Lecture 4: Additional Large Scale Algorithms

D UNCONSTRAINED METHODS

$\min_{x} f(x)$, assuming $f$ is sufficiently smooth.

① (non-linear CG).     Linear: Hestenes, Stiefel '50's.
                       Non-linear: Fletcher, Reeves '60s

<u>Linear case</u>: Solve $Ax = b$ where $A > 0$

In optimization, this is $\min_{x} \frac{1}{2} \| \tilde{A} x - \tilde{b} \|^2$,   $\tilde{A}$ full col. rank
                              $A := \tilde{A}^T \tilde{A}$, $\tilde{b} = \tilde{A}^T \tilde{b}$.

Creates directions $\{p_i\}$ st. $\langle p_i, A p_j \rangle = 0$ if $i \neq j$  "A-orthog.",

   $P_k = b - A x_k + \beta_k \cdot P_{k-1}$
   $X_{k+1} = X_k + \alpha_k P_k$,   $\alpha_k$ via exact linesearch
   Then: $\left[ \begin{array}{l} X_{k+1} = \operatorname{argmin}_{x} \frac{1}{2} \| \tilde{A} x - \tilde{b} \|^2 \\ \phantom{X_{k+1} =} \text{st. } x \in X_0 + \operatorname{span}(p_0, \cdots, p_k) \end{array} \right]$

<u>Non-linear Case</u>: linesearch inexact, Thm not true regardless, no magic, more sensitive.
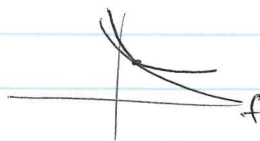   But can still work.                                                    $\underbrace{\phantom{more sensitive}}$
 • Doesn't work well w/ constraints either.                               many parameters.
 • Nemirovsky/Yudin showed it can perform worse than gra.desc.           cf Hager, Zhang
 • Theory "iffy"
 • Q-N. simpler

② Quasi-Newton Methods

Quadratic Approx. of $f$ at $x_k$:

$$q_k(x) = f(x_k) + \langle \nabla f(x_k), x-x_k \rangle + \tfrac{1}{2}\langle x-x_k, B_k \cdot (x-x_k)\rangle$$

and $\tilde{x}_{k+1} := \arg\min_x q_k(x)$

and $x_{k+1} = \tilde{x}_{k+1}$ or $x_{k+1} = x_k + \alpha \cdot (\tilde{x}_{k+1} - x_k)$, line search.

1) $B_k := L \cdot I$, grad desc.

2) $B_k = \nabla^2 f(x_k)$, Newton: $\tilde{x}_{k+1} = -\nabla^2 f(x_k)^{-1}\nabla f(x_k)$

3) $B_k$ — better than $L \cdot I$, cheaper than Newton: "quasi-Newton"

For any $B_k$, $\quad q_k(x_k) = f(x_k), \quad \nabla q_k\big|_{x_k} = \nabla f\big|_{x_k}$

but ask for more: $\quad \nabla q_k\big|_{x_{k-1}} = \nabla f\big|_{x_{k-1}} \qquad (*)$

Defining $\quad s_k = x_{k+1} - x_k$

$y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$, $\quad \nabla q_k\big|_{x_{k-1}} = \nabla f(x_k) + B_k(x-x_k)\big|_{x_{k-1}} = \nabla f(x_k) + B_k(-s_{k-1})$

$\overset{(*)}{=} \nabla f(x_{k-1})$

i.e. "$B_k s_{k-1} = y_{k-1}$" "Secant Eq'n"

Can we solve this eq'n?

(and maintain $B_k > 0$) need $\langle s_{k-1}, B_k s_{k-1}\rangle > 0$, i.e., $\underline{\langle s_{k-1}, y_{k-1}\rangle > 0}$

a necessary condition. "Curvature Condition"

$\Big(> 0$ always if $f$ is convex, i.e., $\nabla f$ monotone $\Big)$

Match gradients at $x_{k-2}, x_{k-3}, \ldots$?

