

Lecture 17 - FW convergence

Tuesday, March 21, 2023 2:30 PM

today: convergence of FW alg.
apply sumstruct obj.

⊗ 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 applications
in ML
where LMO is efficient

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

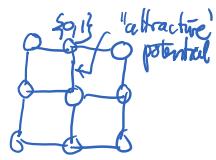
LMO \rightarrow min cost network flow algorithm

$\cdot A \rightarrow \text{clique assignment}$ in graph $(S_{\text{yc}})_{\text{c} \in \mathcal{C}}$

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

LMO \rightarrow max product alg.

or



Graph cut alg. for binary Ising model
with attractive potential
("Associative Markov network")

(\rightarrow "submodular potentials")
(see later)

curvature constant C_f

curvature constant $C_f \triangleq \sup_{\substack{\gamma \in [0, 1] \\ x, s \in M \\ x_\gamma = (1-\gamma)x + \gamma s}} \left[f(x_\gamma) - \left(f(x) + \langle \nabla f(x), x_\gamma - x \rangle \right) \right]$

This is affine invariant

potential FW step update
worst case deviation from linear approximation

C_f is affine invariant

* by descent lemma, if ∇f is L -Lipschitz

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2$$

$$\begin{aligned} & [\| \nabla f(x) - \nabla f(y) \|_* \leq L \|x - y\|] \\ & [\langle d, x \rangle] \leq \|d\|_* \|x\| \end{aligned}$$

$$\Rightarrow C_f \leq \sup_{\gamma} \frac{2}{\gamma^2} \frac{L}{2} \|x_\gamma - x\|^2 \quad \downarrow x + \gamma(s-x)$$

$$\frac{L}{2} \frac{2}{\gamma^2} \|s - x\|^2$$

$$C_f \leq L \cdot \sup_{x, s \in M} \|s - x\|^2$$

$$\text{claim}_{\|\cdot\|}(M) \triangleq \sup_{x \in M} \|x - s\|$$

$C_f \leq L \|\cdot\| \text{claim}_{\|\cdot\|}(M)^2$

depends on $\|\cdot\|$

for any $\|\cdot\|$

affine invariant

④ by def of C_f , we get affine invariant version of descent lemma

$f(x_\gamma) \leq f(x) + \gamma \langle Df(x), s - x \rangle + \frac{\gamma^2}{2} C_f \quad \forall \gamma \in [0, 1] \quad \forall x, s \in M$

let $x = x_t$ and $s = s_t$, FW corner $\langle \nabla f(x_t), s_t - x_t \rangle = -g_t$ FW gap

for FW step of size γ

$f(x_\gamma) \leq f(x_t) - \gamma g_t + \frac{\gamma^2}{2} C_f \quad \forall \gamma \in [0, 1]$

optimize step size for bound (KHS)

$\gamma^* = \min \left\{ \frac{g_t}{C_f}, 1 \right\}$

this gives you an affine-invariant
adaptive step size

$$f(x_{\gamma^*}) \leq f(x_t) - \frac{g_t^2}{2C_f} \quad [\text{when } \frac{g_t}{C_f} \leq 1]$$

$$\text{let } \epsilon_t \triangleq f(x_t) - f^* \leq g_t$$

$$\leq f(x_t) - \frac{\epsilon_t^2}{2C_f}$$

$$\epsilon_{t+1} \leq \epsilon_t - \frac{\epsilon_t^2}{2C_f}$$

thm.: FW alg. with γ_t chosen either as

- (when f is convex)
- yields $\epsilon_t \leq \frac{2C_f}{t+2}$

$\frac{\partial f}{\partial t} = g_t / C_f$
line search

rate:

- non-convex f
 $\min_s g_s \leq O\left(\frac{C_f}{\sqrt{t}}\right)$

f is concave, $C_f = 0$?

$$\min_s g_s \leq O\left(\frac{1}{t}\right)$$

prof: let $x_\gamma = x_t + \gamma(g_t - x_t)$ & apply (+)

$$f(x_\gamma) \leq f(x_t) - \gamma g_t + \frac{\gamma^2}{2} C_f \quad \forall \gamma \in [0, 1]$$

To bound correctly, $g_t \geq \varepsilon_t$
 $g_t \leq -\varepsilon_t$

$$\underbrace{f(x_{t+1}) - f^*}_{\varepsilon_{t+1}} \leq \underbrace{f(x_t) - f^*}_{\varepsilon_t} - \gamma_t \varepsilon_t + \frac{\gamma_t^2 C_f}{2}$$

$$\boxed{\varepsilon_{t+1} \leq (1-\gamma_t) \varepsilon_t + \frac{\gamma_t^2 C_f}{2}}$$

