

# Newton step

$$\Delta x_{\text{nt}} = -\nabla^2 f(x)^{-1} \nabla f(x)$$

interpretations

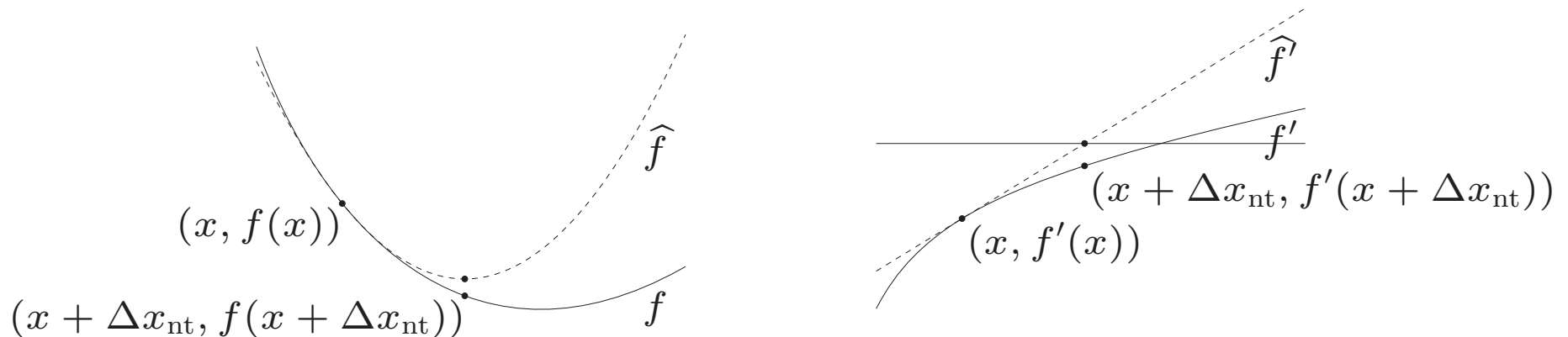
If  $f(x)$  is a quadratic,  $x + \Delta x_{\text{nt}}$  is the minimizer

- $x + \Delta x_{\text{nt}}$  minimizes second order approximation

$$\hat{f}(x + v) = f(x) + \nabla f(x)^T v + \frac{1}{2} v^T \nabla^2 f(x) v$$

- $x + \Delta x_{\text{nt}}$  solves linearized optimality condition

$$\nabla f(x + v) \approx \nabla \hat{f}(x + v) = \nabla f(x) + \nabla^2 f(x) v = 0$$

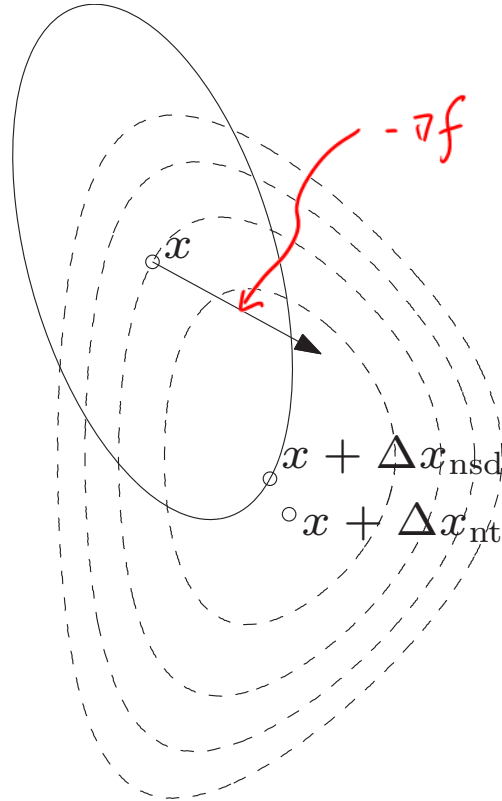


$$\nabla f^T \Delta x_{nt} = -\nabla f^T \nabla^2 f^{-1} \nabla f < 0 \quad \text{unless } \nabla f = 0$$