No, hard to ensure $B_k > 0$ then.

Instead, to still use old information, write $B_k$ as a low-rank update to $B_{k-1}$.
(i.e. "close")

BFGS: most popular, impose $\underset{H_{k+1}}{\underbrace{B_{k+1}^{-1}}}$ close to $\underset{H_k}{\underbrace{B_k^{-1}}}$

$$H_{k+1} = (1-\rho_k s_k y_k^T) H_k (I-\rho_k y_k s_k^T) + \rho_k s_k s_k^T, \quad \rho_k = \frac{1}{\langle y_k, s_k\rangle} < \infty$$

$H_k = \frac{\langle y_k, s_k\rangle}{\langle y_k, y_k\rangle}$ Barzilai-Borwein. When to use...

③ limited-memory BFGS: See Nocedal & Wright.

Saves memory, similar performance. $\boxed{\text{WORKHORSE ALGO.}}$

④ Inexact / Matrix-Free Newton (aka Newton-CG)

to solve $\tilde{x}_{k+1} = x_k - \underbrace{\nabla^2 f(x_k)^{-1} \nabla f(x_k)}$

~~Approximate Hessians~~

ie., solve $\nabla^2 f(x_k) \cdot p = b$.

Use linear CG, which only needs to have a routine

$$p \mapsto \nabla^2 f(x) \cdot p$$

unlike a direct method (Gauss Elim, ie. LU or Cholesky).

"Sensitive, since need a good tolerance, but can be
state-of-the-art.

-Precondition CG w/ quasi-Newton $H_k$.

⑤ Nonlinear Least Squares (not assuming convexity. Ref: §10 in Nocedal & Wright)

$$f(x) = \tfrac{1}{2} \| \vec{r}(x) \|^2 = \tfrac{1}{2} \sum_{i=1}^{m} r_i^2(x), \quad \text{assuming } r_i \text{ smooth.} \quad \underline{Ex}: \text{PDE constr. optim.}$$

Jacobian of $\vec{r}(x)$ ($\vec{r}: \mathbb{R}^n \to \mathbb{R}^m$) is

$$\left( J(x) \right)_{i,j} = \frac{dr_i}{dx_j} \quad \sim \quad J(x) = \begin{bmatrix} \nabla r_1(x)^T \\ \vdots \\ \nabla r_m(x)^T \end{bmatrix} \quad \begin{matrix} m \times n \text{ matrix,} \\ \underline{not} \text{ symmetric.} \end{matrix}$$

Then,

$\left. \right\}$ ie. $f = \tfrac{1}{2} \| J x \|^2$

gradient of $f$ is $\nabla f(x) = J(x)^T \cdot r(x)$       $= 0$ if linear least-sq.

Hessian of $f$ is $\nabla^2 f(x) = \underbrace{J(x)^T J(x)}_{"} + \sum_{i=1}^{m} r_i(x) \overbrace{\nabla^2 r_i(x)}$   $\left( \begin{matrix} \text{In general, annoying} \\ \text{to compute} \end{matrix} \right)$

"for free" from 1st order info.

$\underline{\text{(Gauss-Newton}}$ is like Newton, but instead of $B_k = \nabla^2 f(x_k)$, use $B_k = J(x_k)^T J(x_k)$

$\left[ \begin{matrix} \text{Can derive it by linearizing } \vec{r}(x_k + p) \simeq \vec{r}(x_k) + J_k^T p \\ \text{rather than linearizing } f(x_k + p) \end{matrix} \right]$

Also a
"work-horse"
algo.

\* $\underline{\text{Levenberg-Marquardt}}$ is a trust-region version, $\|p\| \leq \Delta_k$, Lagrangian is a Tikhonov pnd

$\sim$, $\qquad B_k = J_k^T J_k + \lambda I, \quad \lambda > 0.$

▷ CONSTRAINED PROBLEMS

   ① Active-set style methods : "glue" some variables, pretend rest are unconstrained. Eg. L-BFGS-B.

