

Lecture 4

Last week: algorithms for solving 1-d minimization problems e.g. Newton's Method
algorithms for multi-d minimization

Main example for today

Steepest descent:



x_{k+1} minimum point of f in the half-line $\{x_k - s d_k : s \geq 0\}$,
 $d_k = \nabla f(x_k)$

We'll write $x_{k+1} \in A(x_k)$ when \textcircled{X} holds

Goal: ① Use general theory to show "algorithm works"
 ② Analyze its convergence - When does it work well/badly?

Theorem says: $f: E^n \rightarrow \mathbb{R}$ continuous
 If A is a point-to-set mapping $E^n \rightarrow E^n$ and

① A is closed at x if x is not global minimizer

② If $y \in A(x)$ and X is not global minimum then $f(y) < f(x)$.

and if $x_{k+1} \in A(x_k)$ and (x_k) is bounded then any limit of a convergent subsequence is a minimizer

background

To use theorem, must check hypothesis ① and ②.

Let's check for method of steepest descent.

Assume f is C^1 and convex.

Check hypothesis ①

Have to check: If $\begin{cases} x_k \rightarrow x \\ y_k \rightarrow A(x_k) \\ y_k \rightarrow y \end{cases}$ then $y \in H(x)$

True because

$$y_k = x_k - s_k d_k, d_k = \nabla f(x_k)$$

$$s_k \text{ minimizes } g_k(s) = f(x_k - s d_k)$$

Need to show $y = x - s^* d$, $d = \nabla f(x)$
 s^* minimizes $g(s) = f(x - s d)$

Claim 1: $d_k \rightarrow d$ easy because $d_k = \nabla f(x_k)$
 $d = \nabla f(x)$
 $x_k \rightarrow x$ and ∇f is continuous

Claim 2: $s_k \rightarrow s^*$ with $y = x - s^* d$

$$\text{True because } s_k = \frac{|y_k - x_k|}{|d_k|} \rightarrow \frac{|y - x|}{d} = s^*$$

Since $|d| \neq 0$ (since $\nabla f(x) = 0$ only for global minimizers by convexity)

Claim 3: $y = x - s^* d$

True because $y_k = x_k - s_k d_k$

while $y_k \rightarrow y, x_k \rightarrow x, s_k \rightarrow s^*, d_k \rightarrow d$
 everything converges

Claim 4: s^* minimizes $g(s) = f(x - sd)$

$$g(s^*) = f(x - s^*d) = \lim_k f(x_k - s_k d_k) \leq \lim_k f(x_k - s d_k) \text{ for every } s \\ = f(x - sd) = g(s)$$

so s^* minimizes g .

So far we proved A closed.

Also need: $y \in A(\infty) \Rightarrow f(y) < f(x)$ (if x is not global min)

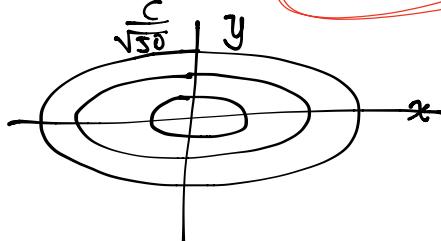
Note: in claim 2 above, definition of s^* was $s^* = \frac{|y-x|}{\|d\|} = \lim_k s_k$

To show $f(y) < f(x)$, enough to check that $f(x - sd) < f(x)$ for some $s > 0$.

~~True by Taylor Thm~~

~~$f(x - sd) = f(x) - sd^T d + ds$ → negligible.
for small s .~~

Consider $f(x, y) = x^2/2 + 50y^2$



level curves: $f(x, y) = c^2$
i.e. $\frac{x^2}{2c^2} + \frac{y^2}{c^2/50} = 1$

very bad convergence

see below

Initial guess: $(x_0, y_0) = (100, 1)$

$$\nabla f(x_0, y_0) = [x_0, 100y_0] = [100, 100]$$

$(x_1, y_1) = (x_0, y_0) - S(100, 100) = (100(1-S), (1-100S))$ where S minimizes: $g(s) = 100^2(1-s)^2/2 + 50(1-100s)^2$

Answer: $S = 2/101$

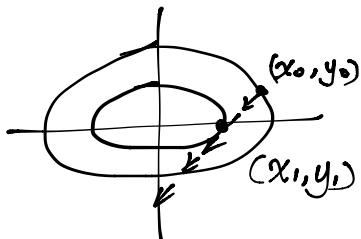
$$(x_1, y_1) = (100 - 200/101, 1 - 200/101) = (98 \cdot 2/101, -99/101) = \frac{99}{101}(100, -1)$$

Repeat to find: $(x_k, y_k) = (\frac{99}{101})^k (100, 1) \dots k \text{ even}$

$$(\frac{99}{101})^k (100, -1) \dots k \text{ odd}$$

e.g. $(x_{100}, y_{100}) \approx (3.5, -135)$ still far from 0.

