

## 5. Duality

- Lagrange dual problem
- weak and strong duality
- geometric interpretation
- optimality conditions
- perturbation and sensitivity analysis
- examples
- generalized inequalities

# Perturbation and sensitivity analysis

## (unperturbed) optimization problem and its dual

$$\begin{array}{ll}\text{minimize} & f_0(x) \\ \text{subject to} & f_i(x) \leq 0, \quad i = 1, \dots, m \\ & h_i(x) = 0, \quad i = 1, \dots, p\end{array}$$

$$\begin{array}{ll}\text{maximize} & g(\lambda, \nu) \\ \text{subject to} & \lambda \succeq 0\end{array}$$

## perturbed problem and its dual

$$\begin{array}{ll}\text{min.} & f_0(x) \\ \text{s.t.} & f_i(x) \leq u_i, \quad i = 1, \dots, m \\ & h_i(x) = v_i, \quad i = 1, \dots, p\end{array}$$

$$\begin{array}{ll}\text{max.} & g(\lambda, \nu) - u^T \lambda - v^T \nu \\ \text{s.t.} & \lambda \succeq 0\end{array}$$

- $x$  is primal variable;  $u, v$  are parameters
- $p^*(u, v)$  is optimal value as a function of  $u, v$
- we are interested in information about  $p^*(u, v)$  that we can obtain from the solution of the unperturbed problem and its dual

## global sensitivity result

assume strong duality holds for unperturbed problem, and that  $\lambda^*, \nu^*$  are dual optimal for unperturbed problem

apply weak duality to perturbed problem:

$$\begin{aligned} p^*(u, v) &\geq g(\lambda^*, \nu^*) - u^T \lambda^* - v^T \nu^* \\ &= p^*(0, 0) - u^T \lambda^* - v^T \nu^* \end{aligned}$$

## sensitivity interpretation

- if  $\lambda_i^*$  large:  $p^*$  increases greatly if we tighten constraint  $i$  ( $u_i < 0$ )
- if  $\lambda_i^*$  small:  $p^*$  does not decrease much if we loosen constraint  $i$  ( $u_i > 0$ )
- if  $\nu_i^*$  large and positive:  $p^*$  increases greatly if we take  $v_i < 0$ ;  
if  $\nu_i^*$  large and negative:  $p^*$  increases greatly if we take  $v_i > 0$
- if  $\nu_i^*$  small and positive:  $p^*$  does not decrease much if we take  $v_i > 0$ ;  
if  $\nu_i^*$  small and negative:  $p^*$  does not decrease much if we take  $v_i < 0$

**local sensitivity:** if (in addition)  $p^*(u, v)$  is differentiable at  $(0, 0)$ , then

$$\lambda_i^* = -\frac{\partial p^*(0, 0)}{\partial u_i}, \quad \nu_i^* = -\frac{\partial p^*(0, 0)}{\partial v_i}$$

proof (for  $\lambda_i^*$ ): from global sensitivity result,  $u = t \cdot \underline{e}_i, v = 0$

$$\frac{\partial p^*(0, 0)}{\partial u_i} = \lim_{t \searrow 0} \frac{p^*(te_i, 0) - p^*(0, 0)}{t} \geq -\lambda_i^*$$

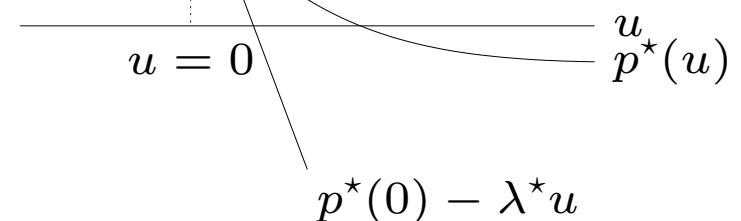
$$\frac{\partial p^*(0, 0)}{\partial u_i} = \lim_{t \nearrow 0} \frac{p^*(te_i, 0) - p^*(0, 0)}{t} \leq -\lambda_i^*$$

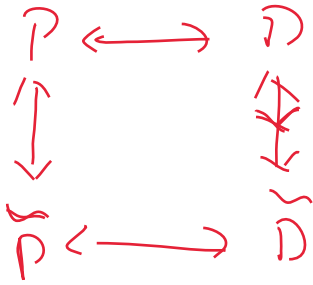
hence, equality

$f_i(\bar{x}) < 0$ , inactive  $\Rightarrow \lambda_i^* = 0$  ( $\lambda_i^* - f_i(\bar{x}) = 0$ )  
 $f_i(\bar{x}) = 0$ , active  $\Rightarrow \lambda_i^* > 0$  large, great effect  
small, little effect

$p^*(u)$  for a problem with one (inequality) constraint:

$$p^*(u) \geq p^*(0) - \lambda^* u$$





## Duality and problem reformulations

- equivalent formulations of a problem can lead to very different duals
- reformulating the primal problem can be useful when the dual is difficult to derive, or uninteresting

### common reformulations

- introduce new variables and equality constraints
- make explicit constraints implicit or vice-versa
- transform objective or constraint functions

*e.g.*, replace  $f_0(x)$  by  $\phi(f_0(x))$  with  $\phi$  convex, increasing

# Introducing new variables and equality constraints

$$\text{minimize } f_0(Ax + b)$$

- dual function is constant:  $g = \inf_x L(x) = \inf_x f_0(Ax + b) = p^*$
- we have strong duality, but dual is quite useless

## reformulated problem and its dual

$$\begin{array}{ll} \text{minimize} & f_0(y) \\ \text{subject to} & Ax + b - y = 0 \end{array}$$

$$\begin{array}{ll} \text{maximize} & b^T \nu - f_0^*(\nu) \\ \text{subject to} & A^T \nu = 0 \end{array}$$

dual function follows from

$$\begin{aligned} g(\nu) &= \inf_{x,y} (f_0(y) - \nu^T y + \nu^T Ax + b^T \nu) \\ &= \begin{cases} -f_0^*(\nu) + b^T \nu & A^T \nu = 0 \\ -\infty & \text{otherwise} \end{cases} \end{aligned}$$

**norm approximation problem:** minimize  $\|Ax - b\|$

$$\begin{array}{ll}\text{minimize} & \|y\| \\ \text{subject to} & y = Ax - b\end{array}$$

can look up conjugate of  $\|\cdot\|$ , or derive dual directly

$$\begin{aligned}g(\nu) &= \inf_{x,y} (\|y\| + \nu^T y - \nu^T Ax + b^T \nu) \\ &= \begin{cases} b^T \nu + \inf_y (\|y\| + \nu^T y) & A^T \nu = 0 \\ -\infty & \text{otherwise} \end{cases} \\ &= \begin{cases} b^T \nu & A^T \nu = 0, \quad \|\nu\|_* \leq 1 \\ -\infty & \text{otherwise} \end{cases}\end{aligned}$$

(see page 5–4)

**dual of norm approximation problem**

$$\begin{array}{ll}\text{maximize} & b^T \nu \\ \text{subject to} & A^T \nu = 0, \quad \|\nu\|_* \leq 1\end{array}$$

