

10. Unconstrained minimization

- terminology and assumptions
- gradient descent method
- steepest descent method
- Newton's method

- 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)

We will assume that $x^* = \text{argmin}_x f(x)$ exists and is unique
Recall $p^* = f(x^*)$

unconstrained minimization methods

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

$$f(x^{(k)}) \rightarrow p^* \quad \text{as } k \rightarrow \text{infinity}$$

$x^{(0)}, x^{(1)}, \dots$ is a minimizing sequence to the problem

Algorithm stops when $f(x^{(k)}) - p^* \leq \varepsilon$, for some tolerance $\varepsilon > 0$

- can be interpreted as iterative methods for solving optimality condition

$$\nabla f(x^*) = 0$$

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 $x, y \in S$,

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

Assume f is twice differentiable

By Taylor's theorem, there exists a z in the line segment from x to y such that

$$\begin{aligned} f(y) &= f(x) + df(x)'(y-x) + \frac{1}{2} (y-x)' d^2 f(z) (y-x) \\ &\geq f(x) + df(x)'(y-x) + \frac{1}{2} (y-x)' (m I) (y-x) && \dots \text{ since } f \text{ is strongly convex} \\ &= f(x) + df(x)'(y-x) + \frac{1}{2} m \|y-x\|_2^2 \end{aligned}$$

(Taylor's theorem is a generalization of the mean value theorem, and is very related to, but is not exactly the same as Taylor series)

Descent methods

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

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

From convexity (slide 3-7):
 $f(x^+) \geq f(x) + df(x)'(x^+ - x)$
 $= f(x) + t df(x)'\Delta x$
Thus: $f(x^+) - f(x) \geq t df(x)'\Delta x$

If $f(x^+) < f(x)$ then:
 $0 > f(x^+) - f(x) \geq t df(x)'\Delta x$
Thus: $df(x)'\Delta x < 0$

General descent method.

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

repeat

1. Determine a descent direction Δx . **(Each algorithm has its own way for choosing Δx)**
2. *Line search.* Choose a step size $t > 0$.
3. *Update.* $x := x + t\Delta x$.

until stopping criterion is satisfied.

Line search types

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

backtracking line search (with parameters $\alpha \in (0, 1/2)$, $\beta \in (0, 1)$)
(one of the many inexact methods)

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

$$f(x + t\Delta x) < f(x) + \alpha t \nabla f(x)^T \Delta x \quad (\text{Armijo-Goldstein condition})$$

Since Δx is a descent direction (see previous slide) then $df(x)' \Delta x < 0$
For small t , we have:

$$f(x + t \Delta x) \approx f(x) + t df(x)' \Delta x < f(x) + \alpha t df(x)' \Delta x$$

Thus, the procedure will eventually terminate.

