


Unconstrained minimization

$$\min_x f(x)$$

f convex, twice continuously differentiable (hence dom f open)
assume $p^* = \inf_x f(x)$ is attained and finite

Solve optimality condition: $\nabla f(x^*) = 0$

1. initial point and sub-level set

$$x^{(0)} \in \text{dom } f$$

$S = \{x \mid f(x) \leq f(x^{(0)})\}$ is closed \Rightarrow {hard to verify, except when all sublevel sets are closed:
equivalent to the condition that: $\text{epi}(f)$ is closed

true if $\text{dom } f = \mathbb{R}^n$

true if $f(x) \rightarrow \infty$ as $x \rightarrow \text{boundary}(\text{dom } f)$

eg. $f(x) = \log \sum \exp(a_i^T x + b_i)$
 $\text{dom } f = \mathbb{R}^n$

eg. $f(x) = -\sum \log(b_i - a_i^T x)$
domain is an open polyhedron
 $b_i - a_i^T x \rightarrow 0$ as $x \rightarrow \text{boundary}$

Strong convexity and implications

f is strongly convex on S if there exists an $m > 0$ such that
 $\nabla^2 f(x) \succeq mI$ for all $x \in S$

$$f(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(z) (y-x) \quad \text{for some } z \in [x, y]$$
$$f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{m}{2} \|y-x\|^2$$

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{m}{2} \|y-x\|^2$$
$$\geq f(x) + \nabla f(x)^T (\tilde{y}-x) + \frac{m}{2} \|\tilde{y}-x\|^2$$
$$= f(x) - \frac{1}{2m} \|\nabla f(x)\|^2$$
$$\tilde{y} = \nabla f(x) + \frac{m \|y-x\|}{m} \cdot \frac{\nabla f(x)}{\|\nabla f(x)\|} = x$$

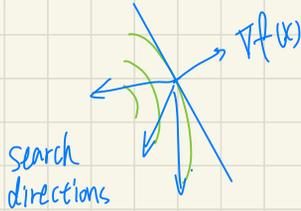
$$f(x) - p^* \leq \frac{1}{2m} \|\nabla f(x)\|^2 \quad (\text{stopping criterion})$$

• Descent methods

$$x^{(k+1)} = x^{(k)} + t^{(k)} \cdot \Delta x^{(k)}$$

with $f(x^{(k+1)}) < f(x^{(k)})$

Δx is the search direction; t is the step size from convexity. $f(x^{(k+1)}) < f(x^{(k)})$ implies that $\nabla f(x)^T \Delta x < 0$



• Line search

1. exact line search

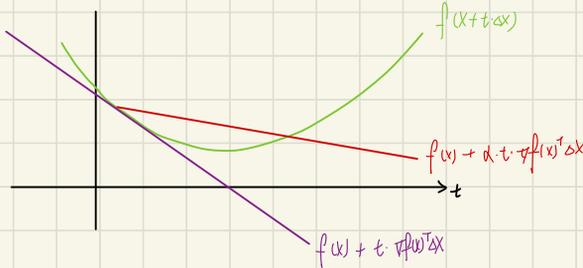
$$t = \underset{t \geq 0}{\operatorname{arg\,min}} f(x + t \cdot \Delta x)$$

2. back-tracking line search (with parameter $\alpha \in (0, 1/2)$ $\beta \in (0, 1)$)

Starting at $t=1$, repeat $t := \beta t$ until

$$f(x + t \Delta x) < f(x) + \alpha \cdot t \cdot \nabla f(x)^T \Delta x$$

$f(x + t \Delta x) < f(x) + t \cdot \nabla f(x)^T \Delta x$
will never be true by convexity



• Gradient descent method

$$\Delta x = -\nabla f(x)$$

giving $x^{(k)} \in \text{dom} f$
repeat $\{$

$$\Delta x := -\nabla f(x)$$

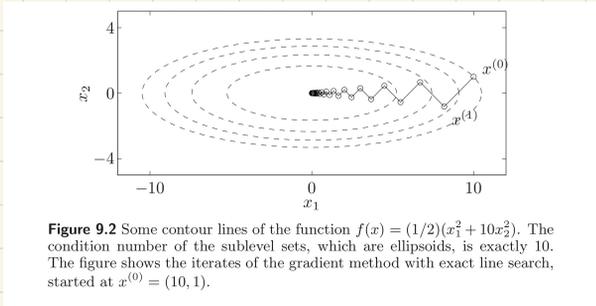
line search to choose t

$$x := x + t \Delta x$$

$\}$

Stopping criterion: $\|\nabla f(x)\|_2 \leq \epsilon$

convergence result: for strongly convex f
 $f(x^{(k)}) - p^* \leq c^k (f(x^{(0)}) - p^*)$
 $c \in (0, 1)$ depends on m , $x^{(0)}$ and line search type