# Implicit constraints

**LP with box constraints:** primal and dual problem

$$\begin{array}{ll} \text{minimize} & c^T x \\ \text{subject to} & Ax = b \\ & -\mathbf{1} \preceq x \preceq \mathbf{1} \end{array} \qquad \begin{array}{ll} \text{maximize} & -b^T \nu - \mathbf{1}^T \lambda_1 - \mathbf{1}^T \lambda_2 \\ \text{subject to} & c + A^T \nu + \lambda_1 - \lambda_2 = 0 \\ & \lambda_1 \succeq 0, \quad \lambda_2 \succeq 0 \end{array}$$

**reformulation with box constraints made implicit**

$$\begin{array}{ll} \text{minimize} & f_0(x) = \begin{cases} c^T x & -\mathbf{1} \preceq x \preceq \mathbf{1} \\ \infty & \text{otherwise} \end{cases} \\ \text{subject to} & Ax = b \end{array}$$

dual function

$$\begin{aligned} g(\nu) &= \inf_{-\mathbf{1} \preceq x \preceq \mathbf{1}} (c^T x + \nu^T (Ax - b)) \\ &= -b^T \nu - \underbrace{\|A^T \nu + c\|_1} \end{aligned} \quad \text{loss} \leftrightarrow l_1$$

**dual problem:** maximize  $-b^T \nu - \|A^T \nu + c\|_1$



## Problems with generalized inequalities

$$\begin{array}{ll} \text{minimize} & f_0(x) \\ \text{subject to} & f_i(x) \preceq_{K_i} 0, \quad i = 1, \dots, m \\ & h_i(x) = 0, \quad i = 1, \dots, p \end{array}$$

$\preceq_{K_i}$  is generalized inequality on  $\mathbf{R}^{k_i}$

**definitions** are parallel to scalar case:

- Lagrange multiplier for  $f_i(x) \preceq_{K_i} 0$  is vector  $\lambda_i \in \mathbf{R}^{k_i}$
- Lagrangian  $L : \mathbf{R}^n \times \mathbf{R}^{k_1} \times \dots \times \mathbf{R}^{k_m} \times \mathbf{R}^p \rightarrow \mathbf{R}$ , is defined as

$$L(x, \lambda_1, \dots, \lambda_m, \nu) = f_0(x) + \sum_{i=1}^m \lambda_i^T f_i(x) + \sum_{i=1}^p \nu_i h_i(x)$$

$\lambda_i \in \mathbf{R}^{k_i}$

- dual function  $g : \mathbf{R}^{k_1} \times \dots \times \mathbf{R}^{k_m} \times \mathbf{R}^p \rightarrow \mathbf{R}$ , is defined as

$$g(\lambda_1, \dots, \lambda_m, \nu) = \inf_{x \in \mathcal{D}} L(x, \lambda_1, \dots, \lambda_m, \nu)$$

**lower bound property:** if  $\lambda_i \succeq_{K_i^*} 0$ , then  $g(\lambda_1, \dots, \lambda_m, \nu) \leq p^*$

proof: if  $\tilde{x}$  is feasible and  $\lambda \succeq_{K_i^*} 0$ , then

$$\begin{aligned}
 \underbrace{f_0(\tilde{x})}_{K^* = \{y \mid y^T x \leq 1, x \in K\}} &\geq f_0(\tilde{x}) + \sum_{i=1}^m \underbrace{\lambda_i^T}_{\succeq_{K_i^*} 0} \underbrace{f_i(\tilde{x})}_{\leq_{K_i} 0} + \sum_{i=1}^p \nu_i h_i(\tilde{x}) \\
 &\geq \inf_{x \in \mathcal{D}} L(x, \lambda_1, \dots, \lambda_m, \nu) \\
 &= g(\lambda_1, \dots, \lambda_m, \nu)
 \end{aligned}$$

minimizing over all feasible  $\tilde{x}$  gives  $p^* \geq g(\lambda_1, \dots, \lambda_m, \nu)$

## dual problem

$$\begin{aligned}
 &\text{maximize} && g(\lambda_1, \dots, \lambda_m, \nu) \\
 &\text{subject to} && \lambda_i \succeq_{K_i^*} 0, \quad i = 1, \dots, m
 \end{aligned}$$

- weak duality:  $p^* \geq d^*$  always
- strong duality:  $p^* = d^*$  for convex problem with constraint qualification (for example, Slater's: primal problem is strictly feasible)

# Semidefinite program

$$\langle a, b \rangle = \underline{a}^T b$$

$$\langle A, B \rangle = \text{tr}(A^T B)$$

primal SDP ( $F_i, G \in \mathbf{S}^k$ )

$$\begin{aligned} & \text{minimize} && c^T x \\ & \text{subject to} && x_1 F_1 + \cdots + x_n F_n \preceq G \end{aligned}$$

$$K^* = K = (\mathbf{S}_+^n)$$

- Lagrange multiplier is matrix  $Z \in \mathbf{S}^k$
- Lagrangian  $L(x, Z) = c^T x + \text{tr}(Z(x_1 F_1 + \cdots + x_n F_n - G))$
- dual function

$$g(Z) = \inf_x L(x, Z) = \begin{cases} -\text{tr}(GZ) & \text{tr}(F_i Z) + c_i = 0, \quad i = 1, \dots, n \\ -\infty & \text{otherwise} \end{cases}$$

dual SDP

$$\begin{aligned} & \text{maximize} && -\text{tr}(GZ) \\ & \text{subject to} && \underline{Z \succeq 0}, \quad \text{tr}(F_i Z) + c_i = 0, \quad i = 1, \dots, n \end{aligned}$$

$p^* = d^*$  if primal SDP is strictly feasible ( $\exists x$  with  $x_1 F_1 + \cdots + x_n F_n \prec G$ )

## 10. Unconstrained minimization

- terminology and assumptions
- gradient descent method
- steepest descent method
- Newton's method
- self-concordant functions
- implementation

# Unconstrained minimization

$$\text{minimize } f(x)$$

- $f$  convex, twice continuously differentiable (hence  $\text{dom } f$  open)
- we assume optimal value  $p^* = \inf_x f(x)$  is attained (and finite)

## unconstrained minimization methods

- produce sequence of points  $x^{(k)} \in \text{dom } f$ ,  $k = 0, 1, \dots$  with

$$f(x^{(k)}) \rightarrow p^*$$

- can be interpreted as iterative methods for solving optimality condition

$$\boxed{\nabla f(x^*) = 0} \quad \text{optimality condition.}$$

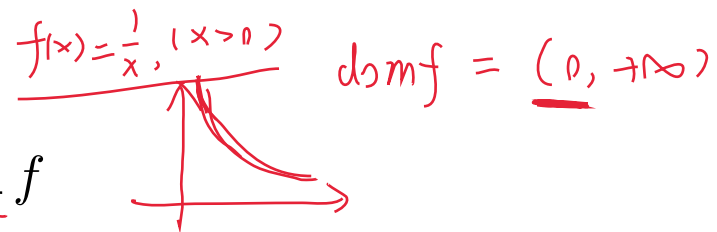
# Initial point and sublevel set

algorithms in this chapter require a starting point  $x^{(0)}$  such that

- $x^{(0)} \in \text{dom } f$
- sublevel set  $S = \{x \mid f(x) \leq f(x^{(0)})\}$  is closed

