# DISCOVERING ACCELERATION

- MOTIVATION
- PROPERTIES of $f$
- ALGORITHM & CONVERGENCE of SD
- BACKTRACKING (ARMIJO)
- DERIVING ACCELERATION
- ACCELERATED GRADIENT

---

## ① MOTIVATION

"THE SIMPLEST OPTIMIZATION METHOD", one of those topic of which you might say "OK, I know this", BUT actually there is MUCH MORE WEALTH BENEATH THE SURFACE !

- I ALREADY MENTIONED THE GRADIENT DESCENT METHOD LAST YEAR IN THE CONTEXT of THE MATRIX COMPLETION PROBLEM.
- THERE ARE MANY VARIANTS OF GRADIENT DESCENT ! e.g. ACCELERATED, CONJUGATE, BUT SIMPLE UNDERLYING COMMON PATTERNS ! PROJECTED STEEPEST DESC. COORDINATEWISE, STOCHASTIC...
- MUCH RESEARCH IN OPTIMIZATION FOCUSES ON <u>CONVERGENCE RATES</u>. BUT OTHER PROPERTIES ARE ALSO IMPORTANT, e.g. <u>ROBUSTNESS</u>.
- BASIC GRADIENT DESCENT IS ROBUST TO NOISE IN SEVERAL IMPORTANT WAYS, WHILE ACCELERATED GRADIENT DESCENT IS MUCH MORE BRITTLE. TRADE-OFF !
- We will fix our ideas on the specific case of : UNCONSTRAINED CONVEX OPTIMIZATION, i.e.,

$$\min_{x \in \mathbb{R}^n} f(x) \quad , \text{ called CONVEX PROGRAM.}$$

## ② SMOOTH AND STRONGLY CONVEX $f$

$f : \mathbb{R}^n \to \mathbb{R}$, TWICE DIFFERENTIABLE AND $(\alpha\text{-})$STRONGLY CONVEX, $(\mathcal{C}^2)$

i.e., $\exists \, \alpha > 0, \; \forall x \quad \nabla^2 f(x) \succeq \alpha I$

$(\Leftrightarrow \exists \alpha \;\; s.t. \;\; \forall x, y \quad f(y) \geq f(x) + \nabla f(x)^T (y - x) + \frac{\alpha}{2} \| y - x \|^2$

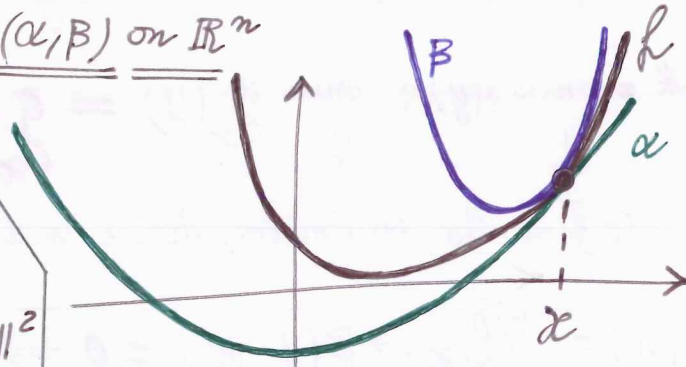Roughly speaking, we want the function $f$ to be "convex enough"

At the same time, we do not want $f$ to be "too convex":

($\beta$-)smoothness: $\quad \nabla^2 f(x) \preccurlyeq \beta I, \quad \beta < \infty$

$$\Leftrightarrow f(y) \leqslant f(x) + \nabla f(x)^T (y - x) + \frac{\beta}{2} \|y - x\|^2$$

In short, $f$ IS OF TYPE $(\alpha, \beta)$ on $\mathbb{R}^n$

Thm 1 (estimate on $p^*$)

Let $f$ be of type $(\alpha, \beta)$.
Then:
$$\frac{1}{2\beta} \|\nabla f(x)\|^2 \leqslant f(x) - p^* \leqslant \frac{1}{2\alpha} \|\nabla f(x)\|^2$$

(THIS GIVES A WAY TO STOP THE ITERATIONS) STOPPING CONDITION OF ITERATIVE ALGORITHM: $\|\nabla f(x_k)\| \leqslant \sqrt{2\alpha \varepsilon} \Rightarrow f(x_k) - p^* \leqslant \varepsilon$

"SMALL GRADIENT $\approx$ SOLVED PROBLEM"

ROUGHLY SPEAKING
$f$ CAN BE SQUEEZE
BETWEEN TWO
PARABOLAS !!