• Steepest descent method

1. normalized steepest descent direction

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

directional derivative along $x+tv$

$$f(x+tv) \approx f(x) + \nabla f(x)^T V \text{ for small } t$$

Δx_{nsd} is unit-norm step with most negative directional derivative

2. unnormalized steepest descent direction

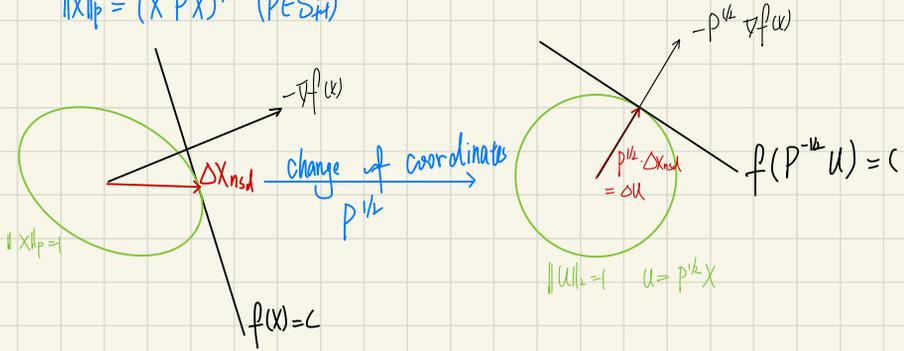
$$\Delta x_{\text{sd}} = \frac{\nabla f(x)}{\|\nabla f(x)\|_2} \cdot \Delta x_{\text{nsd}} \quad \|\nabla f(x)\|_2 = \sup_{\|z\|=1} \{ z^T \nabla f(x) \}$$

$$\nabla f(x)^T \Delta x_{\text{sd}} = -\|\nabla f(x)\|_2$$

Gradient descent is steepest descent in euclidean norm

eg. Quadratic norm

$$\|x\|_p = (x^T P x)^{1/2} \quad (P \text{ PSD})$$



$$\Delta u = -P^{1/2} \cdot \nabla f(P^{-1/2} u) = -P^{1/2} \cdot \nabla f(x)$$

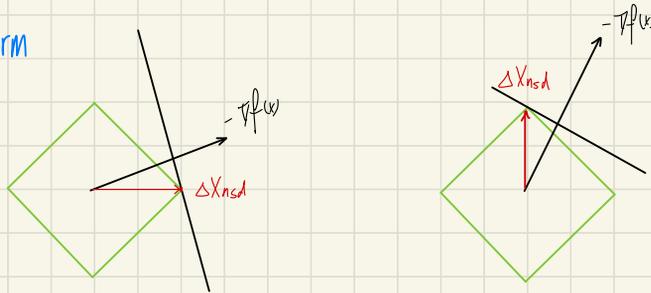
$$\begin{aligned} \Delta u_{nsd} &= -P^{1/2} \cdot \nabla f(x) / \|-P^{1/2} \cdot \nabla f(x)\|_2 \\ &= -P^{1/2} \cdot \nabla f(x) / (\nabla f(x)^T P \nabla f(x))^{1/2} \end{aligned}$$

$$\begin{aligned} \Delta x_{nsd} &= P^{1/2} \cdot \Delta u_{nsd} \\ &= -P^{-1} \cdot \nabla f(x) / (\nabla f(x)^T P \nabla f(x))^{1/2} \end{aligned}$$

$$\Delta u_{sd} = -P^{1/2} \cdot \nabla f(x)$$

$$\Delta x_{sd} = -P \cdot \nabla f(x)$$

eg. L1 norm



$$\Delta x_{nsd} = -\text{sign}\left(\frac{df}{dx_i}\right) \cdot e_i \quad i = \text{argmax}_i \left| \frac{df}{dx_i} \right|$$

Steepest descent in L1: update 1 component each step

3. Choice of norm for steepest descent

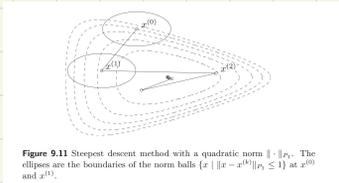


Figure 9.11 Steepest descent method with a quadratic norm $\| \cdot \|_2$. The ellipses are the boundaries of the norm balls $\{x \mid \|x - x^{(0)}\|_2 \leq 1\}$ at $x^{(0)}$ and $x^{(1)}$.

