

Lecture #12



→ We've been thru the basics, culminating in a derivation of KKT.

Example: $\min_x \sum_i \log(1 + e^{x_i}) \text{ s.t. } \sum_i x_i = 1$

mentioned breifly
another deriv based on
 $= f(x) + \max_i f_i(x)$

NNLS: $\min_{\lambda} \frac{1}{2} \|Ax - b\|^2 \text{ s.t. } x \geq 0$

$$\phi(x) = \max \left\{ f(x) - f^*, h_1(x), \dots, h_m(x) \right\}$$

$$L(x, \lambda) = \frac{1}{2} \|Ax - b\|^2 - \lambda^T x$$

Finding direct solution req?

$$\begin{cases} \frac{\partial L}{\partial x} & A^T(Ax - b) - \lambda = 0 \\ & \lambda^T x = 0 \\ & x \geq 0, \lambda \geq 0 \end{cases} \quad \begin{array}{l} \text{---(1)} \\ (\text{nonlinear!}) \end{array}$$

→ Usually, KKT system not easy to solve directly, iterative approaches can be taken

→ Even more simply: consider a func $f(x)$:

$$f(x) = \frac{1}{2} x^T A x + b^T x, \quad A \succ 0$$

• $\nabla f(x) = Ax + b = 0 \Rightarrow$ solve linear sys $Ax = b$.

• If $b \in \text{Ran}(A)$, system has a soln.

Direct methods: Cholesky of A : $\frac{1}{3} n^3$ flops

Iteration methods: (Most area within numerical linear algebra
since well over 40 years!)

Many key Opt. algorithm ideas originally stems from
trying to solve $Ax = b$!

→ At an abstract level:

$$\min f(x)$$

$$(1) \Rightarrow \nabla f(x) = 0 \Rightarrow x = (\nabla f)^{-1}(0) = \nabla f^*(0)$$

~~(2)~~ $\nabla f(x) = 0 \Rightarrow -\alpha \nabla f(x) = 0 \quad \text{if } \alpha > 0 \text{ too}$

$$\Rightarrow x - \alpha \nabla f(x) = x \text{ must hold}$$

Suggests an algorithm! $x_{k+1} = x_k - \alpha_k \nabla f(x_k)$!

Or visit <http://www.cs.cmu.edu/~help>.
For comments or questions, send email to Help@cs.cmu.edu
Software refuses to honor your request.
of paper and other resources, the spooling
Since printing a binary file is usually a waste
So I printed this page instead.
The file that someone asked me to print
really looked like a binary file.

(1.5)

$$-\alpha \nabla f(x) = 0$$

$$\Rightarrow -\alpha \nabla^2 f(x)^{-1} \nabla f(x) = 0 \quad \text{shows if } S \text{ is invertible}$$

we can form $x - \alpha (\nabla^2 f(x))^{-1} \nabla f(x)$: Newton method.

Other "creative" penultimate sometimes

Eg

$$f(x) - g(x)$$

$$\nabla(f(x) - g(x)) = 0$$

$$\nabla f(x) = \nabla g(x)$$

Numerous
penultimates to
play with them

(soon we'll see a practical choice)

→ this is the "algebraic" derivation of gradient descent ⑤

→ Say $f(x)$ is non-diff. Then,

$$0 \in \partial f(x)$$

Now, actually $(\partial f)^{-1}(0)$ is a multivalued map: a set.

How about

$$x \in x - \alpha \partial f(x)$$

as before? less clear. because this generates "

$$x_{k+1} \in x_k - \alpha_k \partial f(x_k)$$

(with careful choosing α_k , using any $g_k \in \partial f(x_k)$ works!)

$$\rightarrow 0 \in \partial f(x) \Rightarrow x \in x + \alpha \partial f(x)$$

$$\Rightarrow (Id + \alpha \partial f)(x) \ni x$$

$$\Rightarrow x = (Id + \alpha \partial f)^{-1}(x)$$

$$\Rightarrow x_{k+1} = (Id + \alpha \partial f)^{-1}(x_k)$$

(Fixed-point iterate).

This is actually Rockafellar's famous "proximal point algorithm" → a method that ultimately includes all sorts of update iterations (Is a proper choice of $\underline{\partial f}$)

→ Q: A priori no reason for above "guesses" to work → but they are a good starting point.