③ STEEPEST DESCENT (SD) [CAUCHY, 1847]

SEARCH DIR.

- Many methods (like SD) are of the form $x_{k+1} = x_k + t_k d_k$.

$t_k > 0$, STEPSIZE

- DESCENT TYPE: $f(x_{k+1}) < f(x_k)$

A) How to choose $d_k$?

For differentiable $f$: $\quad f(x_{k+1}) - f(x_k) = \underbrace{\nabla f(x_k)^T (t_k d_k)}_{\text{WE WANT DESCENT!} \leqslant 0} + o(\|t_k d_k\|)$
(Taylor)

SO: $\nabla f(x_k)^T d_k \leqslant 0$

Greedy choice: most decrease of $f$ at $x$

$$\begin{cases} \max\limits_{d_k} -\nabla f(x)^T d_k \\ \text{s.t. } \|d_k\| \leqslant 1 \end{cases}$$

solution is $d_k = -\nabla f(x)$

DIRECTION of SD

B) How to calculate $t_k$? line-search (LS)

- Exact LS: $\min\limits_{t \geqslant 0} f(x_k + t d_k) \rightsquigarrow t_k^{EX}$ is the unique minimizer if $f$ is stricly convex.

Can sometimes be computed. Good for theory.

Exact L$ is important for theoretical analysis:

**Thm2** Let $f$ be of type $(\alpha, \beta)$. Then $\underline{SD \text{ with exact } L\$}$ satisfies:

$\bigstar$
$$f(x_k) - p^* \leq \gamma^k (f(x_0) - p^*)$$

the optimal value

CONVG. FACTOR $\gamma := 1 - \dfrac{\alpha}{\beta}$

So SD converges to the exact solution $x_*$ for any $x_0$.

__Proof.__ $\quad x_+ = x - t\nabla f(x)$

Def. $\beta$-smooth: $\quad f(x_+) \leq f(x) + \underbrace{\nabla f(x)^T (-t \nabla f(x))}_{-t \|\nabla f(x)\|^2} + \dfrac{\beta}{2} t^2 \|\nabla f(x)\|^2$

Rewrite:

$$f(x - t\nabla f(x)) \leq f(x) + \left(\dfrac{\beta}{2} t^2 - t\right) \|\nabla f(x)\|^2, \qquad \text{valid } \forall t.$$

Minimizing both sides over $t \geq 0$

min. over $t$ gives $\quad \downarrow t = t^*$

$\qquad \qquad \downarrow t = t^*$

$\dfrac{\partial}{\partial t}(\cdots) = 0$
$\Rightarrow (\beta t - 1) \|\nabla f(x)\|^2 = 0$
$\Rightarrow t^* = \dfrac{1}{\beta}$

$$f(x_{k+1}) \leq f(x_k) - \dfrac{1}{2\beta} \|\nabla f(x)\|^2$$

$$f(x_{k+1}) - p^* \leq f(x_k) - p^* - \dfrac{1}{2\beta} \|\nabla f(x_k)\|^2$$

FROM THM.1 WE HAVE $\quad -\|\nabla f(x_k)\|^2 \leq -2\alpha (f(x_k) - p^*)$

$$f(x_{k+1}) - p^* \leq \underbrace{\left(1 - \dfrac{\alpha}{\beta}\right)}_{\gamma ''}(f(x_k) - p^*), \quad \text{use it recursively and get } \bigstar$$