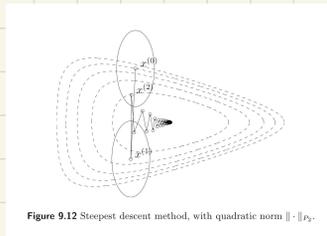


Figure 9.12 Steepest descent method, with quadratic norm $\| \cdot \|_2$.

want the norm to be consistent with the geometry of the sublevel sets

$$f(x) \approx f(x^*) + \underbrace{\nabla f(x^*)^T}_{0} (x-x^*) + (x-x^*)^T \nabla^2 f(x^*) (x-x^*)$$

sublevel sets are ellipsoids (f is nearly quadratic) near the optimal point

Steepest descent in the norm induced by the hessian \Rightarrow Newton's method

• Newton step

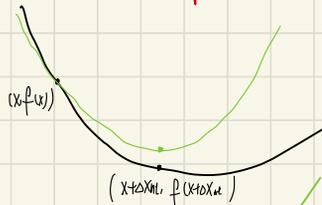
$$\Delta x_{nt} = -\nabla^2 f(x)^{-1} \nabla f(x)$$

1. interpretation 1.

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

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

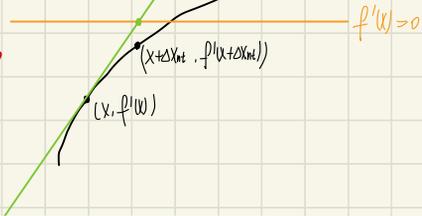
$$v = -\nabla^2 f(x)^{-1} \nabla f(x)$$



2. interpretation 2

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

$$\nabla f(x+v) \approx \nabla f(x) + \nabla^2 f(x) \cdot v := 0$$



3. interpretation 3.

Steepest descent in local hessian norm

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

4. Affine invariance of Newton step.

Suppose $\tilde{f}(y) = f(Ty)$ $T \in S^m$

$$\nabla \tilde{f}(y) = T^T \nabla f(y) \quad \nabla^2 \tilde{f}(y) = T^T \nabla^2 f(y) T$$

$$\Delta y_{nt} = -(T^T \nabla^2 f(y) T)^{-1} \cdot T^T \nabla f(y)$$

$$= -T^{-1} \nabla^2 f(y)^{-1} \nabla f(y)$$

$$= T^{-1} \Delta x_{nt}$$

$x + \Delta x_{nt} = T(y + \Delta y_{nt})$ Newton step is independent of choice of coordinates

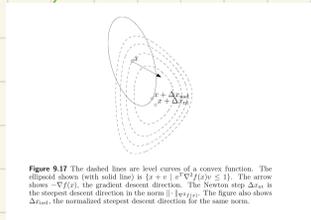


Figure 9.17 The dashed lines are level curves of a convex function. The ellipsoid shown (with solid line) is $\{x + v \mid v^T \nabla^2 f(x) v \leq 1\}$. The arrow shows $-\nabla f(x)$, the gradient descent direction. The Newton step Δx_{nt} is the steepest descent direction in the norm $\|\cdot\|_{\nabla^2 f(x)}$. The figure also shows Δx_{grad} , the normalized steepest descent direction for the same norm.

$$\nabla \tilde{f}(y) = T^T \nabla f(x)$$

gradient method needs prior covariance

5. Newton decrement

$$\lambda(x) = (\nabla^2 f(x)^T \nabla^2 f(x)^{-1} \nabla^2 f(x))^{1/2} \quad \text{the hessian-norm of gradient}$$

λ^2 gives a estimate of decrease

$$\begin{aligned} f(x) - \inf_y f(y) &= f(x) - f(x + \Delta x_{nc}) \\ &= -\nabla f(x)^T \Delta x_{nc} - \frac{1}{2} \Delta x_{nc}^T \nabla^2 f(x) \Delta x_{nc} \\ &= \frac{1}{2} \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) = \frac{1}{2} \lambda(x)^2 \end{aligned}$$

$$\begin{aligned} \Delta x_{nc} &= -\nabla^2 f(x)^{-1} \nabla f(x) \\ &= \frac{1}{2} \Delta x_{nc}^T \nabla^2 f(x) \Delta x_{nc} \end{aligned}$$

