

## 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:

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

## Strong convexity and implications

$f$  is strongly convex on  $S$  if there exists an  $m > 0$  such that

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

### implications

- for  $x, y \in S$ ,

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

hence,  $S$  is bounded

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

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

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

# 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  $f$ ,

$$f(x^{(k)}) - p^* \leq c^k (f(x^{(0)}) - p^*)$$

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

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

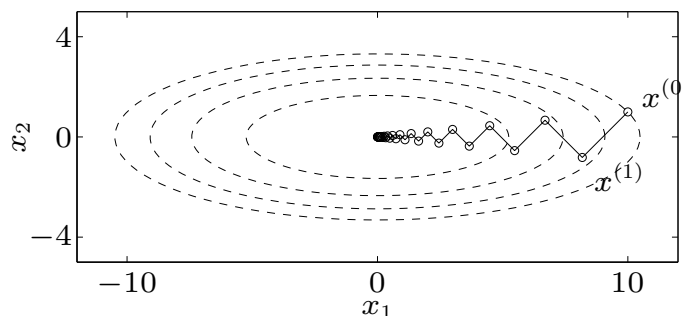
## quadratic problem in $\mathbb{R}^2$

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

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  $\gamma = 10$ :



$$f(x) = x^T A x$$

$$A = \begin{bmatrix} 1/2 & 0 \\ 0 & \gamma/2 \end{bmatrix}$$

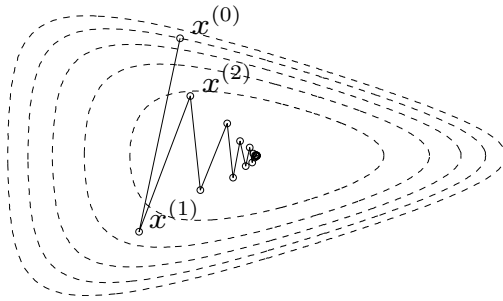
$$c = 1 - \frac{\lambda_{\max}}{\lambda_{\min}}$$

$$= \begin{cases} 1 - \frac{1}{\gamma} & \text{if } \gamma > 1 \\ 1 - \gamma & \text{if } \gamma < 1 \end{cases}$$

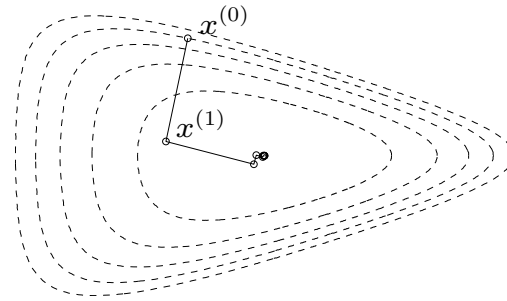
*Problem with gradient descent: curvature info is ignored...*