---

• Let's bound the CONVERGENCE FACTOR. Call $K := \beta/\alpha$ the "condition number", $1 \leq K < \infty$
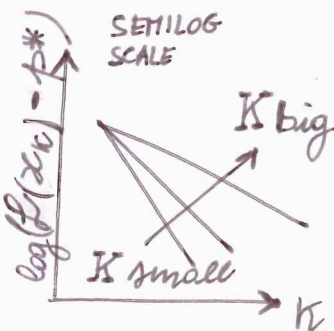
$\gamma^K = \left(1 - \dfrac{1}{K}\right)^K = (e^a)^K, \quad a = \log\left(1 - \dfrac{1}{K}\right) \leq -\dfrac{1}{K}$

$\leq \exp\left(-\dfrac{K}{K}\right)$

$K \to 1$, very good
$K \to \infty$, very bad

SEMILOG SCALE



$K$ big

$K$ small

- So we have **exponential convergence** (for some reason called **linear convergence** in optimization.)

- If we use $t_k = \frac{1}{\beta}$, $\forall k$ ($\underline{SD}$ with $\underline{CONSTANT}$ $\underline{STEPSIZE}$), we would get the same convergence bound.

IN GENERAL, FOR GENERIC $f$ (NOT of TYPE $(\alpha, \beta)$), WE WILL NOT USE EXACT LS !!!

④ **ARMIJO BACKTRACKING.** ( what we do IN PRACTICE ) [LARRY ARMIJO, 1966]

<u>Goal</u>: replace exact LS with something computationally cheaper, but still effective.

$-\nabla f(x_k)$ DIRECTION of $\dot{S}D$

$$\ell_k(t) = f(x_k) + c_1 t \nabla f(x_k)^T d_k$$

$$0 < c_1 < \tfrac{1}{2}$$

$$0 < c_2 < 1$$

<u>NB</u>: if $c_1 = 1$, we get the tangent line to $f(x_k + t d_k)$ at $f(x_k)$

Sufficient decrease condition:

$$\boxed{f(x_k + t_k d_k) < \ell_k(t_k)}$$

$\phi(t) := f(x_k + t d_k)$

$t^{(2)} = c_2 t^{(1)}$

$t^{(3)} = c_2 t^{(2)}$

$t^{(1)}$

$t^{(4)} = \ldots$

$t^{(5)}$

$\ell_k(t)$

$t_k^{EX}$    $t_k^{BT}$    $t = 1$

⇒ SHOW SLIDES ON ARMIJO BACKTRACKING.

<u>Thm 3</u> Let $f$ be of type $(\alpha, \beta)$. Then $\underline{SD}$ with Armijo LS $(c_1, c_2)$ satisfies

$$f(x_k) - p_* \leqslant \gamma^k (f(x_0) - p_*),$$

with $\gamma := 1 - 2 c_1 \min(\alpha, c_2/K)$, $K := \beta/\alpha$.

⟹ SHOW MOVIE $\dot{S}D$ ON QUADRATIC WITH ARMIJO.

# ⑤ SD ON A QUADRATIC $f$

- Let $x \in \mathbb{R}^m$, $A \in \mathbb{R}^{m \times m}$, and consider the quadratic objective function $f(x) = \frac{1}{2} x^T A x - b^T x$ [ WHEN $A > 0$, THIS QUADRATIC IS STRICTLY CONVEX AND HAS A UNIQUE GLOBAL MINIMUM $x^*$ ]

IN PARTICULAR,

- Let $f$ be of type $(\alpha, \beta) \Rightarrow$ SPECTRAL CONDITION $\boxed{\alpha I \preceq A \preceq \beta I}$

Since $\nabla f(x) = Ax - b$, we get $\min f(x) \Leftrightarrow Ax = b$; sol.: $x^* = A^{-1}b$.

$\uparrow$ $\nabla^2 f(x) = A$

- I.e., SD computes the sol. of the linear system $Ax = b$.

- One (not me in this 45 min. talk) can show that SD on this $f$ with constant stepsize $t = \frac{2}{\alpha + \beta}$ satisfies:

  "OPTIMAL"

$$f(x_k) - p^* \leq \left(\frac{K-1}{K+1}\right)^{2k} (f(x_0) - p^*)$$

- Doing the same exercise (BOUND THE CONVERGENCE FACTOR) we get

$$\gamma^k = \left(1 - \frac{2}{K+1}\right)^{2k} \leq \exp\left(-\frac{4k}{K+1}\right)$$

$\Rightarrow$ SHOW SD ON QUADRATIC WITH $K=1$. CONVERGES IN 1 ITERATION!