Q: Say we have $\phi(x) = f(x) + g(x)$

(3)

$$\textcircled{a} \quad x_{\text{new}} = x_u - \alpha_u \nabla \phi(x_u)$$

Sounds to be ok if $\nabla \phi(x_u) = \nabla f(x_u) + \nabla g(x_u)$

- \textcircled{b} Suppose $\phi(x) = f(x) - g(x)$ where both f, g , convex.
How should we go about minimizing such a ϕ ?



Key idea is: • Start with original frame and a guess, say x_u :

- OPT • Build an approx. to $\phi(x)$, centered at x_u , with the aim
that approx is easy to optimize; call it $m(x; x_u)$
• Obtain $x_{\text{new}} = \arg \min m(x; x_u)$

- Gradient descent: $\phi(x + s) \approx \phi(x) + \underbrace{\langle \nabla \phi(x), s \rangle}_{m(x, s)} + \frac{1}{2} \|s\|_2^2$

$$x_{\text{new}} = \arg \min_s x_u + s.$$

$$\begin{aligned} s &= \arg \min_s \phi(x) + \langle \nabla \phi(x), s \rangle + \frac{1}{2} \|s\|_2^2 \\ &= \nabla \phi(x) + \cancel{s} = 0 \\ &\Rightarrow s = -\cancel{\alpha} \nabla \phi(x) \end{aligned}$$

$$\boxed{s \equiv x_{\text{new}} = x_u - \alpha_u \nabla \phi(x_u)} \quad \therefore \text{G.D.}$$

- $\phi(x) = f(x) - g(x)$; f, g, convex

$$g(x+\delta) \geq g(x) + \langle \nabla g(x), \delta \rangle$$

$$\begin{aligned} \text{So } \phi(x+\delta) &= f(x+\delta) - g(x+\delta) \\ &\leq f(x+\delta) - g(x) - \langle \nabla g(x), \delta \rangle \\ &\quad \boxed{m(x+\delta)} \end{aligned}$$

$$x_{\text{new}} = x_u + s$$

$$s = \arg \min_s f(x+\delta) - \langle \nabla g(x), s \rangle$$

(3.5)

Example of CCC

Suppose we wish to minimize

$$\Phi(x) = \min_{x > 0} \sum_i \log \frac{\det(x + a_i)}{\sqrt{\det(x) \det(a_i)}}$$

This problem is called Frobenius-mean:

$$\begin{aligned}\Phi(x) &= \sum_i [\log(x + a_i) - \frac{1}{2} \log(x a_i)] \\ &= \sum_i \left[\underbrace{\frac{1}{2} \log(x a_i)}_{f(x)} - \underbrace{(-\log(x + a_i))}_{g(x)} \right]\end{aligned}$$

Obtain linearization:

$$m(x|s) = \frac{\partial}{\partial s} \sum_i -\frac{1}{2} \log((x + s)a_i) = \frac{1}{2} \langle (x + s)^{-1}, s \rangle$$

$$\text{writing} \quad \sum_i \left(\frac{1}{2} (x + s)^{-1} - \frac{1}{2} (x + a_i)^{-1} \right) = 0$$

$$\Rightarrow n(x + s)^{-1} = \sum_i (x + a_i)^{-1}$$

$$\Rightarrow (x + s)^{-1} = \frac{1}{n} \sum_i (x + a_i)^{-1}$$

$$\Rightarrow x + s = x + s = []^{-1}$$

Does this work? Converse? Etc?

Actually this particular search always do global optimum!

~~# Beyond CC~~

- # However: Not always do.
- # The art of using CCP lies in decomposing a given non-convex function into useful parts that allow easier optimization.
- # It ensues to etc.
- # Reading assignment: "On the convergence of the Concave-Concave Procedure"
Srikanth, Lanctot
→
- # More generally, a fairly large heuristic area: D.C. Programming
etc → we might revisit it, depending on instances.
→
- ~~This is another example~~ For several examples of using CCP in Stats & ML:
 - (a) "The Concave-Concave Procedure" — Yurille & Rangarajan
 - (b)

Note: This method is often a useful baseline to beat when dealing with nonconvex problems!

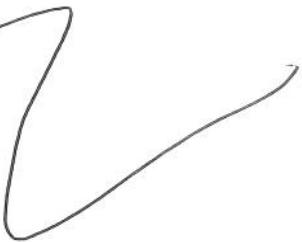
(9)

This repeated model-bound minimization is called
Concave-concave programming!

~~→ convex~~ (one member of a class of methods called Majorize-Minimize, etc.)
 (We'll see some more of this later in class) → But really handy
 tool to have in repertoire.