2nd condition is hard to verify, except when *all* sublevel sets are closed:

- equivalent to condition that  $\text{epi } f$  is closed
- true if  $\text{dom } f = \mathbf{R}^n$
- true if  $f(x) \rightarrow \infty$  as  $x \rightarrow \text{bd dom } f$



examples of differentiable functions with closed sublevel sets:  $-\log(x)$

$$f(x) = \log\left(\sum_{i=1}^m \exp(a_i^T x + b_i)\right),$$

$\text{dom } f = \mathbf{R}^n$

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

$\text{dom } f = \{x \mid b_i - a_i^T x > 0, \forall i\}$

$$\underline{f: \mathbb{R}^n \rightarrow \mathbb{R}} \quad \underline{\nabla f \in \mathbb{R}^n} \quad \underline{\nabla^2 f \in \mathbb{R}^{n \times n}} \quad \text{Hessian matrix } (S^n)$$

$$\underline{x \in \mathbb{R}^n}$$

## Strong convexity and implications

$f$  is strongly convex on  $S$  if there exists an  $(m > 0)$  such that  $(f(x) = \frac{1}{x}, x > 0)$

$$\underline{\nabla^2 f(x) \succeq mI} \quad \text{for all } x \in S$$

X

### implications

- for  $x, y \in S$ ,

$$(K = S_+^n) \quad (\text{convex: } f(y) \geq f(x) + \nabla f(x)^T (y-x))$$

$$f(y) = f(x) + \nabla f(x)^T (y-x) + \left[ \frac{1}{2} (y-x)^T (\nabla^2 f(z)) (y-x) \right]$$

$(z \in [x, y])$

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) + \underline{\frac{m}{2} \|x-y\|_2^2}$$



hence,  $S$  is bounded

- $p^* > -\infty$ , and for  $x \in S$ ,

$$\text{suboptimality: } \underline{f(x) - p^* \leq \frac{1}{2m} \|\nabla f(x)\|_2^2 \leq \zeta}$$

useful as stopping criterion (if you know  $m$ )

$$\underline{\hat{f}(x) \approx \hat{f}(x)}$$

$$= \sum_{n=0}^{\infty} \frac{1}{n!} f^{(n)}(a) \frac{(x-a)^n}{n!}$$

$$x \in \mathbb{R}^n$$

# Descent methods

$f$ : convex

$$\underline{x}^{(k+1)} = \underline{x}^{(k)} + \underline{t^{(k)}} \underline{\Delta x^{(k)}} \quad \text{with } \underline{f(x^{(k+1)})} < \underline{f(x^{(k)})}$$

scalar    vector

- other notations:  $\underline{x^+ = x + t\Delta x}$ ,  $x := x + t\Delta x$  ( $t > 0$ )
- $\Delta x$  is the step, or search direction;  $t$  is the step size, or step length
- from convexity,  $\underline{f(x^+) < f(x)}$  implies  $\underline{\nabla f(x)^T \Delta x < 0}$   
(i.e.,  $\Delta x$  is a descent direction)

$$f(x^+) \geq f(x) + \left( \nabla f(x)^T \frac{(x^+ - x)}{t \Delta x} \right) < 0$$

General descent method.

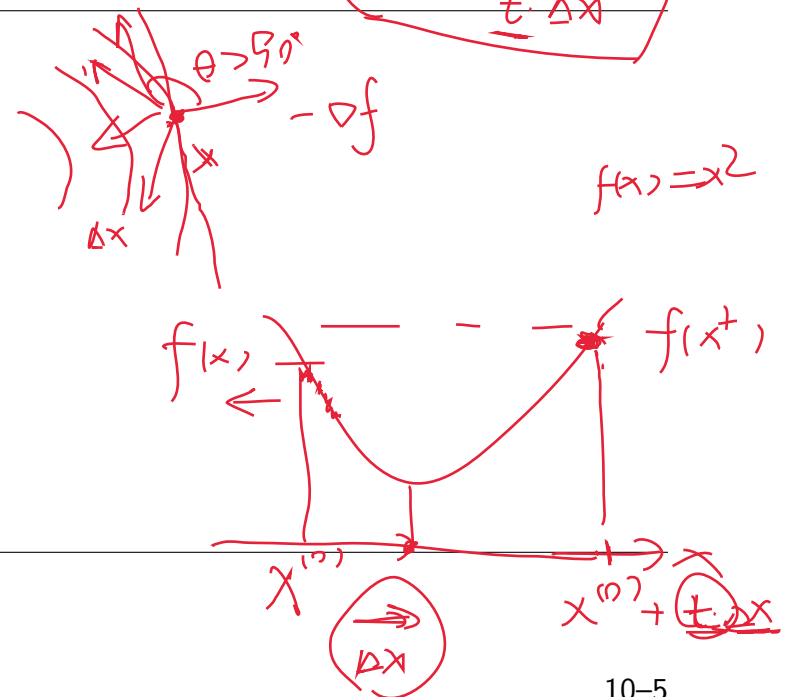
**given** a starting point  $x \in \text{dom } f$ .

**repeat**

1. Determine a descent direction  $\Delta x$ .
2. *Line search*. Choose a step size  $t > 0$ .
3. *Update*.  $x := x + t\Delta x$ .

**until** stopping criterion is satisfied.

$\leftarrow \| \nabla f \|_2 \text{ small}$





$$\min_x f(x)$$

## Line search types

$$f(x) = \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\tilde{f}(t) = \mathbb{R} \rightarrow \mathbb{R}$$

exact line search:  $t = \operatorname{argmin}_{t>0} \underline{f(x + t\Delta x)} = \underline{\tilde{f}(t)}$

backtracking line search (with parameters  $\alpha \in (0, 1/2)$ ,  $\beta \in (0, 1)$ )

- starting at  $t = 1$ , repeat  $t := \beta t$  until

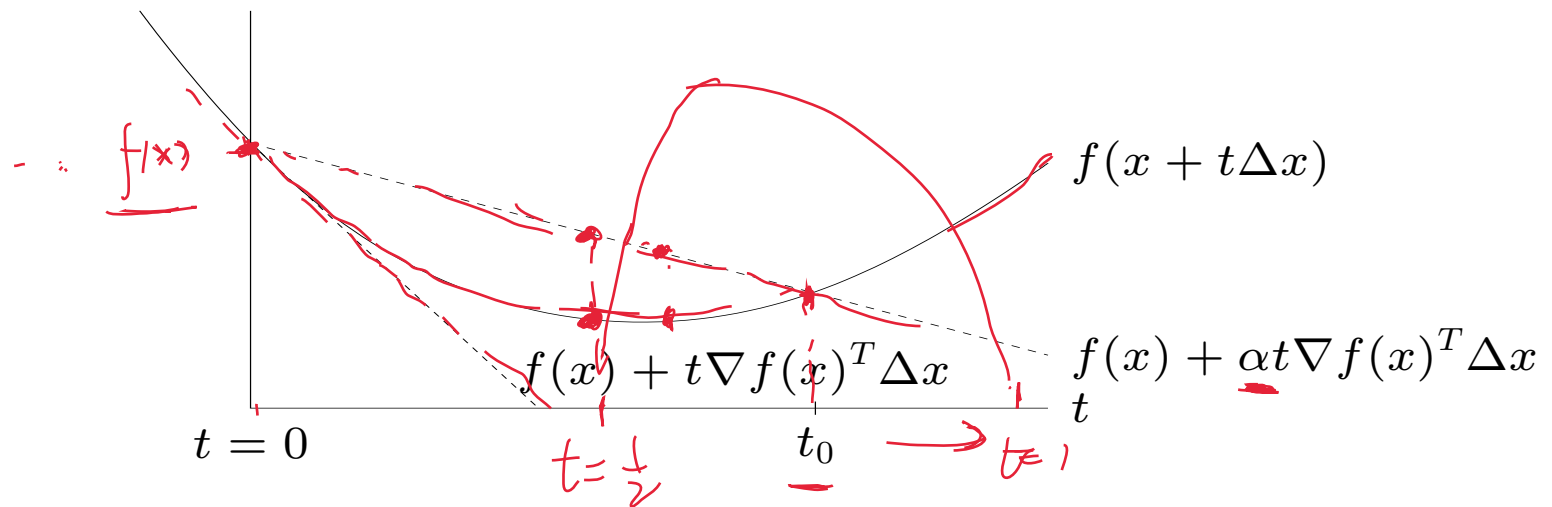
$$(\alpha = 0.1, 0.01, \beta = \frac{1}{2})$$

