

GDI SUMMARY of GRADIENT DESCENT COMPLEXITY

Note the difference!

- Strict Convexity $f(\theta x + (1-\theta)y) < \theta f(x) + (1-\theta)f(y)$
- Strong Convexity $\exists m > 0$ s.t. $\nabla^2 f(x) \succeq mI$ for all $x \in S$.

e^{-x} is strictly convex but not strongly convex.

Big Picture Results for Gradient Descent from last lecture, assuming f strongly convex.

Exact line search at least as good as $t = \frac{1}{M}$,

which gives

$$f(x^{(k)}) - p^* \leq \left(1 - \frac{1}{K}\right)^k (f(x^{(0)}) - p^*)$$

where

$$K = M/m$$

$$M = \sup_{x \in S} \lambda_{\max}(\nabla^2 f(x))$$

$$m = \inf_{x \in S} \lambda_{\min}(\nabla^2 f(x))$$

BUT don't know M, m in practice.

Hence: BACKTRACKING LINE SEARCH

$\alpha \in (0, 1/2)$ ARMISTO COND.

$\beta \in (\alpha, 1)$ BACKTRACK PARAM.

When $\beta = 1/2$, $M \geq 1/2$, set

$$f(x^{(k)}) - p^* \leq \left(1 - \frac{\alpha}{K}\right)^k (f(x^{(0)}) - p^*)$$

⇒ # iters for accy ϵ is $O\left(\frac{1}{\log \epsilon}\right)$

Gradient Descent, cont'd.

Nesterov (Thm 2.1.15) also gives another complexity result for

$$t = \frac{2}{m+M}$$

namely

$$\|x^{(k)} - x^*\| \leq \left(\frac{1 - 1/K}{1 + 1/K} \right)^k \|x^{(0)} - x^*\|$$

which leads to

$$f^{(k)} - p^* \leq K \left(\frac{1 - 1/K}{1 + 1/K} \right)^{2k} (f^{(0)} - p^*)$$

where x^* = minimizer, $p^* = f(x^*)$.

Both BV and Nesterov also give results for the non-(strongly convex) case - these are much weaker.

Newton's Method.

$$\Delta x \text{ is solution of } \nabla^2 f(x^{(k)}) \Delta x = -\nabla f(x^{(k)}).$$

Motivation: minimize "quadratic model"

$$\varphi(v) = \underbrace{\nabla f(x^{(k)})^T}_{[f(x^{(k)})']} v + \frac{1}{2} v^T \nabla^2 f(x^{(k)}) v.$$

To solve equation, use CHOLESKY FACTORIZATION

$$\nabla^2 f(x^{(k)}) = \begin{matrix} L & L^T \\ \Delta & \Delta \end{matrix}$$

$$\underbrace{L L^T}_{y} \Delta x = \Delta x = -g.$$

- 1) forward solve for y
- 2) back solve for Δx .

Cost: $\frac{1}{6} n^3$ adds + mults.

Use same backtracking line search.

(Newton used this for finding zeros of polynomials, not minimization; particularly root of $P(\lambda) = \lambda^2 - a$ i.e. square roots).

Convergence Analysis of Newton's Method.

As before, suppose that $S = \{x: f(x) \leq f(x_0)\}$ is compact and $M I \geq \nabla^2 f(x) \geq m I$ on S , $m > 0$.
Now also need

$$\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq L \|x - y\| \quad \forall x, y \in S$$
 i.e. $\nabla^2 f$ is Lipschitz.

Turns out (BV §9.5) that $\exists \eta > 0, \delta > 0$ s.t.

If $\|\nabla f(x^{(k)})\| \geq \eta$, the backtracking line search returns $t_k \geq \frac{\beta m}{M}$ with

$$f(x^{(k+1)}) - f(x^{(k)}) \leq -\delta \quad (*)$$

while if $\|\nabla f(x^{(k)})\| < \eta$, B.T.L.S. returns $t_k = 1$ with

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

NOTE

"QUADRATIC CONVERGENCE".