How to understand?



Note: $d_0 = \nabla f(x_0, y_0)$ tangent to level curve of f at (x_0, y_0)

$d_1 = \nabla f(x_1, y_1)$ is orthogonal to the level curve of f at (x_0, y_0)

so for steepest descent in 2d, each direction is orthogonal to previous direction.

Facts about convergence

① Suppose f is quadratic: $f(x) = \frac{1}{2}x^T Q x - b^T x + c$
 Then: $|f(x_{k+1}) - f_{\min}| \leq (\frac{1-r}{1+r})^2 (f(x_k) - f_{\min})$

$$r = \frac{\lambda_{\max}}{\lambda_{\min}}$$

λ_{\min} = smallest eigenvalue of Q
 λ_{\max} = largest eigenvalue of Q

Note $0 < r < 1$: $r \geq 1$ at least!

i.e. good convergence if $r=1$

bad convergence if $r=0$

$$\text{in our example, } r=100, \frac{1-r}{1+r} = \frac{99}{101}$$

$$Q = \begin{pmatrix} 1 & 0 \\ 0 & 100 \end{pmatrix} \uparrow$$

Similarly: for more general convex functions, things are similar.

We'll prove these facts later. First some applications

① Suppose I want to minimize convex function f .

2 possible procedures:

a. minimize w/ steepest descent or

b. minimize $F(x) = |\nabla f(x)|^2$ with steepest descent.

(since x^* minimizer $\Leftrightarrow \nabla f(x) = 0 \Leftrightarrow x^*$ minimizes f .)

Which is better?

For simplicity assume $f(x) = \frac{1}{2}x^T Q x - b^T x$

$$\nabla f(x) = x^T Q - b^T$$

$$F(x) = |\nabla f(x)|^2 = (x^T Q - b^T)(Qx - b) = x^T Q^2 x - 2b^T Q x + b^T b$$

Which is easier to minimize via steepest descent?

i.e. which one has bigger condition number

$$\begin{matrix} Q & \text{vs} & Q^2 \\ \text{for } f & & \text{for } |\nabla f|^2 \end{matrix}$$

Recall: eigenvalues for Q^2 = eigenvalues for $Q \cdot Q^T$

so $\lambda_{Q^2} = (\lambda_Q)^2$ where λ_S = cond. # for S .

$$\text{so } \lambda_Q \geq \lambda_{Q^2}$$

Conclusion: better to stick with minimizing f .

Example: Suppose I want to solve:

minimize $f(x)$ with constraint $h(x)=0$

One approach: minimize $f(x) + \frac{1}{2}\mu h(x)^2$

for μ very large, then should happen that h close to 0 at minimum point (so constraint nearly satisfied)

Can we do this well with steepest descent?

For simplicity, assume

$$\begin{aligned} f(x) &= \frac{1}{2}x^T Q x + b^T x && \text{(quadratic)} \\ h(x) &= C^T x && \text{(linear)} \\ h(x)^2 &= x^T C C^T x \end{aligned}$$

$$\text{so } f(x) + \frac{1}{2}\mu h(x)^2 = \frac{1}{2}x^T(Q + \mu C C^T)x - b^T x$$

Now huge what is condition number at $Q + \mu C C^T$?

We guess: condition number is tiny if μ huge.

Conclusion is: steepest descent probably will not work well.

Next convergence analysis

But first, why do we care?

i.e. to minimize $f(x) = \frac{1}{2}x^T Qx - b^T x$
why don't I just find minimizer?

That is: $\nabla f = x^T Q - b^T$

So $\nabla f(x^*) = 0$ when $(x^*)^T Q = b^T$ i.e. $x^* = Q^{-1}b$

Why bother?

- ① To gain insight
- ② Hard to compute Q^{-1}

Now: convergence analysis for steepest descent, when $f(x) = \frac{1}{2}x^T Qx - b^T x$

① Given x find formula for $y = A(x)$ = minimizer of f on half-line $\{x - s d : s > 0\}$
 $d = \nabla f(x)$

But first simplify: since $x^* = Q^{-1}b$ i.e. $b = Qx^*$

$$f(x) = \frac{1}{2}x^T Qx - (x^*)^T Qx = \underbrace{\frac{1}{2}(x-x^*)^T Q(x-x^*)}_{\text{call this } E(x)} - \underbrace{\frac{1}{2}x^* Qx^*}_{\text{a number of independent of } x}$$

So minimizing f = minimizing E

① Formula for $y = A(x)$ $d = (\nabla E)^T = Q(x-x^*)$

$y = x - s \cdot d$ for optimal choice of s .