   ② Penalty Methods, $\min f_0(x)$    ⟶    $\min f_0(x) + M/2 \, h^2(x)$,
                     $h(x) = 0$

     Solve a seq. as $\mu \longrightarrow +\infty$ (and "warm-start" each).
     Quick + Dirty ⟶ it's used a lot, but has many issues.
          Not used in any serious package.

     · Exact penalty methods nicer, but lack smoothness, so subproblem harder.

   ③ Augmented Lagrangian

    (P)    $\min f_0(x)$   ⟺   $\min f_0(x) + M/2 \, h^2(x)$    $(P_\mu)$
            $h(x) = 0$                 $h(x) = 0$

            Lagrangian is $\mathcal{L}(x, \nu) = f_0(x) + M/2 \, h^2(x) + \langle \nu, h(x) \rangle$

      If we knew $\nu^*$, then KKT for $x^*$ means    i) $x^* \in \mathrm{argmin}_x \, \mathcal{L}(x, \nu^*)$

                                         2) $h(x^*) = 0$.
        and if sol'n to (1) is unique, (2) is automatic!

      We "augmented" with $M/2 \, h^2(x)$ because   a) it is allowed
                                     b) it prevents getting $\mathrm{argmin} \, \mathcal{L}(x, \nu) = +\infty \dots$

      <u>Method</u>:   solve   $x_{k+1} \in \mathrm{argmin}_x \, \mathcal{L}(x, \nu_k)$
                      $\nu_{k+1} = \nu_k + \mu \, h(x_k)$   (like a gradient step)

   ie.,
     run gradient <u>ascent</u> on dual problem of $(P_\mu)$. If $f_0$ is strongly convex, this is rigorous.

     · For inequality constraints, see Nocedal + Wright,   Lancelot Software
                                                        GALAHAD (Fortran)

④ SQP: Sequential Quadratic Programming ($18 in Nocedal & Wright) (non-convex)

generalize Newton to allow constraints: linearize constraints,
so each iteration is a quadratic program (QP).
Use a trust-region (like line search)
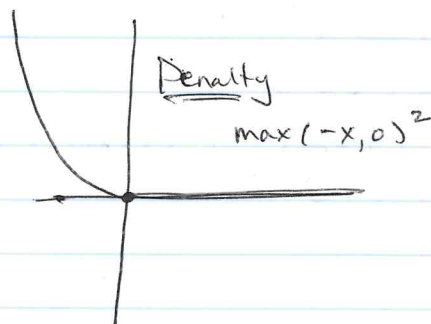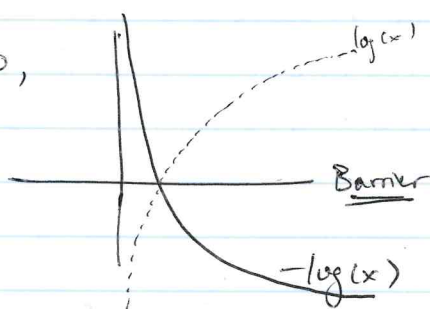
eg. SNOPT, KNITRO, TRON

⑤ IPM: interior pt. methods. (convex)

Equality, for convex problem, is always $Ax=b$, so already like
ie., equality-constrained Newton.

For $x > 0$ (or $X \succeq 0$), use ~~pos~~ special barrier
$-\log(x)$    $-\log \det(X)$

eg., $X > 0$,



Barrier $\qquad$ $-\log(x)$

Penalty
$\max(-x, 0)^2$

Nesterov/Nemirovski - analyzed Newton for self-concordant barriers
ie., $|f'''(x)| \le 2 (f''(x))^{3/2}$

analysis is affine invariant,
just like algo.