$$\underline{f(x + t\Delta x)} < \underline{f(x) + \alpha t \nabla f(x)^T \Delta x}$$

$$x^+ = x + t \cdot \Delta x$$

$$f(x + t\Delta x) \geq \hat{f(x + t\Delta x)} = f(x) + t \nabla f(x)^T \Delta x$$

- graphical interpretation: backtrack until  $t \leq t_0$



$$f(x + t_0 \Delta x) \approx f(x) + \alpha t_0 \nabla f(x)^T \Delta x$$

(GD)

## Gradient descent method

general descent method with  $\Delta x = -\nabla f(x)$

---

**given** a starting point  $x \in \text{dom } f$ .

**repeat**

1.  $\Delta x := -\nabla f(x)$ .
2. *Line search.* Choose step size  $t$  via exact or backtracking line search.
3. *Update.*  $x := x + t\Delta x$ .

**until** stopping criterion is satisfied.

---

- stopping criterion usually of the form  $\|\nabla f(x)\|_2 \leq \epsilon$
- convergence result: for strongly convex  $f$ ,

$$\log \left( f(x^{(k)}) - p^* \right) \leq c^k \left( f(x^{(0)}) - p^* \right) \quad \text{with } c = \frac{1}{2} \log \frac{1}{c}$$

linear -  
convergence

$c \in (0, 1)$  depends on  $m$ ,  $x^{(0)}$ , line search type

- very simple, but often very slow; rarely used in practice

$$\underline{mI \leq D^2 I \leq nI} \quad | \quad \underline{\|K(D^2 f)\| \leq \frac{n}{m}}$$

$$\begin{aligned} \underline{f(y)} &= f(x) + \nabla f(x) \cdot (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(z) (y-x) \\ &\leq f(x) + \nabla f(x) \cdot (y-x) + \frac{m}{2} \|y-x\|_2^2 \end{aligned}$$

GD

$$\Delta x = -\frac{\nabla f(x)}{\|\nabla f(x)\|_2^2}, \quad x^+ = x + t \Delta x$$
$$= x - t \cdot \frac{\nabla f(x)}{\|\nabla f(x)\|_2^2}$$

$$f(x^*) \leq f(x) + t \nabla f(x)^\top \nabla f(x) + \frac{Lt^2}{2} \|\nabla f(x)\|_2^2$$

$$\underline{f(x^*)} \leq f(x) - \frac{1}{2n} \| \nabla f(x) \|_2^2$$

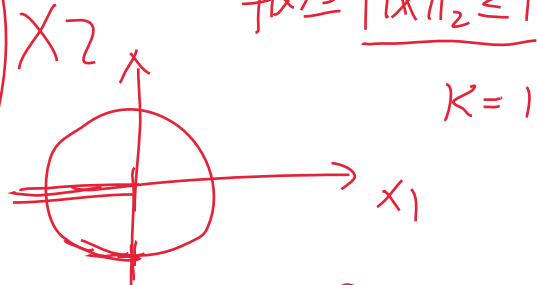
Exact l.s.:  $t_{\text{exact}} = \arg \min_{t \geq 0} f(x - t \nabla f(x_1))$

$$f(x^*) = \check{f}(t_{\text{exact}}) \leq f(x) - \frac{1}{2m} \| \nabla f(x) \|^2$$

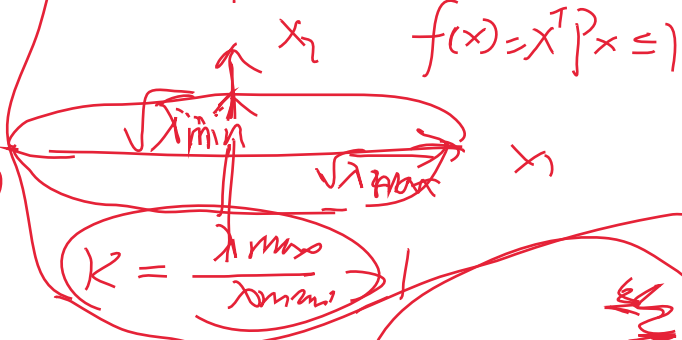
$$f(x^+) - p^* \leq \underbrace{(f(x) - p^*)}_{\leq 0} - \frac{1}{2\alpha n} \| \nabla f \|^2$$

$$f(x^*) - p^* \leq \underbrace{\left(1 - \frac{m}{M}\right)}_{C^-} (f(x) - p^*)$$

$$f: \mathbb{R}^2 \rightarrow \mathbb{R}.$$



$$f(x) = x^T P x \leq 1$$



$$\frac{f(x^{(k)}) - \bar{p}^*}{(1 - \frac{m}{n})} \leq \frac{C^k (f(x^{(1)}) - \bar{p}^*)}{(1 - \frac{m}{n})}$$

$$k \leq \frac{\log((f^{K_n}) - p^*) / 5}{\log(1/c)}$$

$k \uparrow$ , #iters  $\uparrow$   
 $k \downarrow$ , #iters  $\downarrow$

quadratic problem in  $\mathbb{R}^2$

$$f(x) = \frac{1}{2} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 1 & 0 \\ 0 & \gamma \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$P$

$$K(P) = \begin{cases} \gamma, & \gamma \geq 1 \\ \frac{1}{\gamma}, & \gamma < 1 \end{cases}$$