assuming  $f$  is convex.  
 $-\nabla x_{nt}$  may not be descent for nonconvex

- $\Delta x_{nt}$  is steepest descent direction at  $x$  in local Hessian norm

$$\|u\|_{\nabla^2 f(x)} = (u^T \nabla^2 f(x) u)^{1/2}$$



dashed lines are contour lines of  $f$ ; ellipse is  $\{x + v \mid v^T \nabla^2 f(x) v = 1\}$

arrow shows  $-\nabla f(x)$

# Newton decrement

$$\lambda(x) = (\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x))^{1/2}$$

a measure of the proximity of  $x$  to  $x^*$

## properties

- gives an estimate of  $f(x) - p^*$ , using quadratic approximation  $\hat{f}$ :

$$f(x) - \inf_y \hat{f}(y) = \frac{1}{2} \lambda(x)^2$$

- equal to the norm of the Newton step in the quadratic Hessian norm

$$\lambda(x) = (\Delta x_{\text{nt}}^T \nabla^2 f(x) \Delta x_{\text{nt}})^{1/2}$$

- directional derivative in the Newton direction:  $\nabla f(x)^T \Delta x_{\text{nt}} = -\lambda(x)^2$
- affine invariant (unlike  $\|\nabla f(x)\|_2$ )

Newton step: independent of linear change of coordinates

Assume  $T \in \mathbb{R}^{n \times n}$  is nonsingular.

Let  $\bar{f} = f(Ty)$ ,  $x = Ty$ . Then

$$\nabla \bar{f} = T^T \nabla f \quad \nabla^2 \bar{f} = T^T \nabla^2 f T$$

$$\begin{aligned} \Delta y_{nt} &= - (T^T \nabla^2 f T)^{-1} \cdot T^T \nabla f \\ &= - T^{-1} \nabla^2 f \nabla f \\ &= T^{-1} \Delta x_{nt} \end{aligned}$$

$$\Rightarrow x + \Delta x_{nt} = T (y + \Delta y_{nt})$$

i.e. Newton steps of  $f$  and  $\bar{f}$  are related by the same linear transformation

- This property does not hold for gradient descent

# Newton's method

---

**given** a starting point  $x \in \text{dom } f$ , tolerance  $\epsilon > 0$ .

**repeat**

1. *Compute the Newton step and decrement.*

$$\Delta x_{\text{nt}} := -\nabla^2 f(x)^{-1} \nabla f(x); \quad \lambda^2 := \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x).$$

2. *Stopping criterion. quit* if  $\lambda^2/2 \leq \epsilon$ .

3. *Line search.* Choose step size  $t$  by backtracking line search.

4. *Update.*  $x := x + t\Delta x_{\text{nt}}$ .

---

affine invariant, *i.e.*, independent of linear changes of coordinates:

Newton iterates for  $\tilde{f}(y) = f(Ty)$  with starting point  $y^{(0)} = T^{-1}x^{(0)}$  are

$$y^{(k)} = T^{-1}x^{(k)}$$

# Classical convergence analysis

## assumptions

- $f$  strongly convex on  $S$  with constant  $m$
- $\nabla^2 f$  is Lipschitz continuous on  $S$ , with constant  $L > 0$ :

$$\|\nabla^2 f(x) - \nabla^2 f(y)\|_2 \leq L \|x - y\|_2$$

*↑ bound on the third derivative*

( $L$  measures how well  $f$  can be approximated by a quadratic function)

*e.g.,  $L=0$  if  $f(x)$  is a quadratic*

**outline:** there exist constants  $\eta \in (0, m^2/L)$ ,  $\gamma > 0$  such that

- if  $\|\nabla f(x)\|_2 \geq \eta$ , then  $f(x^{(k+1)}) - f(x^{(k)}) \leq -\gamma$
- if  $\|\nabla f(x)\|_2 < \eta$ , then

$$\frac{L}{2m^2} \|\nabla f(x^{(k+1)})\|_2 \leq \left( \frac{L}{2m^2} \|\nabla f(x^{(k)})\|_2 \right)^2 \leftarrow \text{quadratic convergence region}$$