* see notes 2017 for a cool ODE trick + induction
 here, brute force approach to solve recurrence

$$\begin{aligned} \varepsilon_{t+1} &\leq (1-\gamma_t) \varepsilon_t + \frac{\gamma_t^2 C_f}{2} \\ &\leq (1-\gamma_t) [(1-\gamma_{t-1}) \varepsilon_{t-1} + \frac{\gamma_{t-1}^2 C_f}{2}] + \frac{\gamma_t^2 C_f}{2} \end{aligned}$$

$$\boxed{\varepsilon_{t+1} \leq \sum_{s=0}^t (1-\gamma_s) \varepsilon_0 + \frac{C_f}{2} \sum_{s=0}^t \gamma_s^2 \left(\prod_{u=s+1}^t (1-\gamma_u) \right)}$$

initial condition Lipschitz constant

$$\text{use } (1+\gamma) \leq \exp(\gamma) \quad \forall \gamma$$

$$(1-\gamma) \leq \exp(-\gamma)$$

loose?

$$\Rightarrow \varepsilon_{t+1} \leq \varepsilon_0 \exp\left(-\sum_{s=0}^t \gamma_s\right) + \frac{C_f}{2} \sum_{s=0}^t \gamma_s^2 \exp\left(-\sum_{u=s+1}^t \gamma_u\right)$$

$$\gamma_s \approx \frac{1}{s} \Rightarrow \sum_{s=0}^t \gamma_s \approx \log t$$

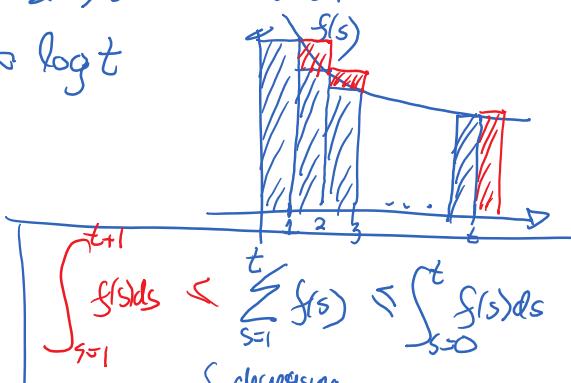
15h37

$$\exp\left(-\sum_{s=0}^t \gamma_s\right) \geq \exp(-\log t)$$

$$\exp(\log \frac{1}{t}) \approx O\left(\frac{1}{t}\right)$$

$$\exp\left(-\sum_{u=s+1}^t \gamma_u\right) \approx \exp(-\log t/s) \approx O\left(\frac{1}{s}\right)$$

$$\sum_{s=0}^t \frac{\gamma_s^2 \exp\left(-\sum_{u=s+1}^t \gamma_u\right)}{s} \approx \frac{1}{t} \sum_{s=0}^t \frac{1}{s} \rightarrow \frac{\log t}{t}$$



$$\begin{aligned} \sum_{s=1}^t \frac{1}{s} &= 1 + \sum_{s=2}^t \frac{1}{s} \leq 1 + \int_{s=1}^t \frac{1}{s} ds \\ &= 1 + [\log s]_1^t \\ &= 1 + \log t \end{aligned}$$

④ in fact if we $\gamma_t = \frac{1}{t+1}$, you do get $O\left(\frac{\log t}{t}\right)$ rate

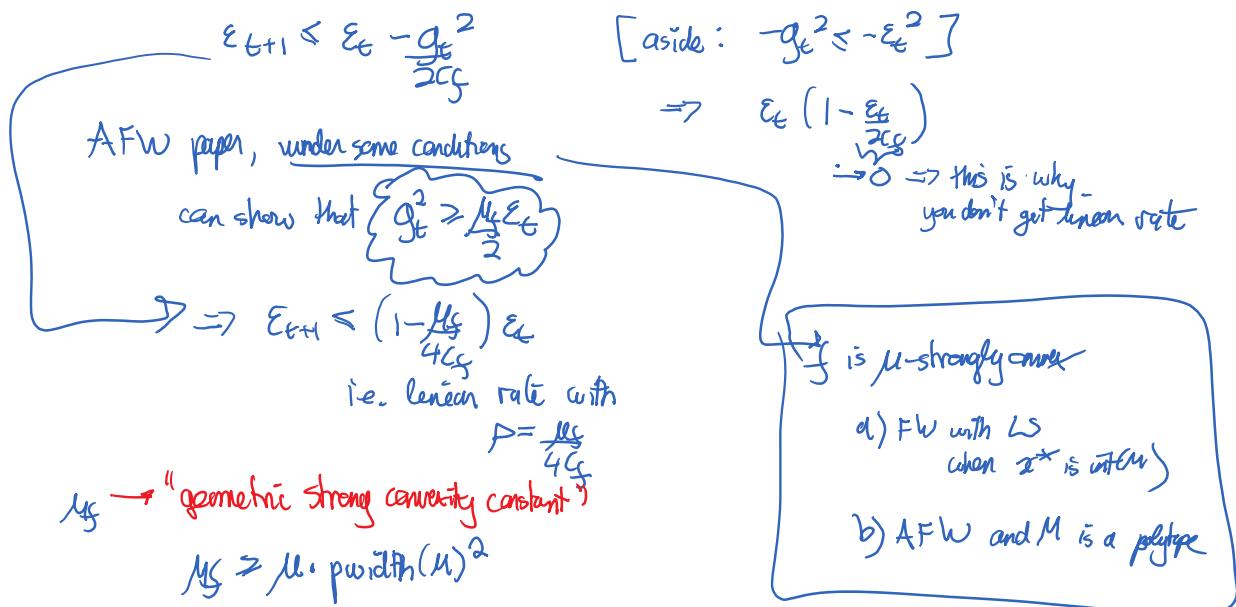
but if $\gamma_t = \frac{2}{t+2}$, here the bound says $O\left(\frac{\log t}{t}\right)$
 but a (tighter) direct analysis yields $O\left(\frac{1}{t}\right)$
 see notes in 2017, for $\gamma_t = \frac{\alpha}{t+\alpha}$ ($O\left(\frac{1}{t}\right)$ for $\alpha \gg 2$)
 [telescoping product]