- BEST RATE FOR SD? NUMERICAL EXPERIMENTS SAY YES!

- IS THIS THE BEST METHOD? $\leadsto$ DEFINE...

$\Rightarrow$ SHOW ZIGZAG BEHAVIOR of SD. THIS MOTIVATES US TO LOOK FOR A BETTER METHOD!!

# ⑥ DERIVING ACCELERATION

- Let's start SD with constant stepsize $t = 1/\beta$ at $x_0 = \frac{1}{\beta} b$.

One can check by induction that

$$x_k = \left(\sum_{j=0}^{k} (I - A')^j\right) b',$$

where $A' = \frac{1}{\beta} A$ and $b' = \frac{1}{\beta} b$.

WHY DOES THIS CONVERGE TO $A^{-1}b$ ?

Recall that for all scalars $|x| < 1$,

$$\sum_{j=0}^{\infty} (1-x)^j = \frac{1}{x} \qquad \qquad \circledast$$

$\alpha I \lesssim A \lesssim \beta I$, $A' = tA$ and $t = \frac{1}{\beta}$; so $\frac{\alpha}{\beta} I \lesssim A' \lesssim I$, i.e., the eigenvalues of $A'$ lie within $(0,1]$. Hence $\circledast$ extends to the matrix case. **I.E., GRADIENT DESCENT IS COMPUTING A DEGREE $K$ (MATRIX POLYNOMIAL) APPROXIMATION of THE INVERSE FUNCTION of $A$ !!!**

- APPROXIMATION ERROR when truncating $\circledast$ to $K$ is $O\left((1-x)^K\right)$

- IN THE MATRIX CASE, this translates to $O\left(\|(I-A')^K\|\right) =$

$$= O\left(\|I - A'\|_2^K\right) = O\left(\left(1-\frac{1}{K}\right)^K\right) \quad \leftarrow \text{This is exactly the convergenc}$$
rate of SD that we
determined earlier !!!

$$\underbrace{\|I - tA\|_2}_{\uparrow} = \lambda_{max}\left(I - \frac{1}{\beta}A\right) = 1 - \frac{1}{\beta}\lambda_{min}(A) = 1 - \frac{\alpha}{\beta} = 1 - \frac{1}{K}$$

$I - A'$ is symm.
and $t = 1/\beta$

- Why we went through this exercise? Because now you see that TO IMPROVE THE CONVERGENCE RATE of GRADIENT DESCENT is equivalent to FIND A BETTER LOW-DEGREE POLYNOMIAL APPROXIMATION TO THE SCALAR FUNCTION $1/x$ !!! we'll be able to save a square root in the degree while achieving the same error!

- WE WANT TO MINIMIZE THE <u>RESIDUAL</u>:

$$r_K := \|(I - A\underline{q_K(A)})b\| \qquad \text{where } q_K(A) \text{ is a matrix polynomial}$$
approximation to $A^{-1}$: $q_K(A) \approx A^{-1}$

$$\leq \underbrace{\|I - Aq_K(A)\|}_{=: p_K(A)} \cdot \|b\| = \max_{\mu \in \lambda(A)} |p_K(\mu)| \cdot \|b\| \quad \boxed{f(A) = Q f(\Lambda) Q^T}$$

$\uparrow$ corresponding scalar polynomial

Relaxing this condition:

$$\leq \max_{\mu \in [\alpha, \beta]} |p_K(\mu)| \cdot \|b\|$$

● MINIMIZE THE RESIDUAL :

$$\min_{\substack{p_k \in \mathbb{P}_k \\ p_k(0)=1}} \max_{\mu \in [\alpha, \beta]} |p_k(\mu)|$$

VERY HARD OPTIMIZATION PROBLEM!
we are looking for a polynomial of degree K that is as small as possible on the location of the eigenvalues of A, namely on the interval $[\alpha, \beta]$.
At the same time, we have the normalization constraint $p_k(0) = 1$.