## nonquadratic example

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



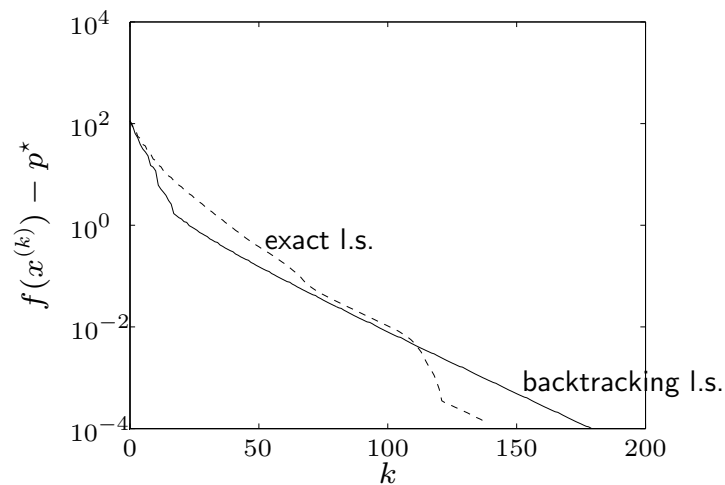
backtracking line search



exact line search

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

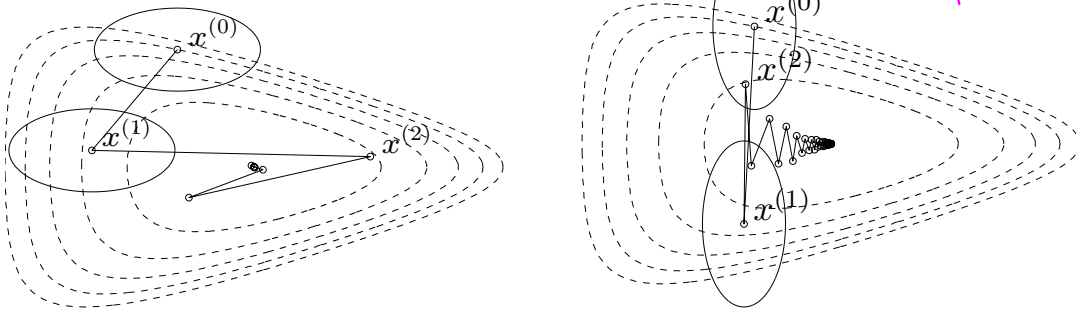
$$f(x) = c^T x - \sum_{i=1}^{500} \log(b_i - a_i^T x)$$



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

## choice of norm for steepest descent

what abt  $\|\Delta x\|_p = 1$  for  $p=1$  or  $\infty$  or matrix induced norm?  
 $\|\Delta x\|_P = (\Delta x^T P \Delta x)^{1/2}$



- steepest descent with backtracking line search for two quadratic norms
- ellipses show  $\{x \mid \|x - x^{(k)}\|_P = 1\}$ : ellipses show search space for  $\Delta x$
- 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

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

### interpretations

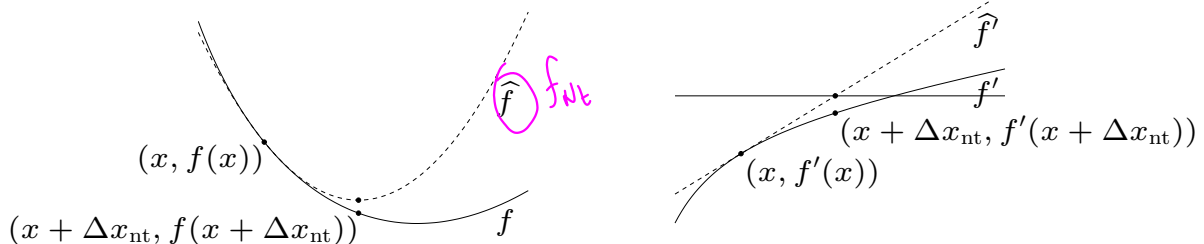
- $x + \Delta x_{nt}$  minimizes second order approximation

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

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

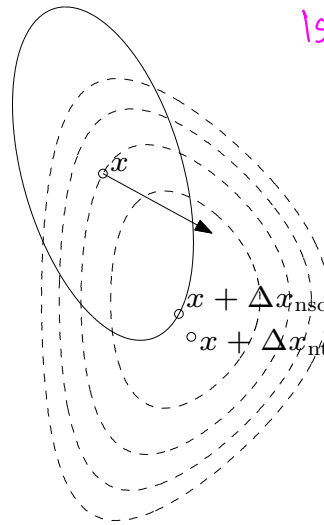
$$\nabla f(x+v) \approx \nabla \hat{f}(x+v) = \nabla f(x) + \nabla^2 f(x)v = 0$$

Recall (4.9) that tried to account for curvature but ONLY in line search stopping criterion



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

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



Is  $\Delta x_{nt}$  a descent direction?  
 $\Delta x_{nt}^T \nabla f(x) < 0$   
 Newton decrement

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

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

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