$$= \min \left\{ \gamma, \frac{1}{\gamma} \right\}$$

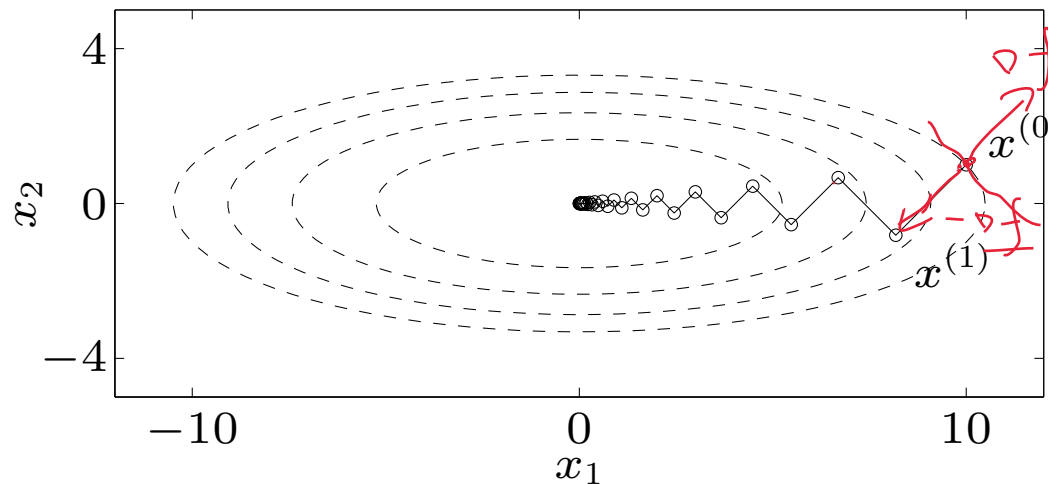
$$f(x) = (1/2)(x_1^2 + \gamma x_2^2) \quad (\gamma > 0)$$

with exact line search, starting at  $x^{(0)} = (\gamma, 1)$ :

$$x_1^{(k)} = \gamma \left( \frac{\gamma - 1}{\gamma + 1} \right)^k, \quad x_2^{(k)} = \left( -\frac{\gamma - 1}{\gamma + 1} \right)^k$$

• very slow if  $\gamma \gg 1$  or  $\gamma \ll 1$

• example for  $\gamma = 10$ :  $\gamma = 1 = K$ ,  $\gamma = 10,000 = K$



GD  $\rightarrow$  exact l.s.

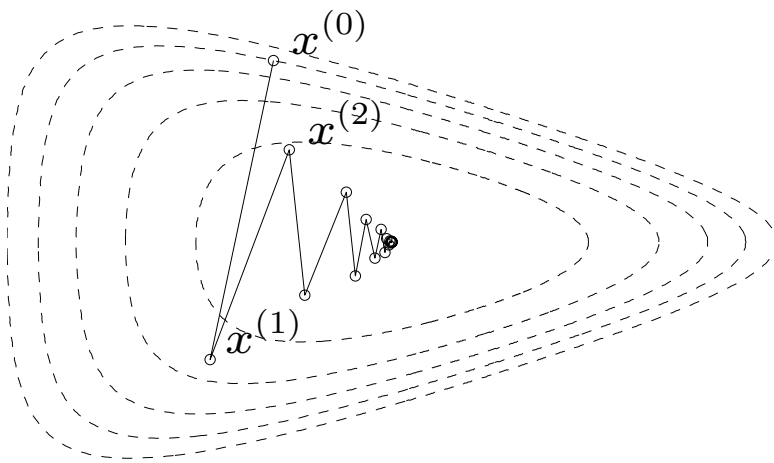
$$x^+ = x - \underline{t} \nabla f$$

# nonquadratic example

$\mathbb{R}^2$

$$f(x_1, x_2) = e^{x_1+3x_2-0.1} + e^{x_1-3x_2-0.1} + e^{-x_1-0.1}$$

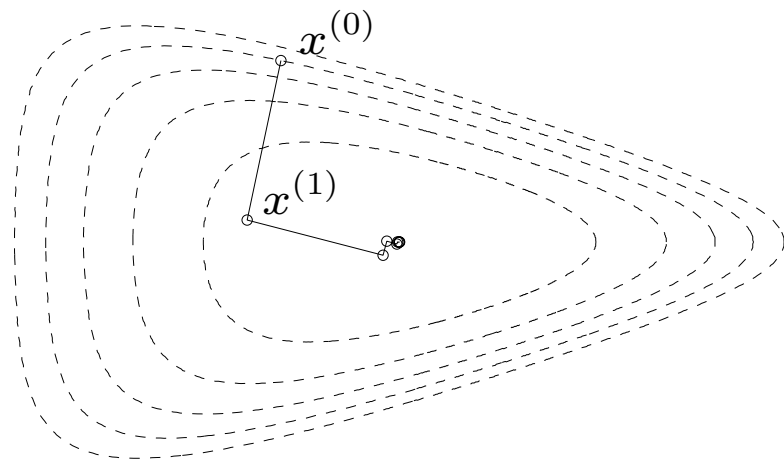
$$x^+ \equiv x - \underline{\tau} \nabla f(x)$$



backtracking line search

$$\text{#iter} = 29, \quad \varepsilon \approx 10^{-7}$$

$$(0.4)$$



exact line search

$$\text{#iter} = 15, \quad \varepsilon \approx 10^{-7}$$

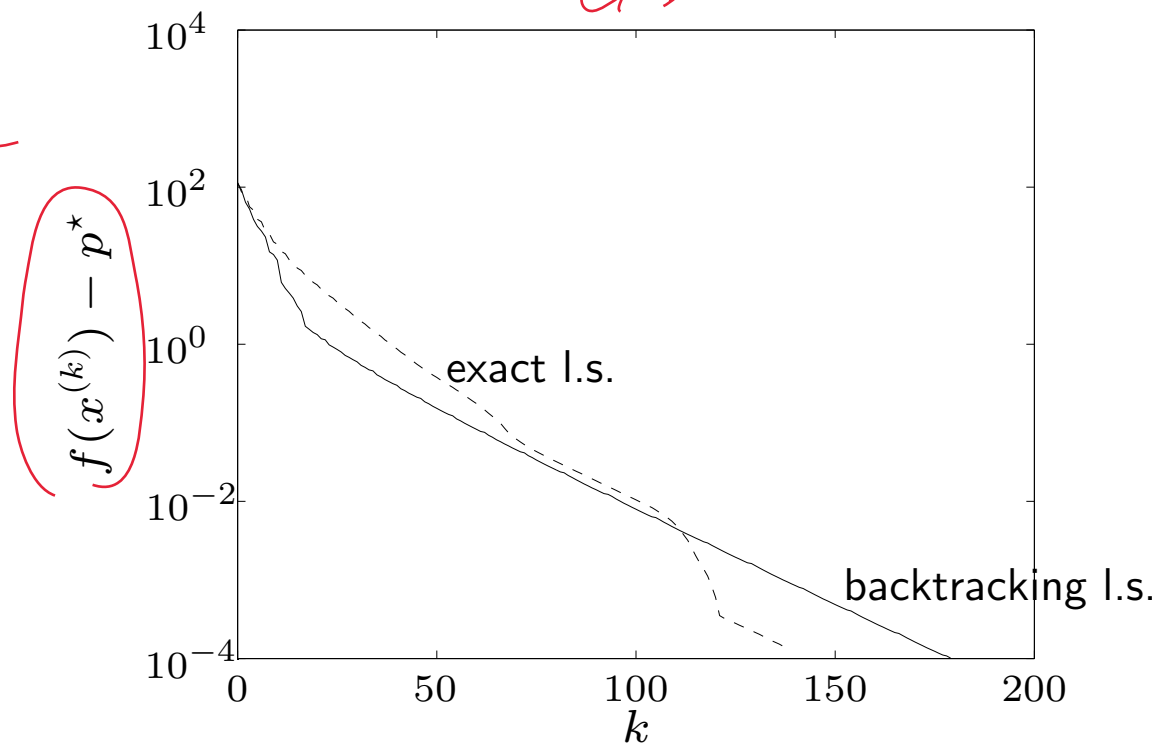
$$(0.2)$$

a problem in  $\mathbf{R}^{100}$

$$f(x) = c^T x - \sum_{i=1}^{500} \log(b_i - a_i^T x)$$

$G \nabla$

log scale



‘linear’ convergence, *i.e.*, a straight line on a semilog plot

$$\begin{aligned} f(x+v) &\approx \hat{f}(x+v) = f(x) + \nabla f(x)^T v \\ (x^T = x + v) \\ f(x+v) &< f(x) \end{aligned}$$

