

today: • properties of FW
• AFW alg.

properties of FW

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

1) $f(x_t) - \min_{x \in M} f(x) \triangleq f^* \leq O(\frac{1}{t})$

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

$$\min_{s \in S} g_s \leq O(\frac{1}{t})$$

3) $x_t = P_0^t x_0 + \sum_{u=1}^t P_u^t s_{u-1}$

$\leadsto x_t$ has a "sparse" expansion in terms of the FW corners $\sum_{u=1}^{t-1} s_u$

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

"sparse method" \rightarrow popular in ML

\rightarrow see later how to apply on dual of SVM struct

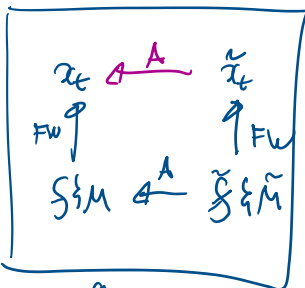
4) there is a $\Omega(\frac{1}{t})$ lower bound for FW-like methods for $t \leq d$

5) FW is affine covariant (like Newton's method):

let \tilde{M} be a new constraint set s.t. $\tilde{M} \xrightarrow{A} M$ ie. $M = A\tilde{M}$

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

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



affine covariance

If run FW on \tilde{f} & \tilde{M} to get \tilde{x}_t as iterates

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

why? \rightarrow inner product with gradient is invariant affine

$$\tilde{s}_t = \arg \min_{\tilde{s} \in \tilde{M}} \langle \tilde{s}, \nabla_{\tilde{x}} \tilde{f}(\tilde{x}) \rangle$$

$$x_t \triangleq A\tilde{x}_t$$

$$\langle \tilde{s}, A^T \nabla_x f(x_t) \rangle$$

$$\langle A\tilde{s}, \nabla_x f(x_t) \rangle$$

$$\Rightarrow \underline{\underline{s_t = A\tilde{s}_t}}$$

(modulo tie breaking)

$$\langle A_s^T, \nabla_x f(x_t) \rangle \quad \rightarrow \quad \text{breaking}$$

$$s = A_s^T$$

$$x_t = \operatorname{argmin}_{s \in M = A\tilde{M}} \langle s, \nabla_x f(x_t) \rangle$$

$$\sum_{\tilde{M}} = A \tilde{S}_t$$

set of ties

\Rightarrow we want an fine invariant analysis

$$C_S \leq L_{\|\cdot\|} (\operatorname{diam}_{\|\cdot\|}(M))^2 \text{ for any } \|\cdot\|$$

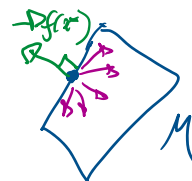
fine invariant constant (we'll see later)

6) convergence of FW for non-convex

aside: necessary first order condition for constrained opt. $\min_{x \in M} f(x)$

$$x^* \text{ is a local min } \Rightarrow \langle \nabla f(x^*), s - x^* \rangle \geq 0 \quad \forall s \in M$$

feasible direction



"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 'non-stationarity'

(if M & f are concave, then this is a sufficient cond. for global min)

see L.-J. 2016 arxiv

$$\min_{s \in T} \operatorname{gap}(x_s) \leq O\left(\frac{1}{J_t}\right) \text{ for FW with line search}$$

"non-convex FW"

f is L -smooth
 M is bounded & convex
 but f is not rec-convex

lower bound of f on M

$$O\left(\frac{1}{J_t}\right) \text{ if } f \text{ is concave}$$

10/19

away step FW

comment:

$$x_t = \sum_u p_u^t s_u$$

"coordinate"

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

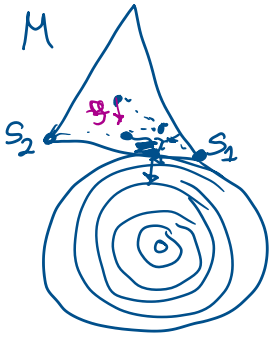
$$= \sum_u \underbrace{p_u^t (1-\delta_t)}_{p_u^{t+1}} s_u + \delta_t s_t$$

$p_u^{t+1} \Rightarrow$ previous coordinates shrank by $(1-\delta_t)$

* FW step moves mass uniformly away from active set $\{s_u \mid \delta_u = 1\}$ to FW corner s_t

→ unless step size $\delta_t = 1$, FW never removes a corner from active set / expansions
 ⇒ zig-zag phenomenon close to boundary on polytopes

(this is why you get $O(\frac{1}{\epsilon})$ rate even if f is μ -strongly convex)



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

⇒ no linear convergence rate for FW
 → sublevel set of a strongly convex f .