## damped Newton phase ( $\|\nabla f(x)\|_2 \geq \eta$ )

- most iterations require backtracking steps
- function value decreases by at least  $\gamma$
- if  $p^* > -\infty$ , this phase ends after at most  $(f(x^{(0)}) - p^*)/\gamma$  iterations

## quadratically convergent phase ( $\|\nabla f(x)\|_2 < \eta$ )

- all iterations use step size  $t = 1$
- $\|\nabla f(x)\|_2$  converges to zero quadratically: if  $\|\nabla f(x^{(k)})\|_2 < \eta$ , then

$$\frac{L}{2m^2} \|\nabla f(x^l)\|_2 \leq \left( \frac{L}{2m^2} \|\nabla f(x^k)\|_2 \right)^{2^{l-k}} \leq \left( \frac{1}{2} \right)^{2^{l-k}}, \quad l \geq k$$

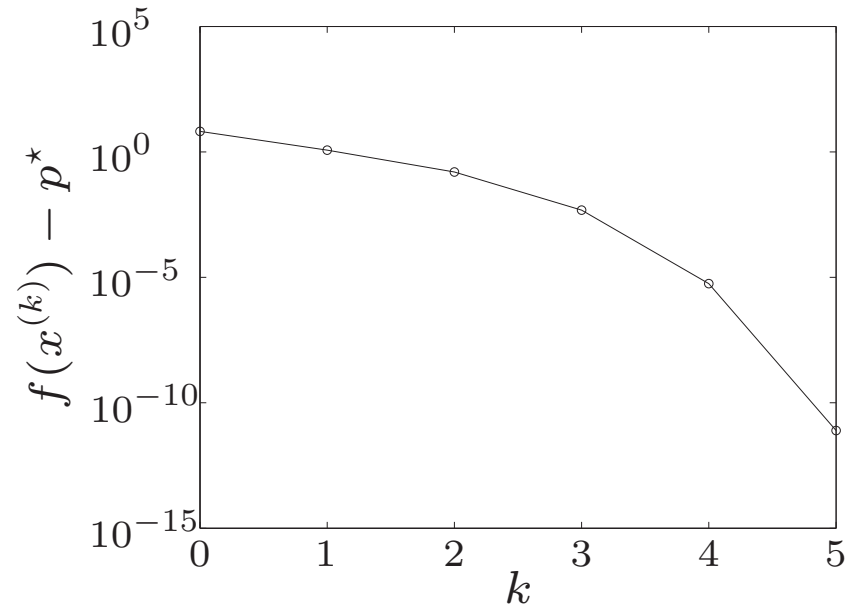
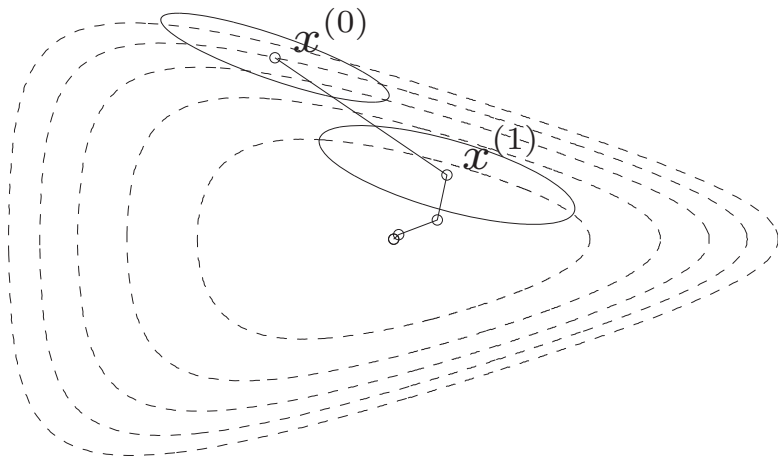
**conclusion:** number of iterations until  $f(x) - p^* \leq \epsilon$  is bounded above by

$$\frac{f(x^{(0)}) - p^*}{\gamma} + \underbrace{\log_2 \log_2(\epsilon_0/\epsilon)}$$