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 ,

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

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

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

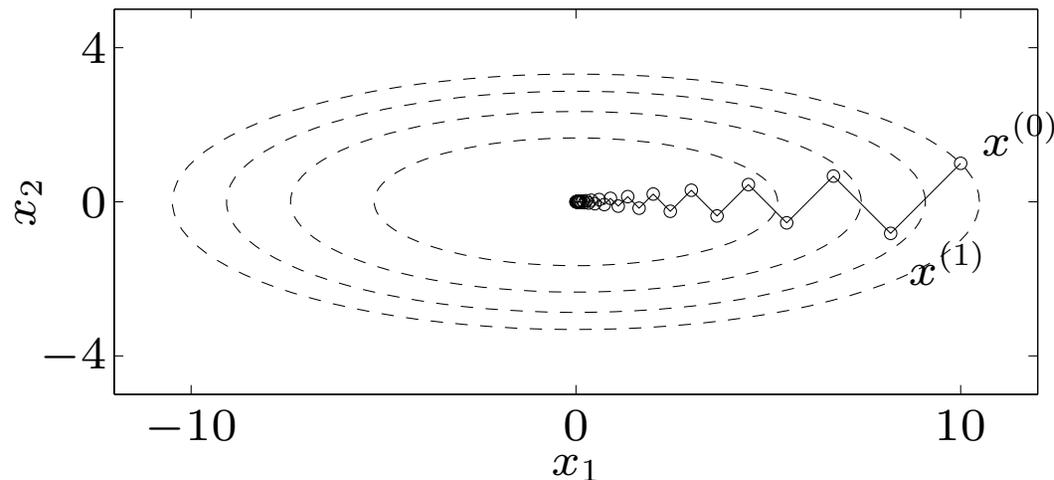
quadratic problem in \mathbf{R}^2

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

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

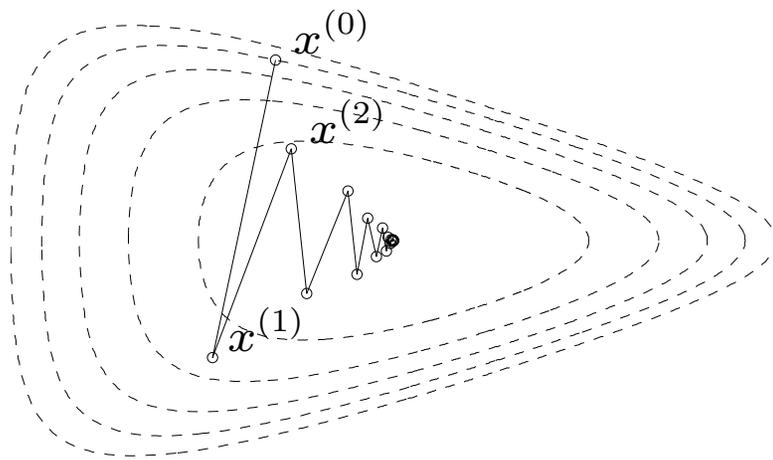
$$x_1^{(k)} = \gamma \left(\frac{\gamma - 1}{\gamma + 1} \right)^k, \quad x_2^{(k)} = \left(-\frac{\gamma - 1}{\gamma + 1} \right)^k$$

- very slow if $\gamma \gg 1$ or $\gamma \ll 1$
- example for $\gamma = 10$:

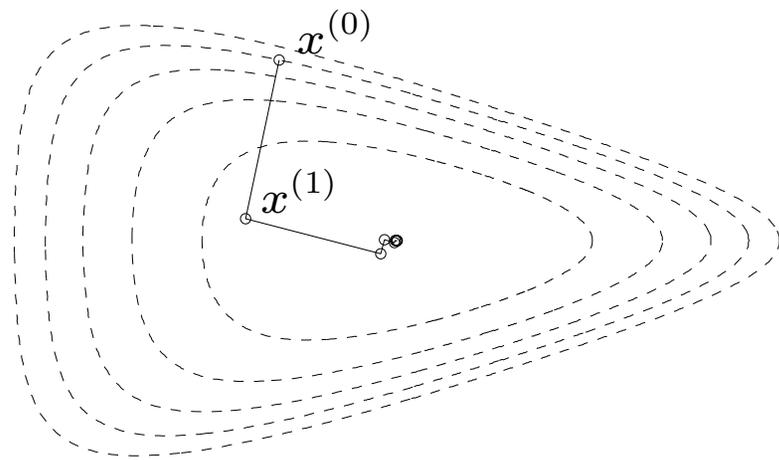


nonquadratic example

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



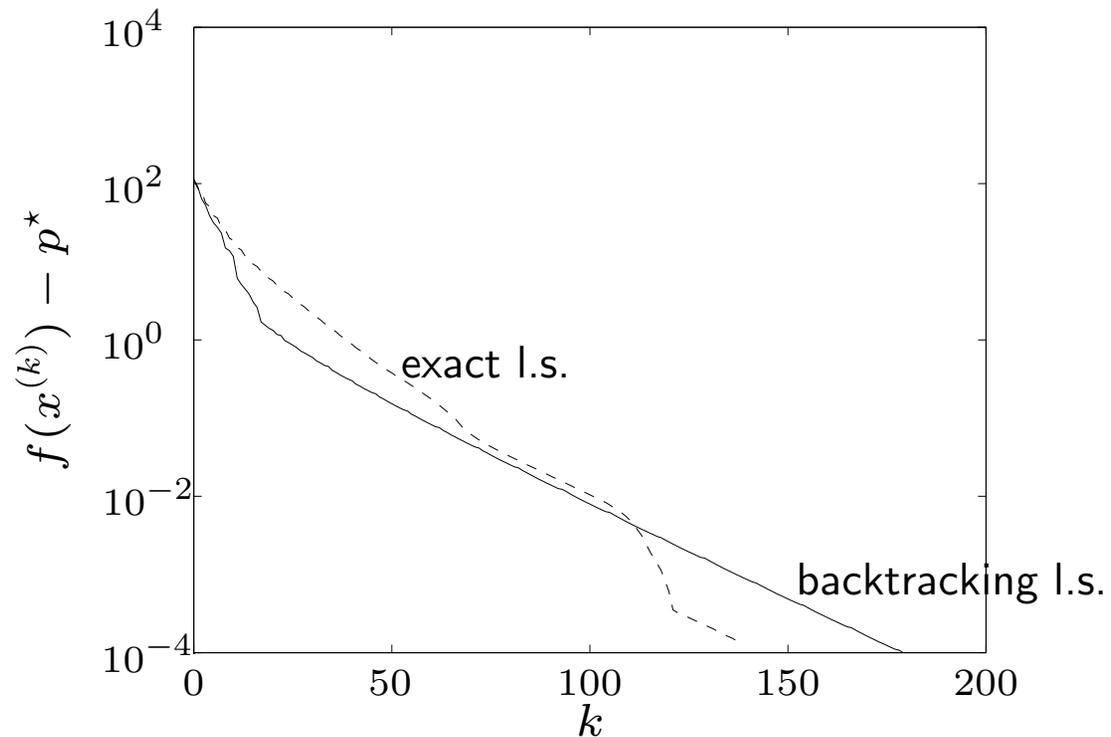
backtracking line search



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

Steepest descent method

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

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

interpretation: for small v , $f(x + v) \approx f(x) + \nabla f(x)^T v$;
direction Δx_{nsd} is unit-norm step with most negative directional derivative