* but FW has no problem

on "strongly convex sets"

⇒ linear convergence rate ($O(\exp(-\mu t))$) when f is strongly convex

when sol'n is the relative interior of M



Away-Step FW f_x : (solves zig-zagging problem)

in addition to compute FW corner

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

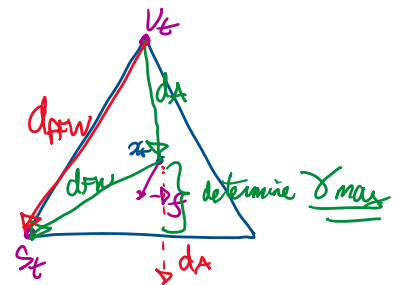
also compute the "away corner"

$$v_t \triangleq \underset{S \in \text{active set}(x_t)}{\text{argmax}} \langle s, \nabla f(x_t) \rangle$$

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

$$d_A \triangleq x_t - v_t$$

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



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

$$\left\{ \begin{array}{l} \text{ie. pick } d_A \text{ if } \underbrace{\langle d_A, -\nabla f(x_t) \rangle}_{\text{"g}_A"} > \underbrace{\langle d_{FW}, -\nabla f(x_t) \rangle}_{\text{g}_{FW}} \\ \text{o.w. pick } d_{FW} \end{array} \right.$$

* if use d_A , let $x_t = \underset{\delta \in [0, \delta_{\max}]}{\text{argmin}} f(x_t + \delta d_A)$

where δ_{\max} depends on μ coefficients

$$\text{ie. suppose } x_t = \sum_u \alpha_u s_u$$

d_A

i.e. suppose $x_t = \sum_u \alpha_u s_u$

$$x_{t+1} = x_t + \underbrace{\gamma_t}_{d_t} (x_t - v_t)$$

$$= \sum_u (1 + \gamma_t) \alpha_u s_u - \gamma_t v_t$$

let α be the coeff. of v_t in expansion of x_t

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

$$(1 + \gamma_{\max})\alpha - \gamma_{\max} = 0 \Rightarrow \gamma_{\max} = \frac{\alpha}{1 - \alpha}$$

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

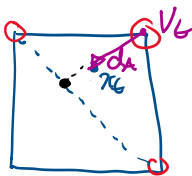
* when $\gamma_t = \gamma_{\max}$

call this a "drop step" \rightarrow removing v_t from expansion

⊕ when run AFW alg.:

either you maintain some expansion $x_t = \sum_u \alpha_u s_u$ } this is "slower" by being conservative

or you have a feasibility oracle + away-step oracle
 γ_{\max} get v_t



[see NeurIPS 2016 paper by Meshi & Garber]

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

\rightarrow assumption needed for convergence result

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

assumption to run alg.

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

⊗ combine d_{FW} & $d_A \rightarrow$ pairwise FW direction

$$d_{AFW} = d_{FW} + d_A = \cancel{\gamma_t} - \cancel{\gamma_t} + \gamma_t - v_t = \underbrace{\gamma_t - v_t}_{d_{AFW}}$$

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

$$g_{AFW} = g_{FW} + g_A \leq 2 \max\{g_{FW}, g_A\}$$

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

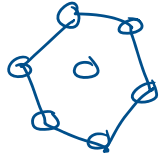
$$g_t \geq \frac{g_{AFW}}{2}$$

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

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

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

↳ lots of applications in ML where LMO is efficient

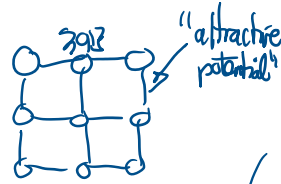


e.g. • $A \rightarrow$ integer flows
 $\text{conv}(A) \rightarrow$ flow polytope

LMO \rightarrow min cost network flow alg.

• $A \rightarrow$ degree assignments on a graph $(\delta_{yc})_{c \in E}$

$\text{conv}(A) \rightarrow$ marginal polytope LMO \rightarrow max product alg.



or graph cut. alg.
 For Ising model with attractive potentials
 (\rightarrow "submodular potentials" (see later))