- $\gamma, \epsilon_0$  are constants that depend on  $m, L, x^{(0)}$
- second term is small (of the order of 6) and almost constant for practical purposes
- in practice, constants  $m, L$  (hence  $\gamma, \epsilon_0$ ) are usually unknown
- provides qualitative insight in convergence properties (*i.e.*, explains two algorithm phases)

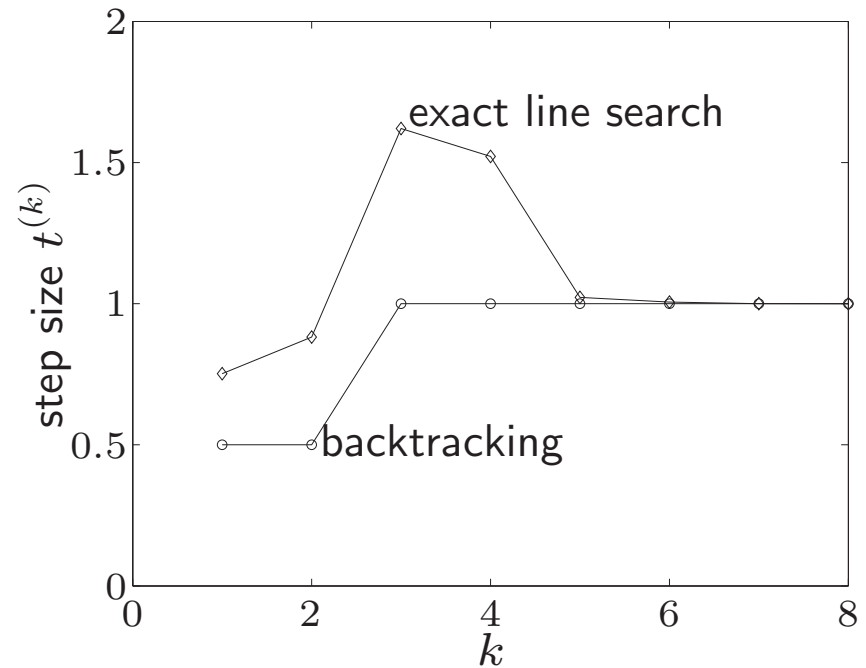
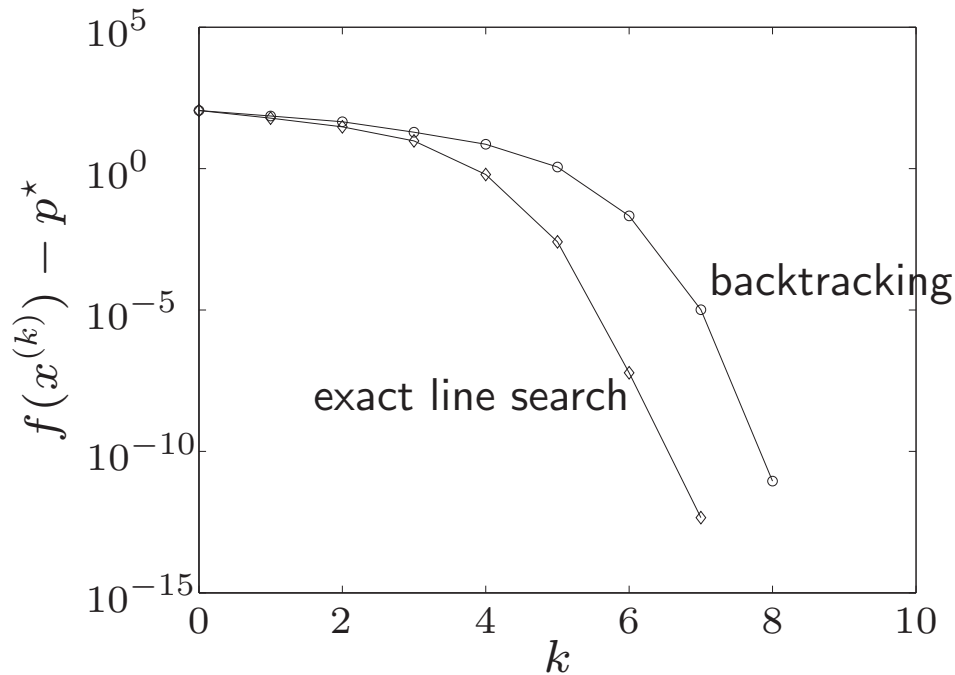
# Examples

example in  $\mathbf{R}^2$  (page 10–9)