$$\left( \min_v \nabla f(x)^T v \text{ s.t. } \|v\|=1 \right) \Rightarrow \Delta x_{\text{nsd}} = \underset{\|v\|=1}{\operatorname{argmin}} \nabla f(x)^T v$$

## Steepest descent method

normalized steepest descent direction (at  $x$ , for norm  $\|\cdot\|$ ):

$$\Delta x_{\text{nsd}} = \underset{v \in \mathbb{R}^n}{\operatorname{argmin}} \{ \nabla f(x)^T v \mid \|v\| = 1 \}$$

$$\|z\|_* = \left\{ \begin{array}{l} \sup_{\|x\| \leq 1} z^T x \\ \inf_{\|x\| \leq 1} z^T x \end{array} \right\}$$

interpretation: for small  $v$ ,  $f(x+v) \approx f(x) + \nabla f(x)^T v$ ;

direction  $\Delta x_{\text{nsd}}$  is unit-norm step with most negative directional derivative

**(unnormalized) steepest descent direction**

$$\nabla f(x)^T \Delta x_{\text{nsd}} = -\|\nabla f(x)\|_*$$

$$\Delta x_{\text{nsd}} = \underset{v}{\operatorname{argmin}} \{ \nabla f(x)^T v \mid \|v\| \leq 1 \}$$

$$\Delta x_{\text{sd}} = \|\nabla f(x)\|_* \Delta x_{\text{nsd}}$$

satisfies  $\nabla f(x)^T \Delta x_{\text{sd}} = -\|\nabla f(x)\|_*^2$

$$\begin{aligned} &\nabla f(x)^T v \\ &= \min_{\|v\| \leq 1} \langle \nabla f(x), v \rangle \end{aligned}$$

$$\begin{aligned} &= -\max_{\|v\| \leq 1} \langle -\nabla f(x), v \rangle = -\|-\nabla f(x)\|_* \\ &= -\|\nabla f(x)\|_* \end{aligned}$$

**steepest descent method**

- general descent method with  $\Delta x = \Delta x_{\text{sd}}$
- convergence properties similar to gradient descent

## examples

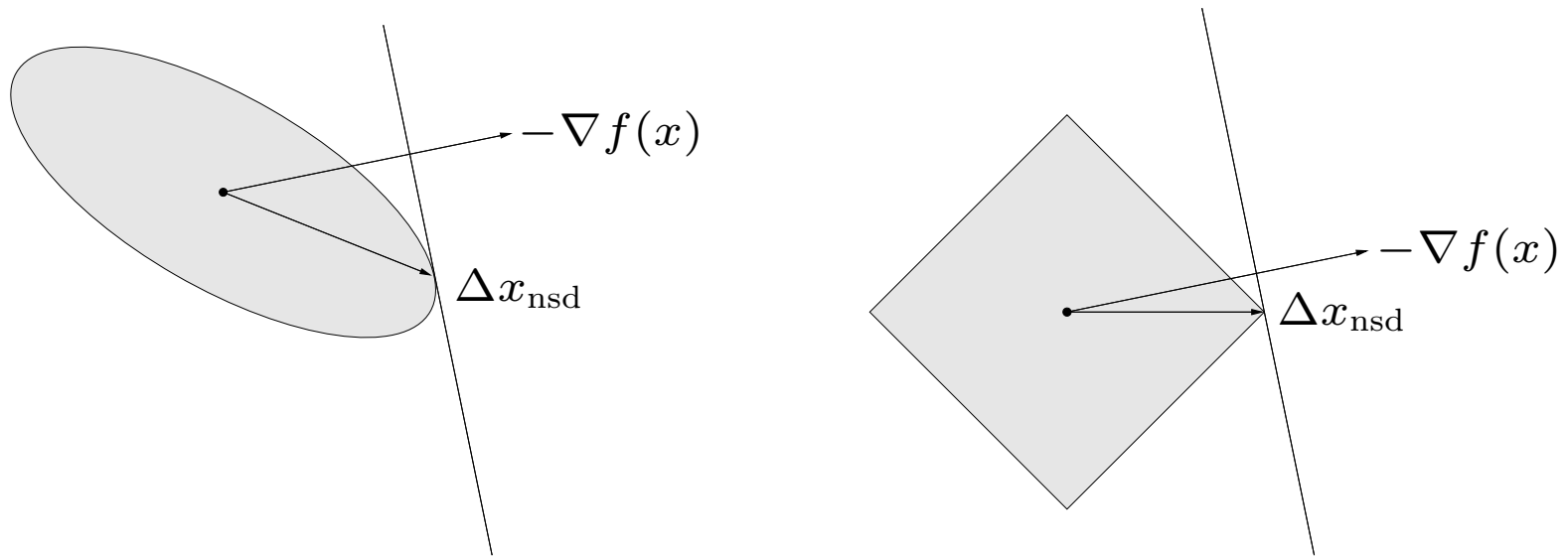
- Euclidean norm:  $\Delta x_{\text{sd}} = -\nabla f(x)$
- quadratic norm  $\|x\|_P = (x^T P x)^{1/2}$  ( $P \in \mathbf{S}_{++}^n$ ):  $\Delta x_{\text{sd}} = -P^{-1} \nabla f(x)$