See Boyd & Vandenberghe's book

ie. PPA

6) ADMM , Douglas-Rachford          See Boyd's 2011 monograph

$$\min_x f(x) + g(x) \quad \Longleftrightarrow \quad \min_{x,z} f(x) + g(z)$$
$$\text{assume}$$
$$(\text{convex})$$
$$\text{st. } x - z = 0 \quad \longrightarrow \text{more generally, } Ax + Bz = c.$$

$$\Updownarrow$$
$$\min_{x,z} f(x) + g(z) + \rho/2 \|x-z\|^2$$
$$\text{st. } x - z = 0.$$

Aug. Lagr. idea
$$\left\{\binom{x}{z}\right\}_{k+1} \in \underset{\binom{x}{z}}{\arg\min} \left( \mathcal{L}_\rho(x,z,\overset{\overbrace{\text{dual (was "}y\text{" before)}}}{y_k}) = f(x) + g(z) + \rho/2 \|x-z\|^2 + \langle y_k, x-z \rangle \right)$$

$$y_{k+1} = y_k + \rho (x_{k+1} - z_{k+1}).$$

ADMM approximates this à la Gauss-Seidel

a)  $x_{k+1} \in \underset{x}{\arg\min} \ \mathcal{L}(x, z_k, y_k) = \underset{x}{\arg\min} \ f(x) + \rho/2 \|x - (z_k - \frac{1}{2\rho} y_k)\|^2$
     $= \text{prox}_{\rho^{-1}f} (z_k - \frac{1}{2\rho} y_k)$

b)  $z_{k+1} \in \underset{z}{\arg\min} \ \mathcal{L}(x_{k+1}, z, y_k)$

c)  $y_{k+1} = y_k + \rho (x_{k+1} - z_{k+1}).$

- Stronger Convergence Results than Aug. Lag.
- Slow Sometimes
- $\rho$ is "magic" parameter $\rightarrow$ for fast convergence, it must be chosen with

Douglas-Rachford   Bauschke + Combetts ed-2 §28.3
   $f, g \in \Gamma_0(\mathbb{R}^n)$, assume ∃ $x$ st. $0 \in \partial f(x) + \partial g(x)$  ie. Cor 27-6 ∃ sol'n
(P) $\min_x f(x) + g(x)$                                              and ri(dom g)∩ri(dom f) ≠ ∅
(D) $\min_u f^*(-u) + g^*(u)$,                                        or polyhed.
       $0 < \lambda < 2, \gamma > 0 \ (\simeq \rho^{-1})$, any $y_0$,
       $x_k = \text{prox}_{\gamma g}(y_k)$              then   $y_k \to y$, and if
       $z_k = \text{prox}_{\gamma f}(2x_k - y_k)$            $x = \text{prox}_{\gamma g}(y)$ , $x_n \to x$,
       $y_{k+1} = y_k + \lambda(z_k - x_k)$              $x$ is primal optimal, $z_k \to x$ too.

cf. Cevher, Becker, Schmidt '14 review

## ⑦ PRIMAL-DUAL METHODS

$$\min \ g(x) + \tilde{h}(\underbrace{Ax}_{z}) \quad \longrightarrow \quad \min_{x,z} \ g(x) + \tilde{h}(z)$$
$$Ax = z$$

$$h(x) = \tilde{h}(Ax)$$

→ apply ADMM, update requires $\text{prox}_{\tilde{h} \circ A}$

i.e. $\text{prox}_{\tilde{h}}$ cheap $\not\Rightarrow$ $\text{prox}_h$ cheap.

Chambolle-Pock "primal-dual hybrid gra." or "preconditioned ADMM"

    in prox update for $z$, perform in a <u>scaled</u> norm,

$$\|z - z_k\|_M^2, \ \text{w,} \ \|z\|_M^2 \equiv \langle z, Mz \rangle, \quad M := \sigma^{-1} I - A^T A, \ \text{so} \ \sigma < \frac{1}{\|A\|_2^2} \Rightarrow M > 0.$$

    Cancels out non-separable term.