- backtracking parameters  $\alpha = 0.1, \beta = 0.7$
- converges in only 5 steps
- quadratic local convergence

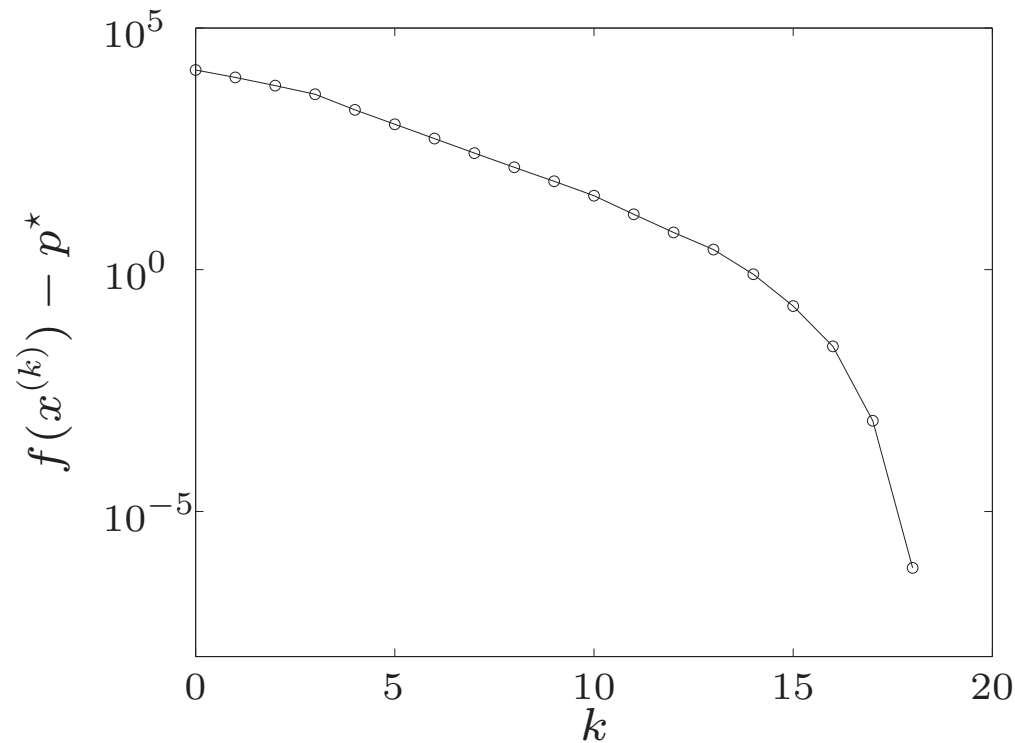
## example in $\mathbf{R}^{100}$ (page 10–10)



- backtracking parameters  $\alpha = 0.01$ ,  $\beta = 0.5$
- backtracking line search almost as fast as exact l.s. (and much simpler)
- clearly shows two phases in algorithm

example in  $\mathbf{R}^{10000}$  (with sparse  $a_i$ )

$$f(x) = - \sum_{i=1}^{10000} \log(1 - x_i^2) - \sum_{i=1}^{100000} \log(b_i - a_i^T x)$$



- backtracking parameters  $\alpha = 0.01$ ,  $\beta = 0.5$ .
- performance similar as for small examples

# Summary of Newton's method

- Convergence is typically fast once the algorithm enters the quadratic convergence phase; which depends on  $\eta \leq m^2/2$ . If  $\eta$  is very small, quadratic convergence regime can be small
- Affine invariant
- Scales well with prob. size