(unnormalized) steepest descent direction

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

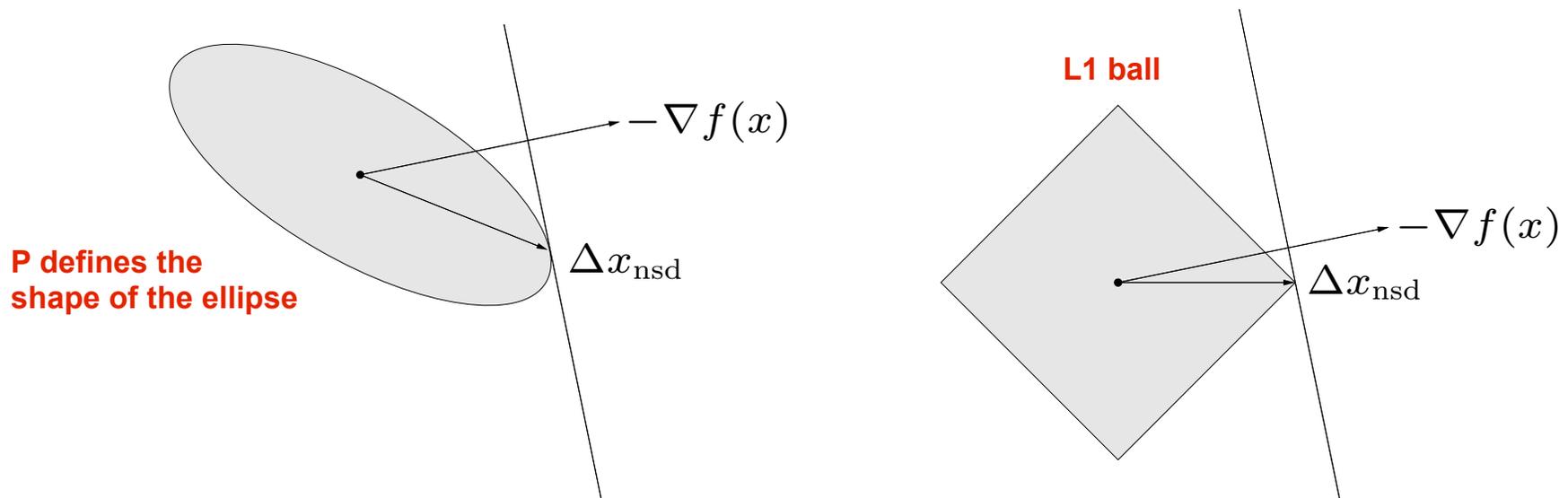
steepest descent method

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

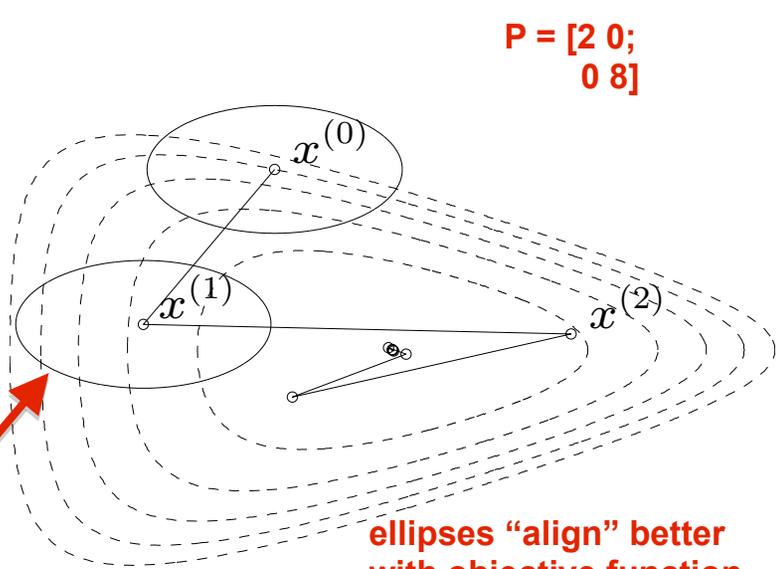
examples

- Euclidean norm: $\Delta x_{\text{sd}} = -\nabla f(x)$
- quadratic norm $\|x\|_P = (x^T P x)^{1/2}$ ($P \in \mathbf{S}_{++}^n$): $\Delta x_{\text{sd}} = -P^{-1} \nabla f(x)$
- ℓ_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 ℓ_1 -norm:

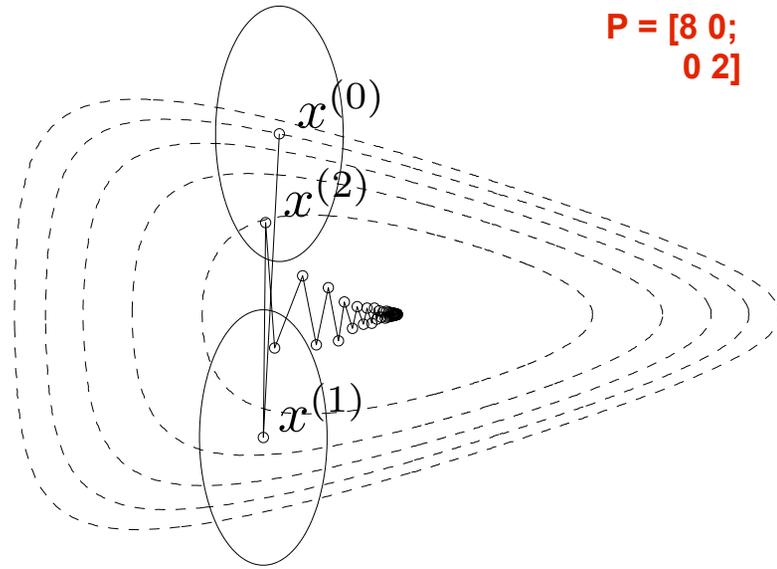


choice of norm for steepest descent



$$P = \begin{bmatrix} 2 & 0 \\ 0 & 8 \end{bmatrix}$$

ellipses "align" better with objective function thus convergence is faster



$$P = \begin{bmatrix} 8 & 0 \\ 0 & 2 \end{bmatrix}$$

- steepest descent with backtracking line search for two quadratic norms (two different P's)
- ellipses show $\{x \mid \|x - x^{(k)}\|_P = 1\}$

See Figure 9.13

- equivalent interpretation of steepest descent with quadratic norm $\|\cdot\|_P$: gradient descent after change of variables $\bar{x} = P^{1/2}x$

See Figures 9.14, 9.15

shows choice of P has strong effect on speed of convergence

Newton step

(Uses the Hessian as a good ellipse, see previous slide)