Consequences If $\eta \leq \frac{m^2}{L}$ and, for some k , $\|\nabla f(x^{(k)})\| \leq \eta$

then

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

so this applies recursively and hence (*) holds for all $l \geq k$, so

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

| | | | | |
|-------|---------------|---------------|----------------|-----------------|
| $l=k$ | $l=k+1$ | $l=k+2$ | $l=k+3$ | quad. conv. |
| RHS | $\frac{1}{2}$ | $\frac{1}{4}$ | $\frac{1}{16}$ | $\frac{1}{256}$ |

But what are β, η, γ ? Turns out (BV p. 489-491) that ~~if~~ using BTLN with Newton step, we always have $t_k \geq \frac{\beta m}{M}$ and that consequently

$$f(x^{(k+1)}) - f(x^{(k)}) \leq -\alpha \beta \eta^2 \frac{m}{M^2} \underbrace{\hspace{2cm}}_{\gamma}$$

so (†) holds if we set ~~the~~ η

It also turns out that if $\eta \leq 3(1-2\alpha) \frac{m}{L}$

then $t_k = 1$, i.e. $f(x^{(k)} + \underbrace{\Delta x_{NT}}_{\text{Newton step}})$ satisfies the

"sufficient decrease" condition in the BTLN.

We will now show that (*) holds as a consequence.
(QUAD. CONTR.)

Proof of (*) (quadratic contraction) assuming $t_k = 1$.

$$\|\nabla f(x^{(k)} + \Delta x_{NT})\| = \underbrace{\|\nabla f(x^{(k)} + \Delta x_{NT}) - \nabla f(x^{(k)}) - \nabla^2 f(x^{(k)}) \Delta x_{NT}\|}_{\text{ZERO BY DEF.}}$$

$$\int_0^1 \nabla^2 f(x^{(k)} + s \Delta x_{NT}) \Delta x_{NT} ds$$

find the calc.

$$= \left\| \int_0^1 (\nabla^2 f(x^{(k)} + s \Delta x_{NT}) - \nabla^2 f(x^{(k)})) \Delta x_{NT} ds \right\|$$

$$\leq \int_0^1 L \|\Delta x_{NT}\| ds$$

$$= \frac{L}{2} \|\Delta x_{NT}\|^2$$

$$= \frac{L}{2} \|(\nabla^2 f(x^{(k)}))^{-1} \nabla f(x^{(k)})\|^2$$

$$\leq \frac{L}{2m^2} \|\nabla f(x^{(k)})\|^2 \equiv (*)$$

Note: for this to apply recursively, also need $\eta \leq \frac{m^2}{L}$ as explained on NM2 (bottom). So, need

$$\eta = \min(1, 3(1-2d)) \frac{m^2}{L}$$

Total # iterations

Initial Phase with $\|\nabla f(x)\| \geq \eta$

$$\# \text{ steps} \leq \frac{f(x_0^{(l)}) - p^*}{\gamma} \text{ immediately from (1)}$$

Quadratically convergent phase with $\|\nabla f(x)\| < \eta$:

$$\begin{aligned} f(x_0^{(l)}) - p^* &\leq \frac{1}{2m} \|\nabla f(x_0^{(l)})\|^2 && \leftarrow \text{(from 1')} \\ &\leq \frac{1}{2m} \frac{4m^4}{L^2} \left(\frac{1}{2}\right)^{2^{l-k+1}} && \text{(square both sides of (7)) (BV p457, line 5)} \\ &= \frac{2m^3}{L^2} \left(\frac{1}{2}\right)^{2^{l-k+1}} \end{aligned}$$