"THERE IS ONLY ONE BULLET IN THE GUN: IT'S CALLED THE CHEBYSHEV POLYNOMIAL."

(G.P.'s) "Sur les questions de minima qui se rattachent à la représentation approximative des fonctions"

CHEBYSHEV POLYNOMIALS (of 1$^{st}$ kind)     Чебышёв [1859]

Def. $T_0(x) = 1$, $T_1(x) = x$, $T_{m+1}(x) = 2x T_m(x) - T_{m-1}(x)$, $m \geq 1$.

LEMMA Let $m \geq 1$ and $q(x) = 2^{m-1} x^m + b_{m-1} x^{m-1} + \dots + b_0 \neq T_m(x)$,

"OPTIMALITY PROPERTY of CHEBYSHEV POLYNOMIALS"

THEN     $\max_{x \in [-1,1]} |q(x)| > \max_{x \in [-1,1]} |T_m(x)| = 1$.

PUT DIFFERENTLY, THE POLYNOMIAL $\overset{\text{HAVING THE FORM of } q(x)}{\text{THAT}}$ DEVIATES THE LEAST POSSIBLE FROM ZERO ON $[-1; 1]$ IS THE CHEBYSHEV POLYNOMIAL !!!

⟹ Show slide on G.P.'s.

● Suitably rescaled, G.P.'s MINIMIZE THE ABSOLUTE VALUE of $p_k$ IN A DESIRED INTERVAL $[\alpha, \beta]$ WHILE SATISFYING $p_k(0) = 1$:

$$P_m(x) := T_m\left(\frac{\alpha + \beta - 2x}{\beta - \alpha}\right) \Big/ T_m\left(\frac{\alpha + \beta}{\beta - \alpha}\right)$$

WHICH IS EXACTLY WHAT WE WANTED !!!

⟹ Show plot of this polynomial.

⑦ ACCELERATED GRADIENT METHOD

ERROR BOUND THAT COMES OUT OF THE CHEBYSHEV POL!'s:
(RATE OF NESTEROV'S FGM or AGM)

$$\|x_{k+1} - x^*\| \leq 2 \exp\left(-k \sqrt{\frac{2}{K}}\right) \|x_0 - x^*\|$$

(quite technical derivation)

THIS MEANS THAT, FOR LARGE K, WE GET QUADRATIC SAVINGS IN THE DEGREE WHILE ACHIEVING THE SAME ERROR !

Due to the recursive def. of G. Pol's, we get an iterative algorithm out of it. Transferring the recursive def. to our rescaled G. Pol's, we have:

$$P_{K+1}(\mu) = (t_K \mu + \gamma_K) P_K(\mu) + \delta_K P_{K-1}(\mu)$$

(the coeff's $t_K$, $\gamma_K$, $\delta_K$ can be worked out from the recurrence def.). Moreover, since $P_K(0) = 1$, we must have $\gamma_K + \delta_K = 1$. $\forall K$

$\Longrightarrow$ <u>UPDATE RULE</u>:

$$x_{K+1} = x_K - \underbrace{t_K(A x_K - b)}_{\nabla \ell(x_K)} + \underbrace{\delta_K(x_K - x_{K-1})}_{\text{ONLY THIS ADDITIONAL TERM!}}$$

$\Longrightarrow$ SHOW VIDEO of ACCELERATED GRADIENT!

RKs

- FINDING THE BEST POSSIBLE coeff.'s LEADS TO THE ABOVE CONVERGENCE RATE. WE USED ONLY $1^{st}$ ORDER INFORMATION!

- WORKS FOR ANY FUNCTION of TYPE $(\alpha, \beta)$, AND NOT JUST THE SPECIAL of A QUADRATIC OBJECTIVE $\ell$ THAT WE SHOW HERE!

- THIS IS WHAT NESTEROV SHOWED IN 1983!

- THE POLYNOMIAL APPROX. METHOD WAS KNOWN MUCH EARLIER IN THE CONTEXT of EIGENVALUE METHODS!

## REFERENCES