A few words about gradient-descent methods:

Say $\min_{x \in \mathbb{R}^n} f(x)$ s.t. $x \in \mathbb{R}^n$
 s.t. $\nabla f(x) \neq 0$;
 Look at $x_d = x - \alpha \nabla f(x)$



Methods without explicit derivatives

$$\min_{x \in \mathbb{R}^n} f(x) = f(x_1, \dots, x_n)$$

at each iteration optimize over just 1 coordinate
 (or max it handle);
 Called coordinate descent or Block-CG

If $f(x) = f(\underline{x}_0, \underline{x}_1)$, $\underline{x}_0 \in \mathbb{R}^{n_1}$, $\underline{x}_1 \in \mathbb{R}^{n_2}$

then $\underline{x}_{0,n} \leftarrow \arg \min_u f(u, \underline{x}_1)$
 $\underline{x}_{1,n} \leftarrow \arg \min_v f(\underline{x}_{0,n}, v)$

Observe that $\{f(\underline{x}_{0,n}, \underline{x}_1)\}$ is ↓.

The two-block case of alternating minimization is special →
 we will return to it later in the course.

(5)

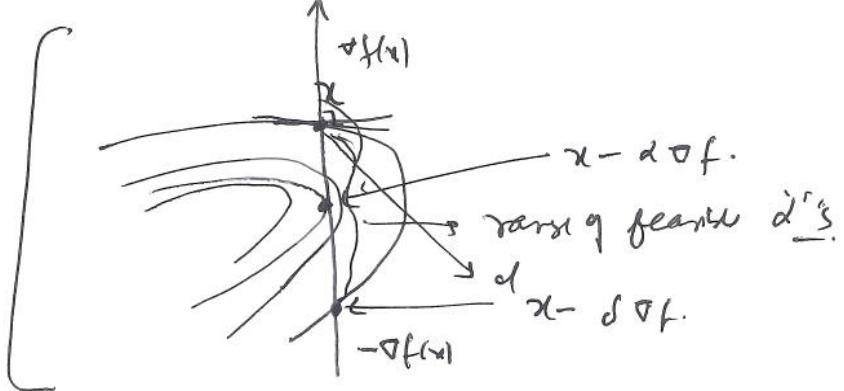
Back to gradient-methods:

descent

$$x \leftarrow x - \alpha \nabla f$$

for small enough $\alpha > 0$
we have descent

more generally $x \leftarrow x + \alpha d$ ← descent direction

This picture is
important.

Code! while !done

$$x_{\text{new}} \leftarrow x_{\text{old}} + \alpha_d d_{\text{new}}$$

Update curr, curr Sondow.

end.

dm chosen to ensure $f(x_{\text{new}}) < f(x_{\text{old}})$ (Unless other strategy)# d_{new} : descent direct. $\langle \nabla f(x_{\text{old}}), d_{\text{new}} \rangle < 0$ # Read "Numerical Optimization" by Nocedal & Wright to learn much more about different strategies for selecting $\underline{\alpha}$ and $\overline{\alpha}$ # How about breaking the rules a little bit?# Do NOT insist on $f(x_{\text{new}}) < f(x_{\text{old}})$

allows more freedom in choosing step sizes.

But why? Maybe empirical evidence!

Observe: $m(x, \alpha) = f(x) + \langle \nabla f(x), \alpha \rangle + \frac{\alpha}{2} \| \alpha \|^2$ was ① model

How about the model:

$$m(x, 1) = f(x) + \langle \nabla f(x), 1 \rangle + \frac{1}{2} \| 1 \|^2$$

But now optimising over $\alpha \geq 0$

(6)

$$\bullet x_{k+1} = x_k - \alpha S_k \nabla f(x_k) \quad : \text{Newton, Quasi-Newton}$$

~~What is~~ the issue is: How to select S_k efficiently... (LSFGS, etc.)