If want RHS $\leq \epsilon$, or LHS $\leq \epsilon$, need

$$\left(\frac{1}{2}\right)^{2^{l-k+1}} \leq \frac{\epsilon L^2}{2m^3}$$

$$2^{2^{l-k+1}} \geq \frac{\epsilon_0}{\epsilon} \quad \text{where } \epsilon_0 = \frac{2m^3}{L^2}$$

$$2^{l-k+1} \geq \log_2 \frac{\epsilon_0}{\epsilon}$$

$$\# \text{ steps} = l - k + 1 \geq \log_2 \log_2 \frac{\epsilon_0}{\epsilon}$$

$$\text{e.g. } \frac{\epsilon_0}{\epsilon} = 10^{25}$$

$$\log_2 10^{25} \approx \log_2 2^{50} = 50$$

$$\log_2 50 \leq 6$$

ACCURATE DIGITS
DOUBLES EVERY
ITERATION.

Very few steps once quadratic
convergence starts.

Lower Complexity Bounds (Nesterov Sec 2.14)

Assume as before that f is strongly convex & C^2 with

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

↙ μ in Nesterov

↖ L in Nesterov.

Assume that at each point $x^{(k)}$, a "first-order oracle" or "black box" computes $f(x^{(k)})$ and $\nabla f(x^{(k)})$.

Assume also that for $k=1, 2, \dots$

(ASSUMPTION) $x_k \in \text{Linear span} \{ \nabla f(x_0), \dots, \nabla f(x_{k-1}) \}$.

A

For simplicity, assume $\text{dom } f = \mathbb{R}^\infty \equiv \ell_2 =$

$$\left\{ x = (x_i)_{i=1}^\infty : \|x\|^2 = \sum_{i=1}^\infty x_i^2 < \infty \right\}.$$

Now we define a "difficult" function F by

$$F(x) = \frac{M-m}{8} \left\{ (x_1)^2 + \sum_{i=1}^\infty (x_i - x_{i+1})^2 - 2x_1 \right\} + \frac{m}{2} \|x\|^2.$$

We have

$$\frac{\partial F}{\partial x_1} = \frac{M-m}{8} (2x_1 + 2(x_1 - x_2) - 2) + mx_1$$

$$\begin{aligned} j \geq 1: \quad \frac{\partial F}{\partial x_j} &= \frac{M-m}{8} (2(x_j - x_{j+1}) - 2(x_{j-1} - x_j)) + mx_j \\ &= \frac{M-m}{8} (4x_j - 2x_{j+1} - 2x_{j-1}) + mx_j \end{aligned}$$

$$\text{so } \nabla^2 F(x) = \frac{M-m}{4} \begin{bmatrix} 2 & -1 \\ -1 & 2 \\ & & \ddots \\ & & & 2 & -1 \\ & & & -1 & 2 \end{bmatrix} + mI$$

with $\nabla^2 F(x) \leq 4I, \geq 0$

$$\left. \begin{aligned} \lambda_{\max}(\nabla^2 F(x)) &\leq M \\ \lambda_{\min}(\nabla^2 F(x)) &\geq m \end{aligned} \right\} \text{as required.}$$

and

$$\nabla F(x) = \left(\frac{M-m}{4} T + mI \right) x - \frac{M-m}{4} e_1 \quad \text{with } \begin{bmatrix} 1 \\ 0 \\ \vdots \end{bmatrix}$$

The solution x^* is given by $\nabla F(x^*) = 0$:

$$\frac{M-m}{4} (2x_1 - x_2) + m x_1 = \frac{M-m}{4}$$

$$2x_1 - x_2 + \frac{4}{M-m} m x_1 = 1$$

$$x_2 - \frac{(2(M-m) + 4m)x_1 + 1}{M-m} = 0$$

$$x_2 - 2 \frac{M+m}{M-m} x_1 + 1 = 0$$

and, for $j=2, 3, \dots$

$$x_{j+1} + 2 \frac{M+m}{M-m} x_j + x_{j-1} = 0$$

This difference equation can be solved by plugging in $x_j = q^j$ and solving for q :