More generally, L-Condat '11 and "A forward-backward view of some primal-dual —"
                     by Combettes, Condat, Pesquet, Vu

$$\min_x \ f(x) + g(x) + h(Ax),$$

       assume $\nabla f$ Lipschitz, $A$ a matrix, $g, h$ have easy prox, $f, g \in \Gamma_0(\mathbb{R}^n)$
                            $(m \times n)$                                           $h \in \Gamma_0(\mathbb{R}^m)$

assuming CQ,

$$0 \in \nabla f(x) + \partial g(x) + A^T \underbrace{\partial h(Ax)}_{y}$$

                     $y \in \partial h(Ax) \iff Ax \in \partial h^*(y)$ since $\partial h^* = \partial h^{-1}$.

So, solve KKT conditions / saddle-pt conditions

1) $0 \in \nabla f(x) + \partial g(x) + A^T y$

2) $Ax \in \partial h^*(y)$

$\Big\}$   ∈

[abusing notation since nonlinear]

$$-\underbrace{\begin{bmatrix} -\nabla f & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}}_{T_2} \in \underbrace{\begin{bmatrix} \partial g & A^T \\ -A & \partial h^* \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}}_{T_1} \quad \tilde{x} = \begin{bmatrix} x \\ y \end{bmatrix}, \ \text{find}$$

$$0 \in T_1 x + T_2 x$$

<u>Careful!</u> multiply a row by $-1$ won't change sol'n, but makes hard to find,
    Since operator no longer monotone!

<u>More general than optimization.</u>

to solve $0 \in T_1 x + T_2 x$ via Forward-Backward,

$$\vec{x}_{k+1} = \underbrace{(I + T_1)^{-1}}_{\text{like prox}} \underbrace{(I - T_2)}_{\text{like grad.}} \vec{x}_k \qquad (\text{assuming we've scaled to make}$$
$$T_2 \ 1\text{-Lipschitz})$$

Condat's Primal-dual:

$$X_{k+1} = (V + T_1)^{-1}(V - T_2)\vec{x}_k, \qquad V = \begin{bmatrix} \tau^{-1} I & -A^T \\ -A & \sigma^{-1} I \end{bmatrix} > 0 \text{ iff}$$
$$\sigma \tau > \|A\|^{-2}$$

$$\left( \begin{bmatrix} \tau^{-1} I & -A^T \\ -A & \sigma^{-1} I \end{bmatrix} + \begin{bmatrix} \partial g & A^T \\ -A & \partial h^* \end{bmatrix} \right)^{-1} = \left( \begin{bmatrix} \tau^{-1} I + \partial g & 0 \\ -2A & \sigma^{-1} I + \partial h^* \end{bmatrix} \right)^{-1}$$

decoupled! solve for $x$ first,
then $y$. (Back substitution!)

⑧ Alternating min, coordinate descent $\boxed{\text{Often } f(\cdot, y), f(x, \cdot) \text{ convex, not jointly cvx}}$

$$\min_{x, y} f(x, y) \longrightarrow X_{k+1} \in \operatorname{argmin}_{x \in \Omega_x} f(x, y_k) \qquad (\text{or a gradient step})$$
$$\text{`or more`}$$
$$(x, y) \in \Omega_x \times \Omega_y$$
$$Y_{k+1} \in \operatorname{argmin}_{y \in \Omega_y} f(x_{k+1}, y)$$

Convergence weak. Efficiency depends on problem structure.

Proximal methods better, e.g. see discussion in
PALM, Bolte, Sabach, Teboulle, i.e., $X_{k+1} = \operatorname{argmin}_{x \in \Omega_x} f(x, y_k) + \frac{M}{2} \|x - x_k\|^2$

$$Y_{k+1} = \operatorname{argmin}_{y \in \Omega_y} f(x_{k+1}, y) + \frac{M}{2} \|y - y_k\|^2$$