Newton's method

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

compute the newton step and decrement

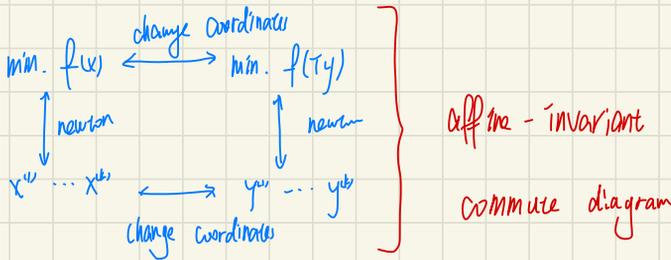
$$\begin{aligned} \Delta x_{nc} &= -\nabla^2 f(x)^{-1} \nabla f(x) \\ \lambda^2 &= \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) \end{aligned}$$

stop if $\lambda^2 \leq \epsilon$

line search to choose step size t

$$x := x + t \Delta x_{nc}$$

$\}$



Newton method convergence analysis

Assume f is strongly convex on S with constant $m > 0$ s.t. $\nabla^2 f(x) \succeq mL$.

$\nabla^2 f$ is Lipschitz continuous on S , with constant $L > 0$

$$\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq L \|x - y\| \quad (\text{affine transformation changes norm!!!})$$

L is a constraint on third derivative; L measures how well f is approximated by a quadratic function

(L can be taken 0 for a quadratic function; newton method works perfectly (1 step) for quadratic function)

(when the 3rd derivative is small, i.e. the hessian is slowly changing, newton method works well)

there exist constants $\eta \in (0, m/L)$, $\gamma > 0$ such that
 if $\|\nabla f(x^k)\|_2 \geq \eta$, then $f(x^{k+1}) - f(x^{k+1}) \geq \gamma$ when gradient large, ^{certain constant} guaranteed to decrease by a
 if $\|\nabla f(x^k)\|_2 < \eta$ then

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

1. damped Newton phase ($\|\nabla f(x^k)\|_2 \geq \eta$)
 most iterations require backtracking steps
 function value decreases by at least γ each step
 if $\gamma^* > -\infty$, the phase takes at most $(f(x^{(0)}) - \gamma^*)/\gamma$ iterations

2. quadratically convergent phase ($\|\nabla f(x^k)\|_2 < \eta$)
 all iterations use step size $t=1$
 $\|\nabla f(x^k)\|_2$ converges to zero quadratically; if $\|\nabla f(x^{k_0})\|_2 < \eta$, then

$$\frac{L}{2m^2} \|\nabla f(x^k)\|_2 \leq \left(\frac{L}{2m^2} \|\nabla f(x^{k_0})\|_2\right)^{2^{k-k_0}} \leq \underbrace{\left(\frac{1}{2}\right)^{2^{k-k_0}}}_{(L \geq k)}$$

number of accurate digits doubles every step.

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

$$\underbrace{\frac{f(x^{(0)}) - p^*}{\gamma}}_{\text{first stage:}} + \log_2 \log_2 (t_0/\epsilon)$$

first stage:
 depends on how suboptimal $x^{(0)}$ is

eg. \mathbb{R}^2

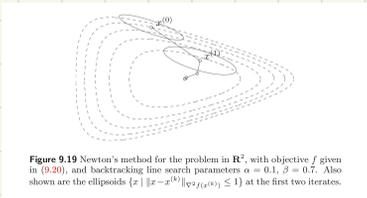


Figure 9.19 Newton's method for the problem in \mathbb{R}^2 , with objective f given in (9.20), and backtracking line search parameters $\alpha = 0.1$, $\beta = 0.7$. Also shown are the ellipsoids $\{x \mid \|x - x^{(k)}\|_{\text{Hess} f(x^{(k)})} \leq 1\}$ at the first two iterates.

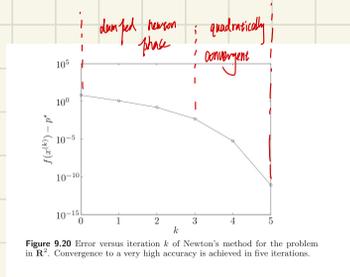


Figure 9.20 Error versus iteration k of Newton's method for the problem in \mathbb{R}^2 . Convergence to a very high accuracy is achieved in five iterations.

each newton step computes $Hv = -g$