• Let us instead choose $S_k = \beta I$; $S_k = \beta I$;

Find a good β !

$$S_k \approx [\nabla^2 f(x_k)]^{-1};$$

$$H \Delta x = \Delta g$$

$$H^{-1} \Delta g = \Delta x$$

$$\min_{\beta} \|\beta \Delta g - \Delta x\|^2 : \quad \beta^2 \langle \Delta g, \Delta g \rangle - 2\beta \langle \Delta g, \Delta x \rangle \\ \Rightarrow \beta = \frac{\langle \Delta g, \Delta x \rangle}{\|\Delta g\|^2}.$$

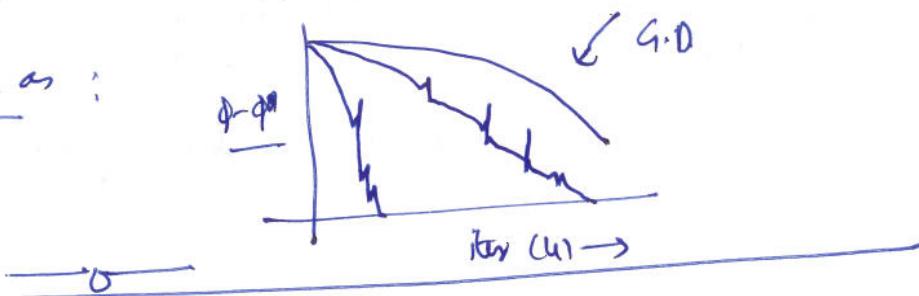
alternatively we can search for

$$\|\beta \Delta x - \Delta g\|^2 \Rightarrow \beta = \underbrace{\frac{\langle \Delta x, \Delta g \rangle}{\|\Delta x\|^2}}_{\text{Step is } \frac{1}{\beta}} \Rightarrow \beta^{-1} = \frac{\|\Delta x\|^2}{\langle \Delta x, \Delta g \rangle}.$$

→ These steps are known as "Barzilai-Borwein" step sizes!

→ Empirically remarkable! "Barzilai & Borwein": Two-point step sizes gradient methods. IMA J. Num. Anal.; 1988.

Lead to picture such as:



Let's begin now a more careful study of theoretically evaluate performance of methods. But always, for a real problem, from computer to get a better idea! This is just a guideline.

7

$$O\left(\frac{1}{e^k}\right) \cdot \frac{1}{\sum_{i=0}^k} = O\left(\frac{1}{e^k}\right) \Rightarrow \log k = -2^k \log 2 \quad (y_k = 2^k \log 2)$$

$$2^k = \frac{\log k}{\log 2}$$

To judge methods / their performance, we'll ~~use~~ could
~~use~~ use:

* Analytical complexity:

* Arithmetic Complexity:

What does it mean to solve? a problem?

Alg.: obtain a point \bar{x} s.t. $f(\bar{x}) \leq f^* + \epsilon$

When running depends on $\text{poly}\left(\frac{1}{\epsilon}\right)$ apply $\log\left(\frac{1}{\epsilon}\right)$, $\log_2\left(\frac{1}{\epsilon}\right)$

Alternatively: $\|\nabla f(\bar{x})\| \leq \epsilon$. maybe measured.

Or even: $\|\bar{x}_k - (\bar{x}_k - \alpha_k \nabla f(\bar{x}_k))\|$. and so on.

—

Model of Computation:

Types of oracles and examples of these oracles.

Determ	[Zeroth	—	More generally: oracles with <u>error</u> . that does not satisfy $E[\eta] = 0$. (<u>Biased oracle</u>)
		First	—	
		Second	—	
Stoch	[

8

Class of functions: $C_L^{k,p}(\mathcal{X})$: k times cont. diff func.when p th deriv is Lipschitz cont. with const L s.t.

$$\|f^{(p)}(x) - f^{(p)}(y)\| \leq L\|x-y\|, \quad \forall x, y \in \mathcal{X}.$$

Useful Lemma (Hw): $f \in C_L^{k,1}(\mathcal{X})$ iff $\|f''(x)\| \leq L$ ~~s.t. $\forall x$~~ What is l.e. of $f(x) = x^{\alpha}$?; $f(x) = \log(1+x)$.

Note: $-\log x \notin C_L^{k,1}(\mathbb{R}_+)$ for any L | we can
 $\max(0, -x) \in C_L^{0,1}(\mathbb{R})$. | $C_L^1 \equiv C_L^{1,1}$

~~Let us finally analyze ~~the~~ to~~# Theorem: Let f be convex, $C_L^1(\mathbb{R}^n)$; let $\alpha_k = 1/L$.Let $x_{k+1} = x_k - \frac{1}{L} \nabla f(x_k)$; then

$$f(x_k) - f(x^*) \leq \frac{2L\|x_0 - x^*\|^2}{Lk} = O\left(\frac{1}{k}\right)$$

i.e. $O(1/k)$ iter to get ϵ' acc.