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

interpretations

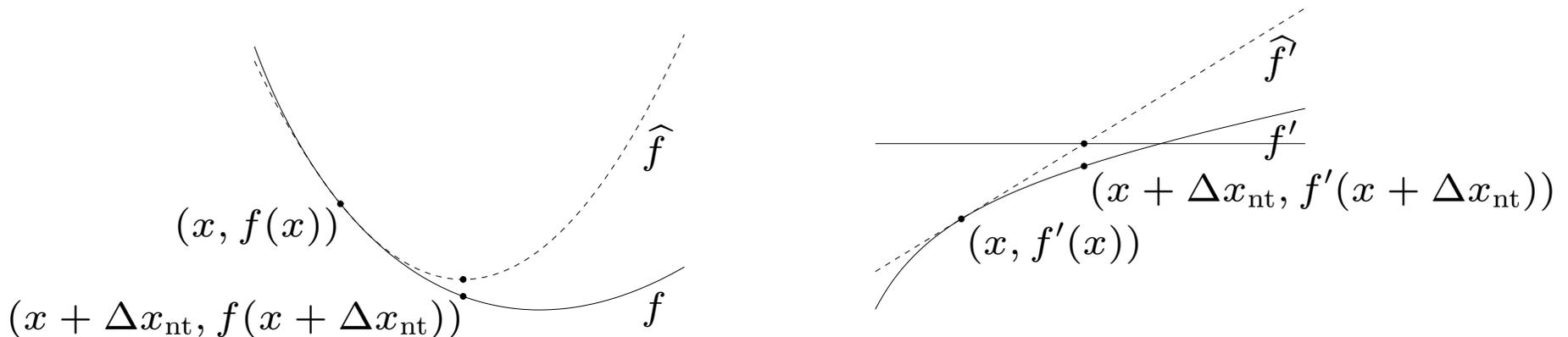
- $x + \Delta x_{\text{nt}}$ minimizes second order approximation

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

Second order
Taylor series
approximation
(we are discarding
the remainder term)

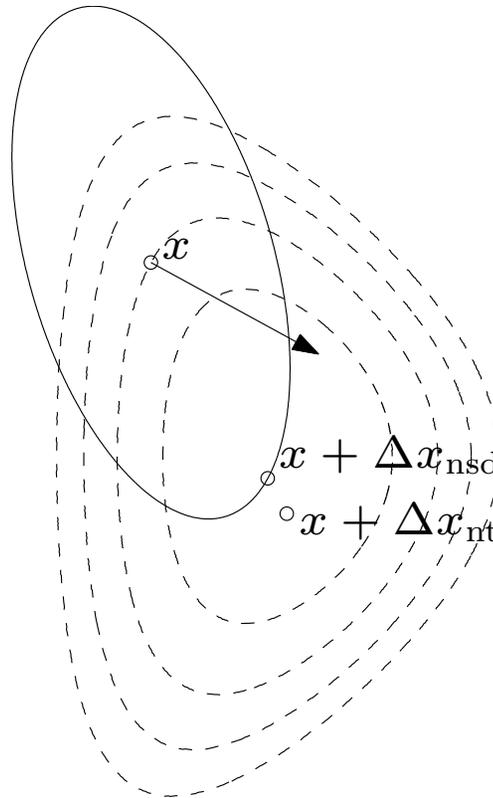
- $x + \Delta x_{\text{nt}}$ solves linearized optimality condition

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



- Δx_{nt} is steepest descent direction at x in local Hessian norm

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



Let $H = d^2 f(x)$
 $d = df(x)$

From slide 10-11 we have:

$\min d'u$
s.t. $u' H u = 1$

$L(u,v) = d'u + v (u'H u - 1)$
 $dL/du = d + 2 v H u = 0$

Then:

$u^{\wedge*} = -1/(2v) H^{-1} d$

Now, the objective function is:
 $d'u^{\wedge*} = -1/(2v) d' H^{-1} d$

If f is strongly convex, then H is positive definite, $d' H^{-1} d > 0$.
Then $v > 0$ since otherwise $d'u^{\wedge*}$ would not be minimized.

$u^{\wedge*}$ has the direction of Δx_{nt} !

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^*

properties

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


Remember $p^* = \inf_y f(y)$

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

Let $H = d^2 f(x)$

$d = df(x)$

$\lambda = \lambda(x)$

$\Delta x = \Delta x_{nt} = -H^{-1} d$

$\inf_y \hat{f}(y) = \hat{f}(x + \Delta x)$

$= f(x) + d' \Delta x + \frac{1}{2} \Delta x' H \Delta x$

$= f(x) - \frac{1}{2} d' H^{-1} d$

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

Thus $\lambda = \sqrt{d' H^{-1} d}$

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

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

$$x = T y$$

$$\text{Let } H_{\tilde{f}}(y) = d^2 f(y)$$

$$d\tilde{f}(y) = T' df(T y) = T' df(x)$$

$$H_{\tilde{f}}(y) = T' H_f(T y) T = T' H_f(x) T$$

$$y = T^{-1} x$$

$$\Delta y = -H_{\tilde{f}}(y)^{-1} d\tilde{f}(y) = -(T' H_f(x) T)^{-1} T' df(x)$$

$$= -T^{-1} H_f(x)^{-1} df(x) = T^{-1} \Delta x$$

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

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 Δx from first equation; compute w and Δx 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 $L_0^T A D^{-1} A^T L_0$)