LCB3

$$q^{j+1} - 2 \frac{M+m}{M-m} q^j + q^{j-1} = 0$$

$$q^2 - 2 \frac{M+m}{M-m} q + 1 = 0$$

Claim: roots are $\frac{M+m \pm 2\sqrt{Mm}}{M-m}$.

check: sum of roots is then $2 \frac{M+m}{M-m}$ ✓

product of roots is $\frac{(M+m)^2 - 4Mm}{(M-m)^2} = 1$ ✓

smaller root is

$$q = \frac{M+m - 2\sqrt{Mm}}{M-m} = \frac{(\sqrt{M} - \sqrt{m})^2}{(\sqrt{M} - \sqrt{m})(\sqrt{M} + \sqrt{m})}$$

$$= \frac{\sqrt{M} - \sqrt{m}}{\sqrt{M} + \sqrt{m}}$$

$$= \frac{1 - \sqrt{1/k}}{1 + \sqrt{1/k}}$$

where $k = \frac{M}{m}$

NOTE THE SQRT.

We get Theorem For any $x^{(0)} \in \mathbb{R}^{\infty}$ and any $m > 0$,
 $M > m$, \exists function F with $mI \leq \nabla^2 F \leq MI$
 (quadratic)
 with (under ASSUMPTION)
 A

$$\|x^{(k)} - x^*\|^2 \geq \left(\frac{1 - \sqrt{1/k}}{1 + \sqrt{1/k}} \right)^{2k} \|x^{(0)} - x^*\|^2$$

where x^* minimizes F and $k = M/m$, and hence

$$f(x^{(k)}) - f^* \geq \frac{m}{2} (\text{same})^{2k} \|x^{(0)} - x^*\|^2$$

PP. WLOG take $x^{(0)} = 0$. Then we use F as defined already, getting:

$$\|x^{(0)} - x^*\|^2 = \sum_{i=1}^{\infty} (x^*)_i^2 = \sum_{i=1}^{\infty} q^{2i} = \frac{q^2}{1 - q^2}$$

since T is tridiagonal and $\nabla f(x^{(0)}) = e_1$, we can show by

induction that $x_j \in \text{span}(e_1, e_2, \dots, e_j)$,

$$\begin{aligned} \text{so } \|x^{(j)} - x^*\|^2 &\geq \sum_{i=j+1}^{\infty} (x^*)_i^2 = \sum_{i=j+1}^{\infty} q^{2i} \\ &= \frac{q^{2(j+1)}}{1 - q^2} = \frac{q^{2j}}{1 - q^2} \|x^{(0)} - x^*\|^2 \end{aligned}$$

with $q = \frac{1 - \sqrt{1/k}}{1 + \sqrt{1/k}}$. (Last inequality follows from our original (1), just Taylor's Thm.)
 (BV (9.8))

LCBS

In comparison, the gradient method with $t \equiv \frac{1}{M}$ gave us

$$f(x^{(k)}) - f^* \leq \left(1 - \frac{1}{K}\right)^k (f(x^{(0)}) - f^*)$$

Compare $\frac{1 - \sqrt{1/K}}{1 + \sqrt{1/K}} \approx \left(1 - \sqrt{1/K}\right)^2 \approx 0.998$

if $K = 10^6$, while $1 - \frac{1}{K} \approx 0.999999$.

So lower bound indicates we may be able to do much better. The rest of Nesterov's Chapter 2 derives "optimal gradient" methods, but the argument is very

complicated! In the end, the simplest is:
NESTEROV OPTIMAL GRADIENT ALGORITHM (p. 81)

Choose $y^{(0)} = x^{(0)} \in \mathbb{R}^n$

For $k = 0, 1, 2, \dots$

$$\begin{cases} \text{let } x_{k+1} = y_k - \frac{1}{M} \nabla f(y_k) \\ \text{let } y_{k+1} = x_{k+1} + \rho (x_{k+1} - x_k) \end{cases}$$

where $\rho = (1 - \sqrt{1/K}) / (1 + \sqrt{1/K})$