Disadvantage:

- computation requires

{ evaluation of  $\nabla^2 f$   
solve  $\nabla^2 f d = -\nabla f$

# Self-concordance

## shortcomings of classical convergence analysis

- depends on unknown constants ( $m, L, \dots$ )
- bound is not affinely invariant, although Newton's method is

## convergence analysis via self-concordance (Nesterov and Nemirovski)

- does not depend on any unknown constants
- gives affine-invariant bound
- applies to special class of convex functions ('self-concordant' functions)
- developed to analyze polynomial-time interior-point methods for convex optimization

# Self-concordant functions

## definition

- convex  $f : \mathbf{R} \rightarrow \mathbf{R}$  is self-concordant if  $|f'''(x)| \leq 2f''(x)^{3/2}$  for all  $x \in \text{dom } f$
- $f : \mathbf{R}^n \rightarrow \mathbf{R}$  is self-concordant if  $g(t) = f(x + tv)$  is self-concordant for all  $x \in \text{dom } f, v \in \mathbf{R}^n$

chosen for convenience  
if  $|f'''(x)| \leq K f''(x)^{3/2}$   
 $\tilde{f}(x) = K^2/4 f(x)$   
is s.c. following definition

$f: \mathbf{R}^n \rightarrow \mathbf{R}$  is self-concordant if  
it is s.c. along every line

## examples on $\mathbf{R}$

- linear and quadratic functions
- negative logarithm  $f(x) = -\log x$
- negative entropy plus negative logarithm:  $f(x) = x \log x - \log x$

**affine invariance:** if  $f : \mathbf{R} \rightarrow \mathbf{R}$  is s.c., then  $\tilde{f}(y) = f(ay + b)$  is s.c.:

$$\tilde{f}'''(y) = a^3 f'''(ay + b), \quad \tilde{f}''(y) = a^2 f''(ay + b)$$

# Self-concordant calculus

## properties

- preserved under positive scaling  $\alpha \geq 1$ , and sum
- preserved under composition with affine function
- if  $g$  is convex with  $\text{dom } g = \mathbf{R}_{++}$  and  $|g'''(x)| \leq 3g''(x)/x$  then

$$f(x) = \log(-g(x)) - \log x$$

is self-concordant

**examples:** properties can be used to show that the following are s.c.

- $f(x) = -\sum_{i=1}^m \log(b_i - a_i^T x)$  on  $\{x \mid a_i^T x < b_i, i = 1, \dots, m\}$
- $f(X) = -\log \det X$  on  $\mathbf{S}_{++}^n$
- $f(x, y) = -\log(y^2 - x^T x)$  on  $\{(x, y) \mid \|x\|_2 < y\}$

# Convergence analysis for self-concordant functions

**summary:** there exist constants  $\eta \in (0, 1/4]$ ,  $\gamma > 0$  such that

- if  $\lambda(x) > \eta$ , then

$$f(x^{(k+1)}) - f(x^{(k)}) \leq -\gamma$$

- if  $\lambda(x) \leq \eta$ , then

$$2\lambda(x^{(k+1)}) \leq \left(2\lambda(x^{(k)})\right)^2$$

( $\eta$  and  $\gamma$  only depend on backtracking parameters  $\alpha, \beta$ )