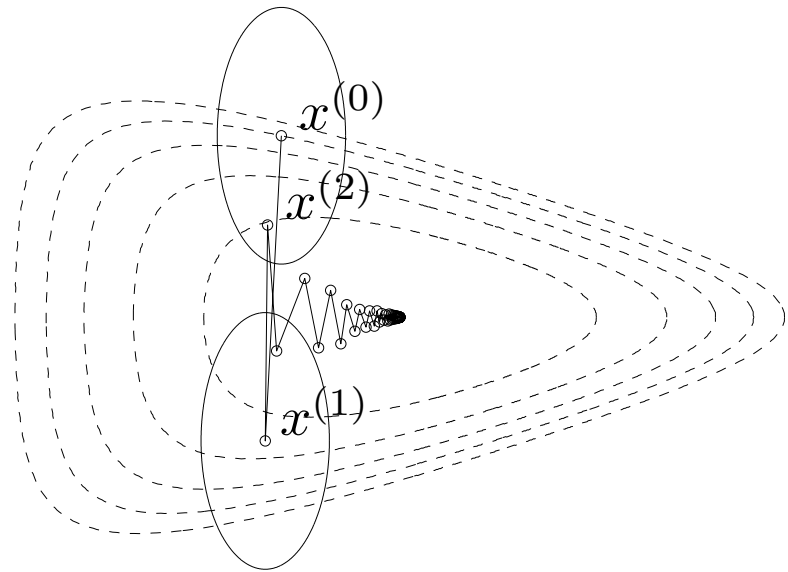
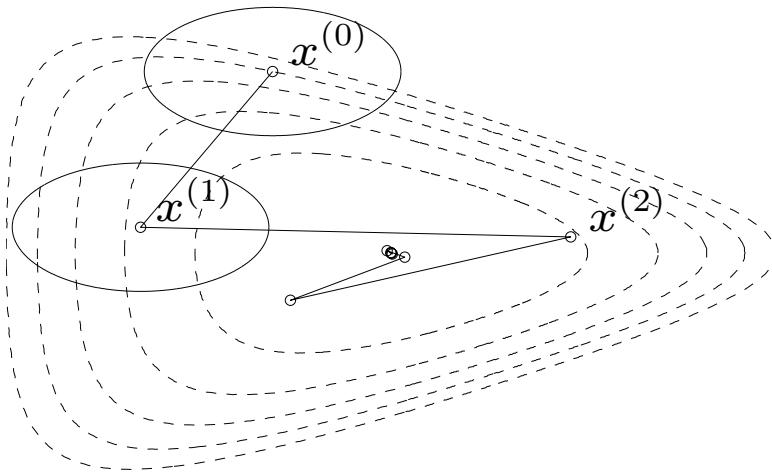
## examples

- Euclidean norm:  $\Delta x_{\text{sd}} = -\nabla f(x)$
- quadratic norm  $\|x\|_P = (x^T P x)^{1/2}$  ( $P \in \mathbf{S}_{++}^n$ ):  $\Delta x_{\text{sd}} = -P^{-1} \nabla f(x)$
- $\ell_1$ -norm:  $\Delta x_{\text{sd}} = -(\partial f(x)/\partial x_i)e_i$ , where  $|\partial f(x)/\partial x_i| = \|\nabla f(x)\|_\infty$

unit balls and normalized steepest descent directions for a quadratic norm and the  $\ell_1$ -norm:



## choice of norm for steepest descent



- steepest descent with backtracking line search for two quadratic norms
- ellipses show  $\{x \mid \|x - x^{(k)}\|_P = 1\}$
- equivalent interpretation of steepest descent with quadratic norm  $\|\cdot\|_P$ : gradient descent after change of variables  $\bar{x} = P^{1/2}x$

shows choice of  $P$  has strong effect on speed of convergence

# Newton step

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

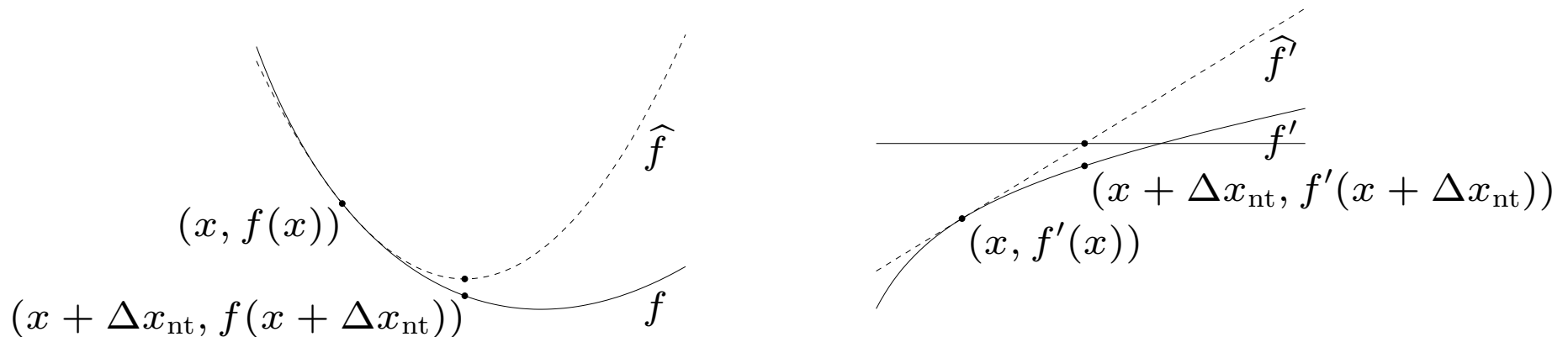
## interpretations

- $x + \Delta x_{\text{nt}}$  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

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



- $\Delta x_{\text{nt}}$  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

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

a measure of the proximity of  $x$  to  $x^\star$

## properties

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

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

# Newton decrement

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

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

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

- 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'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  $L > 0$ :

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

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

**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

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

### **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

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

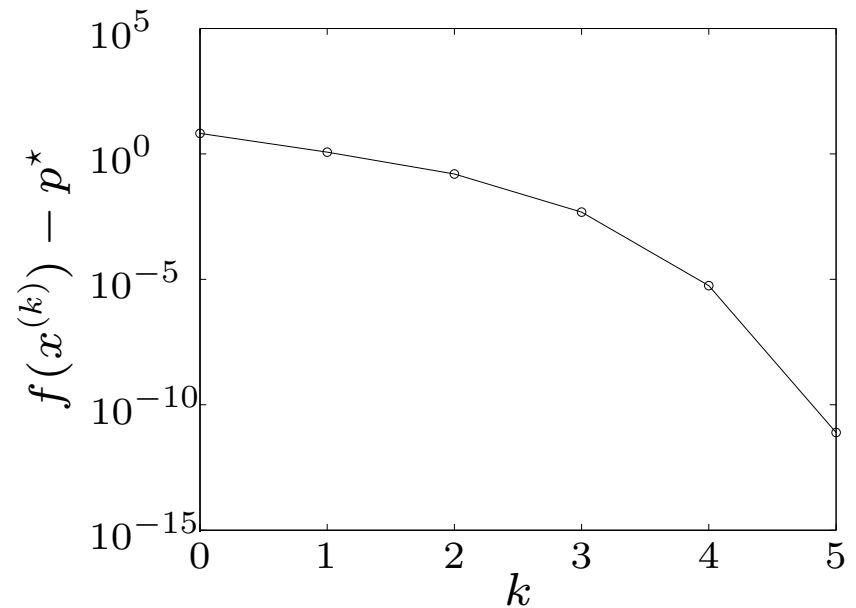
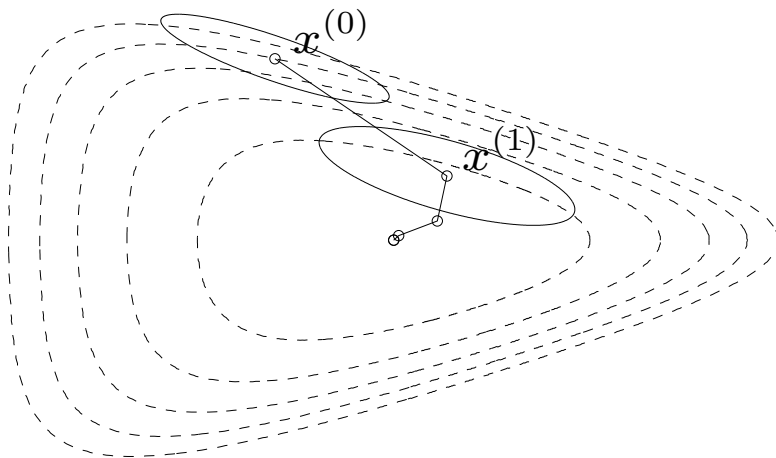
**conclusion:** number of iterations until  $f(x) - p^* \leq \epsilon$  is bounded above by