$$\begin{aligned} E(x - s \cdot d) &= \frac{1}{2}(x-x^*-s \cdot d)^T Q(x-x^*-s \cdot d) = \frac{s^2}{2}(\quad) + s(\quad) \\ &= \frac{s^2}{2}(d^T Q d) - s(x-x^*)^T Q d + (\text{independent of } s) \\ &= g(s) \end{aligned}$$

$$g'(s) = s \cdot d^T Q d - (x-x^*)^T Q d$$

$$g'(s) = 0 \text{ if } s = \frac{(x-x^*)^T Q d}{d^T Q d} = \boxed{\frac{d^T d}{d^T Q d}}$$

We have used: $\nabla E(x) = \nabla f(x) = x^T Q - b^T = (x-x^*)^T Q$

So $d = (\nabla E)^T = Q(x-x^*)$

Thus $y = x - s^* d = x - \left(\frac{d^T d}{d^T Q d}\right) d \quad \leftarrow \text{This is } A(x)$

$$\begin{aligned} d &= \nabla E(x)^T \\ &= \nabla f(x)^T \\ &= \alpha x - b \end{aligned}$$

② Formula for $E(y)$ compared to $E(x)$

$y = A(x)$

Idea: just plug into definition of E & compute:

This leads to:

$$E(y) = \underbrace{\left(1 - \frac{(d^T d)^2}{(d^T Q d)(d^T Q^{-1} d)}\right)}_{\text{happy if this is small } \ll 1} E(x)$$

happy if this is small $\ll 1$

$$\begin{aligned} ③ \text{ I claim that } 1 - \frac{(d^T d)^2}{(d^T Q d)(d^T Q^{-1} d)} &\leq \left(\frac{1-r}{1+r}\right)^2 = \left(\frac{A-a}{A+a}\right)^2 \\ &= 1 - \frac{4ra}{(A+a)^2} \end{aligned}$$

$r = \frac{a}{A}$ = condition number

I need to show that

$$\frac{(d^T d)^2}{(d^T Q d)(d^T Q^{-1} d)} \geq \frac{4\alpha A}{(A + \alpha)^2}$$

a) divide numerator and denominator by $(d^T d)^2 = \|d\|^4$

$$\text{Let } y = \frac{d}{\|d\|} = \text{unit vector}$$

Then left-hand side is

$$\frac{1}{(y^T Q y)(y^T Q^{-1} y)}$$

Let w_1, \dots, w_n orthonormal basis of eigenvectors for Q .

Any unit vector y can be written $y = a_1 w_1 + \dots + a_n w_n$ where $1 = \|y\|^2 = y^T y = a_1^2 + \dots + a_n^2 = \theta_1 + \dots + \theta_n$
 $\theta_i = a_i^2$

Then substitute & expand $Qw = \lambda w$

$$\frac{1}{(y^T Q y)(y^T Q^{-1} y)} = \frac{1}{(\sum_{i=1}^n \theta_i \lambda_i)(\sum_{i=1}^n \theta_i / \lambda_i)} = \frac{\psi(\sum \theta_i \lambda_i)}{\sum \theta_i \psi(\lambda_i)} \quad \psi(s) = \frac{1}{s}$$

The point is: geometric reasoning tells us that

smallest value of $\frac{1}{(y^T Q y)(y^T Q^{-1} y)}$ for $\|y\|=1$ occurs when

$\theta_1, \dots, \theta_n$ nonzero

$\theta_2, \dots, \theta_{n-1} = 0$ i.e. only smallest + largest eigenvalues "archive"

if $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n = A$

i.e. minimum of $\frac{(d^T d)^2}{(d^T Q d)(d^T Q^{-1} d)}$ occurs for θ minimizing

$$\frac{1}{(\theta a + (1-\theta)A)(\frac{\theta}{a} + \frac{1-\theta}{A})} \quad 0 \leq \theta \leq 1$$

In fact, this is minimized if $\theta = \frac{1}{2}$

This finally leads to our formula.

$$E(y) \leq \left(\frac{1-r}{1+r}\right)^2 E(x) \quad \text{if } y = Ax$$

One other fact: Suppose f is not quadratic, but $[a] \leq \nabla^2 f(x) \leq A[I]$ for all x

Where $Q \leq R$ means $v^T Q v \leq v^T R v$ for all v

So $Q \leq A[I]$

means $v^T Q v \leq v^T (A[I]) v = A(v)^2$

(Then $\frac{\alpha}{A}$ is analog of condition number)

When this holds, $[f(x_{k+1}) - f_{\min}] \leq (1 - \frac{\alpha}{A})[f(x_k) - f_{\min}] = (1-r)[f(x_k) - f_{\min}]$
 for method of steepest descent.