**complexity bound:** number of Newton iterations bounded by

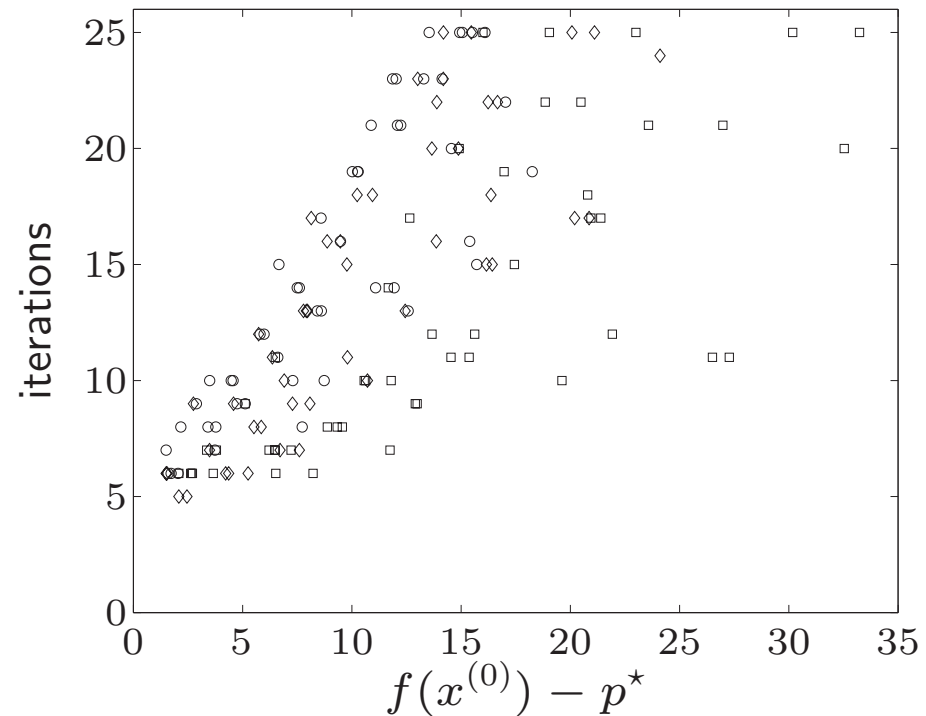
$$\frac{f(x^{(0)}) - p^*}{\gamma} + \log_2 \log_2(1/\epsilon)$$

for  $\alpha = 0.1$ ,  $\beta = 0.8$ ,  $\epsilon = 10^{-10}$ , bound evaluates to  $375(f(x^{(0)}) - p^*) + 6$

**numerical example:** 150 randomly generated instances of

$$\text{minimize } f(x) = -\sum_{i=1}^m \log(b_i - a_i^T x)$$

- :  $m = 100, n = 50$
- :  $m = 1000, n = 500$
- ◇:  $m = 1000, n = 50$



- number of iterations much smaller than  $375(f(x^{(0)}) - p^*) + 6$
- bound of the form  $c(f(x^{(0)}) - p^*) + 6$  with smaller  $c$  (empirically) valid

# Implementation

main effort in each iteration: evaluate derivatives and solve Newton system

$$H\Delta x = g$$

where  $H = \nabla^2 f(x)$ ,  $g = -\nabla f(x)$

## via Cholesky factorization

$$H = LL^T, \quad \Delta x_{\text{nt}} = L^{-T}L^{-1}g, \quad \lambda(x) = \|L^{-1}g\|_2$$

- cost  $(1/3)n^3$  flops for unstructured system
- cost  $\ll (1/3)n^3$  if  $H$  sparse, banded

## example of dense Newton system with structure

$$f(x) = \sum_{i=1}^n \psi_i(x_i) + \psi_0(Ax + b), \quad H = D + A^T H_0 A$$

- assume  $A \in \mathbf{R}^{p \times n}$ , dense, with  $p \ll n$
- $D$  diagonal with diagonal elements  $\psi_i''(x_i)$ ;  $H_0 = \nabla^2 \psi_0(Ax + b)$

**method 1:** form  $H$ , solve via dense Cholesky factorization: (cost  $(1/3)n^3$ )

**method 2** (page 9–15): factor  $H_0 = L_0 L_0^T$ ; write Newton system as

$$D\Delta x + A^T L_0 w = -g, \quad L_0^T A \Delta x - w = 0$$

eliminate  $\Delta x$  from first equation; compute  $w$  and  $\Delta x$  from

$$(I + L_0^T A D^{-1} A^T L_0)w = -L_0^T A D^{-1} g, \quad D\Delta x = -g - A^T L_0 w$$

cost:  $2p^2 n$  (dominated by computation of  $L_0^T A D^{-1} A^T L_0$ )