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

- $\gamma, \epsilon_0$  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)

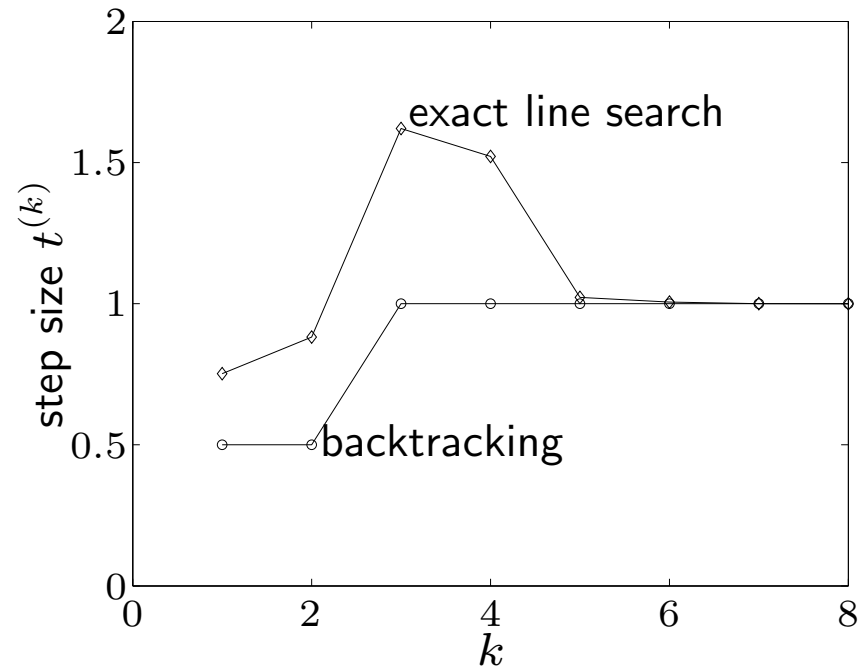
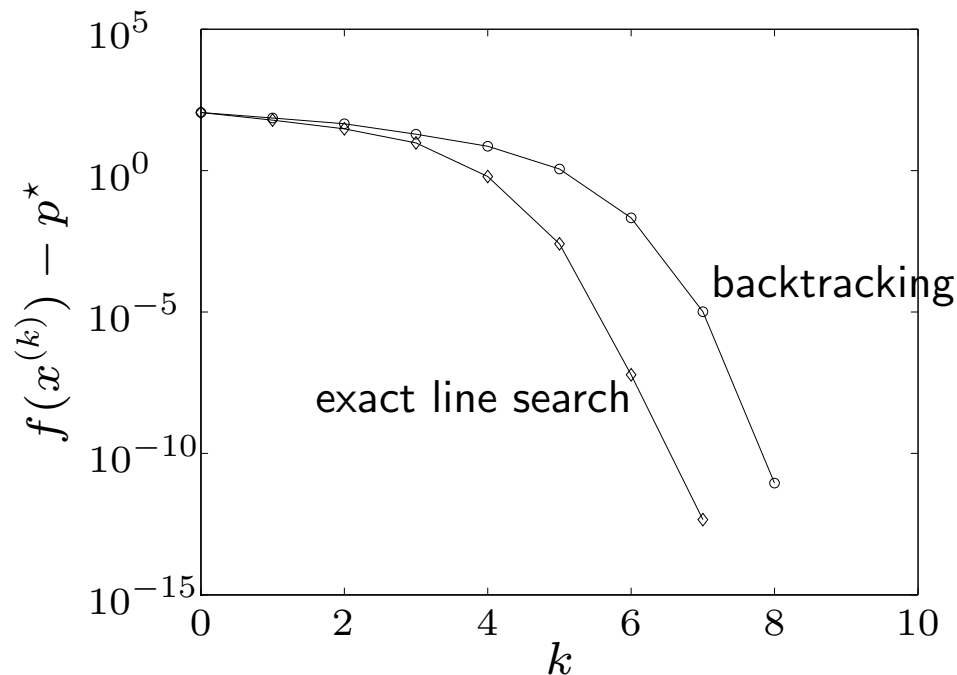
# Examples

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



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

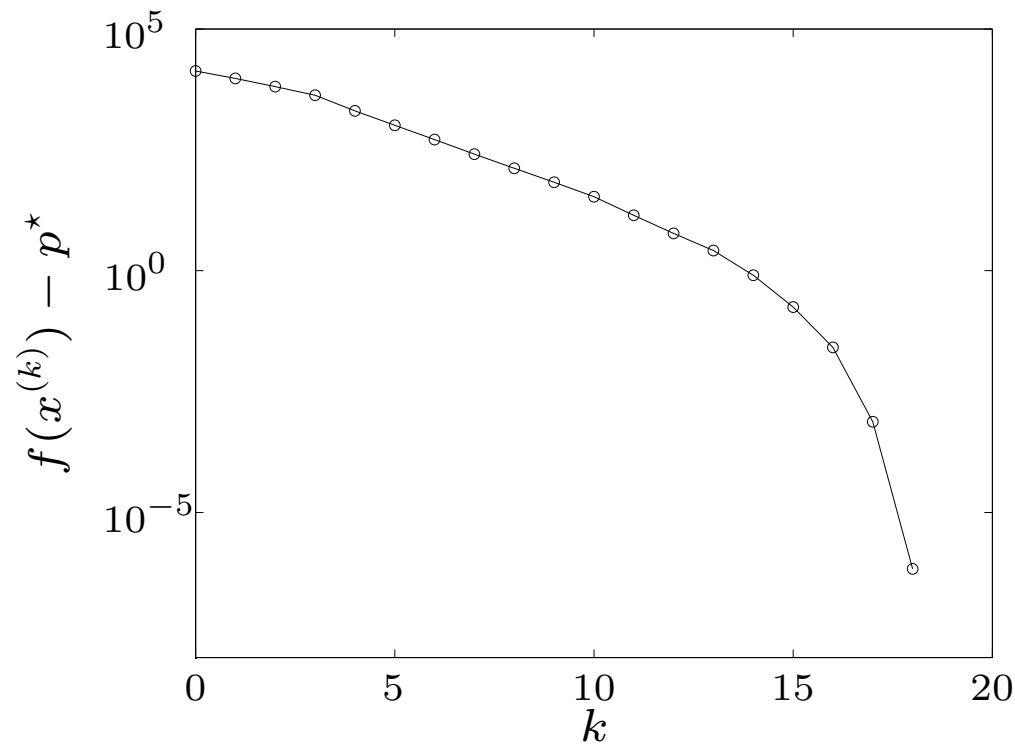
## example in $\mathbf{R}^{100}$ (page 10–10)



- 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