

Lecture 16 - FW & AFW

Tuesday, March 16, 2021 14:19

- today:
- properties of FW
- AFW alg.

Properties of FW

$$1) f(x_t) - \min_{x \in M} f(x) = O\left(\frac{1}{t}\right)$$

$$2) \text{FW-gap } g_t \geq f(x_t) - f^* \rightarrow \text{certificate of subopt.}$$

$$\min_{x \in M} g_t \leq O\left(\frac{1}{t}\right) \quad [\text{i.e. will stop in } O\left(\frac{1}{\epsilon}\right) \text{ iterations.}]$$

$$3) x_t = p_0^t x_0 + \sum_{u=1}^t p_u^t s_{u-1} \quad \rightsquigarrow x_t \text{ has a "sparse" expansion in terms of the FW-vectors } \{s_u\}_{u=1}^{t-1}$$

where $\sum_u p_u^t = 1$
 $p_u^t \geq 0$

⊗⊗ "sparse method"
 → popular in ML

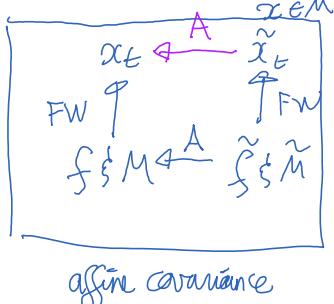
→ see later how to apply on dual of SVM struct

4) there is a $\Omega\left(\frac{1}{t}\right)$ lower bound for FW-like methods for $t \leq d$

5) FW is affine covariant (like Newton's method):
 linear trans \nRightarrow subjective
 Let \tilde{M} be a new constraint set s.t. $\tilde{M} \xrightarrow{A} M$ i.e. $M = A\tilde{M}$

$$\text{define } \tilde{f}(\tilde{x}) \triangleq f(A\tilde{x})$$

$$\min_{\tilde{x} \in \tilde{M}} \tilde{f}(\tilde{x}) = \min_{\tilde{x} \in \tilde{M}} f(A\tilde{x}) = \min_{x \in M} f(x)$$



if run FW on $\tilde{f} \setminus \tilde{M}$ to get \tilde{x}_t as iterates

then $x_t \triangleq A\tilde{x}_t$ corresponds to running FW on $f \setminus M$ (modulo tie breaking)

why? → inner product with gradient!

$$\tilde{s}_t = \operatorname{arg\,min}_{\tilde{s}} \langle \tilde{s}, D_{\tilde{x}} \tilde{f}(\tilde{x}) \rangle \quad x_+ \triangleq A\tilde{x}_+$$

$$\begin{aligned} \tilde{s} \in \tilde{M} & \quad \langle \tilde{s}, A^T \nabla_x f(x_t) \rangle \\ & \quad \langle A\tilde{s}, \nabla_x f(x_t) \rangle \\ s_t = \arg \min_{\substack{s \in M \\ = A\tilde{M}}} & \quad \langle s, \nabla_x f(x_t) \rangle \Rightarrow \underline{s_t} = \tilde{A}\tilde{s}_t \text{ (modulo tie breaking)} \\ & \quad \overbrace{s_t}^{\text{set of firs}} = \tilde{A}\tilde{s}_t \end{aligned}$$

\Rightarrow we want an affine invariant analysis

$$\Rightarrow G_S \leq L_{1/1} (\text{diam}_{1/1}(M))^2 \text{ for any } 1/1$$

affine invariant constant
(we'll see later)

6) convergence for non-convex f

considers necessary first order condition for constrained opt.

$$\min_{x \in M} f(x) \quad x^* \text{ is a local min}$$

$$\Rightarrow \langle \nabla f(x^*), s - x^* \rangle \geq 0 \quad \forall s \in M$$



"stationary pt." for const. opt. problem

$$\Leftrightarrow \min_{s \in M} \langle \nabla f(x^*), s - x^* \rangle \geq 0$$

$$\Leftrightarrow \max_{s \in M} \langle -\nabla f(x^*), s - x^* \rangle \leq 0$$

FW gap(x^*)

quantify the "non-stationarity"

(if $f \in M$ are convex,
then this is a sufficient condition for global min)

See L.-J. 2016 aktiv

$$\min_{s \in M} \text{gap}(x_s) \leq 0 \quad \text{for FW with line search}$$

"non-convex FW"

f is L -smooth and lower bounded
 M bounded $\&$ convex
but f is not nec. convex

15h23

away-step FW

comment: $x_t = \sum_u \mu_u^t s_u$

"coordinate"

$$x_{t+1} = (1-\gamma_t)x_t + \gamma_t s_t$$

$$= \sum_u \mu_u^t (1-\gamma_t) s_u + \gamma_t s_t$$

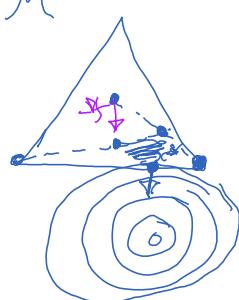
μ^{t+1}

$\mu^{t+1} \rightarrow$ previous coordinates shrunk by $(1-\delta_2)$

* FW step moves mass uniformly away from active set S_t to FW corner s_t

→ unless step-size $\gamma_t=1$, FW never removes a corner from active set / expansion

⇒ zig-zag phenomenon close to boundary on polytopes



(this is why you get $\Omega(1/t)$ rate even if f is μ -strongly convex)

