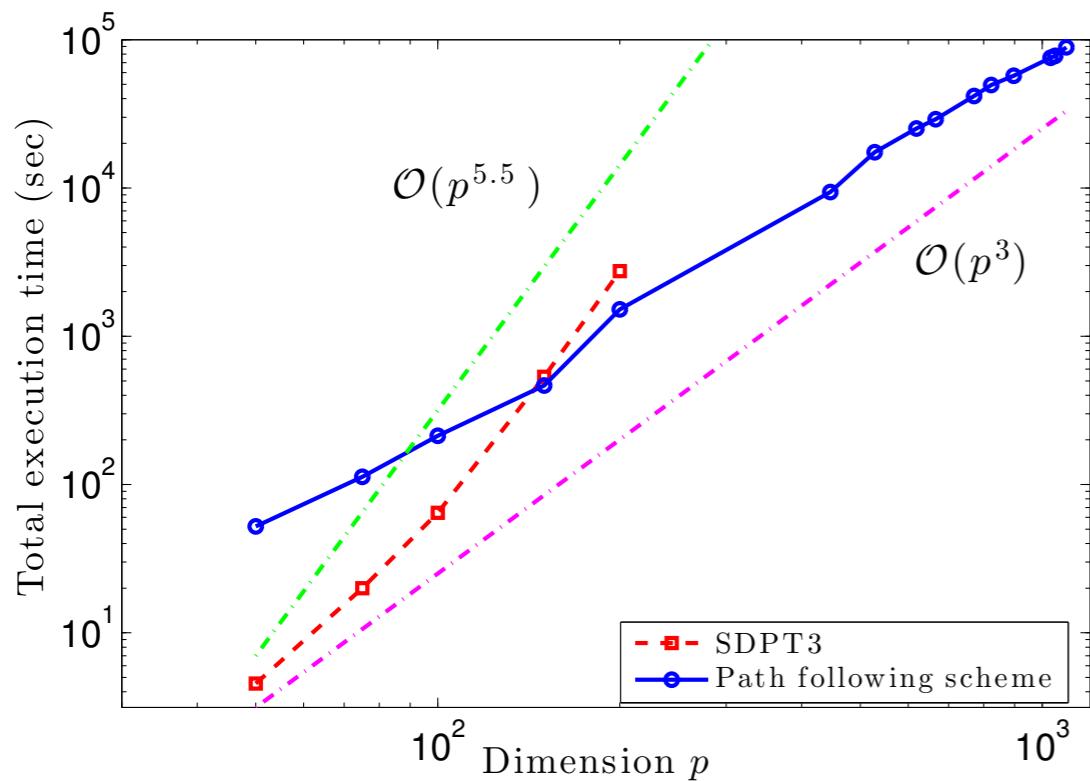
$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \Omega} g(\mathbf{x})$   
 $\mathbf{x}_{t_k}^* = \arg \min \{g(\mathbf{x}) + t_k \cdot f(\mathbf{x})\}$   
 $\mathbf{x}^k$ :  $\delta$ -approximate solution with  
 one proximal Newton step!



$$\underset{\mathbf{x} \in \text{int}(\Omega)}{\text{minimize}} \quad g(\mathbf{x}) + t \cdot f(\mathbf{x})$$

# Application #3: Max-norm clustering

$$\begin{aligned}
 & \min_{\mathbf{L}, \mathbf{R}, \mathbf{K}} \quad \|\text{vec}(\mathbf{K} - \mathbf{A})\|_1 \\
 \text{s.t.} \quad & \begin{bmatrix} \mathbf{L} & \mathbf{K} \\ \mathbf{K}^T & \mathbf{R} \end{bmatrix} \succ 0, \quad \mathbf{L}_{ii} \leq 1, \quad \mathbf{R}_{ii} \leq 1, \quad i = 1, \dots, p.
 \end{aligned}$$



DCO:

$p$	SDPT3		PF scheme	
	variables	constraints	variables	
50	15.1	2.6	10	
75	33.9	5.8	22.5	
100	60.2	10.2	40	
150	135.3	22.8	90	
200	240.4	40.4	160	

(in thousands)

# Overview

- ❑ New results in composite convex optimization
- ❑ Non-smooth + self-concordant objective function
- ❑ Extensions to path-following schemes for constrained optimization
- ❑ Several extensions to be exploited

Thank you!