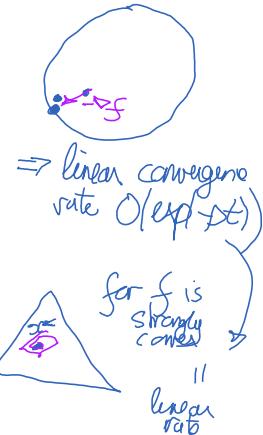
$$\left\langle -\frac{\nabla f(x_t)}{\|\nabla f(x_t)\|}, \frac{d_t}{\|d_t\|} \right\rangle \xrightarrow{t \rightarrow \infty} 0$$

* but FW has no problem

on "strongly convex" sets

↪ sublevel set
of a strongly convex
convex set.

when sol'n is
in the relative
interior of M



away-step FW fix: (solves zig-zagging problem)

in addition to compute FW corner

$$s_t = \underset{s \in M}{\operatorname{argmin}} \langle s, \nabla f(x_t) \rangle$$

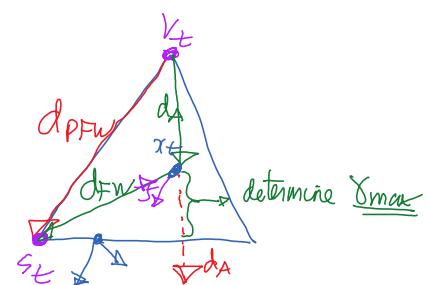
also compute the "away corner"

$$v_t = \underset{s \in \text{activeset}(x_t)}{\operatorname{argmax}} \langle s, \nabla f(x_t) \rangle$$

$$d_{FW} = s_t - x_t$$

 $d_A \triangleq x_t - v_t$

→ (aside: Wolfe's original alg.
Used whole $M \Rightarrow$ non-convergent
alg.)



* AFW picks direction with best inner product with $-\nabla f$

{ i.e. pick d_A if $\underbrace{\langle d_A, \nabla f(x_t) \rangle}_{\parallel Q_A \parallel} > \underbrace{\langle d_{FW}, \nabla f(x_t) \rangle}_{\parallel Q_{FW} \parallel}$

re. pick d_A if $\langle \alpha_A, \nabla f(x_e) \rangle < \langle d_F w, \nabla f(x_e) \rangle$
 o.w. pick $d_F w$

* if use \hat{d}_A , let $\hat{x}_t = \arg \min_{x \in [0, x_{\max}]} f(x_t + \hat{d}_A)$
where x_{\max} depends on β^+ coefficients

also suppose $\alpha_t = \sum_u \alpha_u s_u$
 $x_{t+1} = x_t + \gamma_t (\underbrace{x_t - v_t}_{d_A})$

$$= \sum_u (1+\gamma_t) q_{uy} s_u - \gamma_t v_t$$

let α be coeff. of v_6 in expansion for x_6

$$(1+\gamma)\alpha - \gamma \geq 0$$

$$(1 + \gamma_{\max}) \alpha - \gamma_{\max} = C$$

$$\Rightarrow \gamma_{\max} = \frac{q}{1-\alpha}$$

* when $\gamma_t = \gamma_{\max}$

(call this a "drop step") \rightarrow removing $\forall c$ from expansion

(*) when run AFW alg's

either you maintain some expansion $x_t = \sum_u p_u s_u$

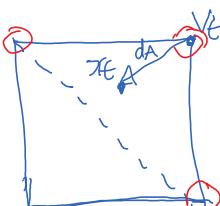
or
you have a feasibility oracle + away-step oracle
get δ_{\max} get ν

[see NIPS 2016 paper by Meshi & Cohen]

↳ they assume $\text{corners}(u) \subseteq \{0, 1\}^d$

\rightarrow assumption needed for convergence result

and M is described as $Ax=b$
 $x \geq 0$



AFW has linear convergence rate on polytopes when f is strongly convex
 (vs. FW $\rightarrow \mathcal{O}(\frac{1}{\epsilon})$)

* Combine ∂_{FW} & $d_A \rightarrow$ pairwise FW direction

$$d_{FW} = d_{FW} + d_A = s_t - \cancel{x_6} + \cancel{x_4} - v_t = s_t - v_t$$

$\cancel{x_6}$
 $\cancel{x_4}$
 v_t

$$\underbrace{\langle -\nabla f(x_t), d_{FW} \rangle}_{g_{FW}} = g_{FW} + g_A$$

$$g_{FW} + g_A \leq 2 \max\{ \dots \}$$

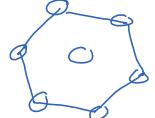
$$\text{during AFW: } g_t = \langle -\nabla f(x_t), d_t \rangle = \max \{ g_{FW}, g_A \}$$

$$g_t \geq \frac{g_{FW}}{2}$$

* note: if $M = \text{conv}(A)$ where A is some finite set (called "atoms")

$$\text{LMO}(r) : \min_{S \in M = \text{conv}(A)} \langle S, r \rangle = \min_{a \in A} \langle a, r \rangle$$

↳ lot of application
in ML



where LMO is efficient

e.g.: $f_t \hookrightarrow$ integer flows

$\text{conv}(A) \hookrightarrow$ flow polytope

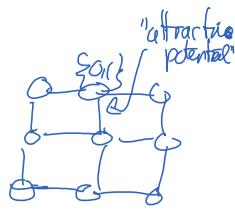
LMO \rightarrow min cost network
algorithm

$A \rightarrow$ clique assignments
in graph $(S_{yc})_{c \in \zeta}$

LMO \rightarrow max product
alg.

$\text{conv}(A) \rightarrow$ marginal
polytope

or



Graph cut alg.
for Ising model with attractive potential
("Associative Markov network")
(\rightarrow submodular potentials)
(see later)