Lecture 11-- 2017/2/20 -- http://www.iro.umontreal.ca/~slacoste/teaching/ift6085/W17/protected/notes/lecture11_scribbles.pdf

Linear rate of AFW

"linear rate" constant $(1-p) \leq \exp(-pt)$
 linear rate : $\varepsilon_{t+1} \leq (1-p)\varepsilon_t \leq (1-p)^t\varepsilon_0 \leq \varepsilon_0 \exp(-pt)$
 (terminology)
 for gradient descent $p = \frac{\mu}{L} = \frac{1}{K} \leftarrow \text{condition\#}$
 sublinear rate $\varepsilon_t \leq O\left(\frac{1}{t} \text{ some power}\right)$

recall for FW with line search:



linear rate : $p = \frac{\mu_g}{4C_f} \geq \frac{1}{4} \frac{\mu \cdot \text{width}(M)^2}{\text{claim}(M)^2}$

$\frac{1}{K_f} \quad \frac{1}{\text{claim}(M)^2} \quad \text{"condition\# of } M\text{"}$

FW for SVM struct dual

dual of SVM struct : $\min_{\alpha_i \in \Delta_{\{y_i\}}} \frac{1}{2} \|A\alpha\|^2 - b^T \alpha$

(i.e. $M = \bigcup_{y_i} \Delta_{\{y_i\}}$) $A\alpha = \sum_{i=1}^n \alpha_i y_i$ $w(\alpha) = w(\alpha)$

let $\alpha^{(0)} = S_{y^{(1)}}$ $\Rightarrow w(\alpha^{(0)}) = 0$

$w_t \triangleq w(\alpha_t)$

FW steps :

$$s_t = \underset{s \in \mathcal{M}}{\operatorname{argmin}} \langle s, \nabla f(\alpha_t) \rangle \quad \nabla f(\alpha_t) = A A^T \alpha_t - b$$

$$(\nabla f(\alpha_t))_{i,\tilde{y}} = \frac{A^T A \alpha_t}{n} - \frac{b}{n}$$

$$= -\frac{1}{n} [l_i(\tilde{y}) - w^T y_i(\tilde{y})]$$

$$\min_{S \in \mathcal{M}} \langle s, \nabla f(\alpha_t) \rangle = \min_{\{s_i \in \mathcal{M}_i\}} \sum_{i=1}^n \langle s_i, \nabla f(\alpha_t) \rangle$$

$$= \sum_i \min_{S_i \in \mathcal{M}_i} \langle s_i, \nabla_i f(\alpha_t) \rangle \quad (\text{block-separable structure})$$

$$\alpha_t = \sum_u \alpha_{u,t} s_u$$

$$M_i = \Delta_{\{y_i\}}$$

$$\min_{\tilde{y}} \langle \delta_{\tilde{y}}, \nabla_i f(\alpha_t) \rangle$$

$$\nabla_i y f(\alpha_t) = -\frac{1}{n} H_i(\tilde{y}; w_t)$$

thus $\hat{s}_t \triangleq (\hat{s}_i)_{i=1}^n$ where $\hat{s}_i = \delta_{\tilde{y}_i}(w_t)$ where $y_i(w_t) = \operatorname{argmax}_{\tilde{y} \in S_i} H_i(\tilde{y}; w_t)$

$$\hat{s}_i(\tilde{y}) = \mathbb{1}\{\tilde{y} = \hat{y}_i\}$$



loss-augmented decoding?

$$\alpha^{(t)} \xrightarrow{A} w^{(t)}$$

$$\alpha_i^{(t+1)} = (1-\gamma) \alpha_i^{(t)} + \gamma \hat{s}_i^{(t)}$$

[here, need to maintain active set $\{S_i^{(t)}\}_{i=0}^n$]

$$\alpha^{(t+1)} = (1-\gamma) \alpha^{(t)} + \gamma s^{(t)}$$

$$w^{(t+1)} = (1-\gamma) A \alpha^{(t)} + \gamma A s^{(t)}$$

you can choose via analytic LSS on dual obj.

recall primal obj:

$$p(w) = \frac{1}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n H_i(w) \xrightarrow{\text{max } \tilde{y}} \max_{\tilde{y}} l_i(\tilde{y}) - w^T y_i(\tilde{y})$$

$$p'(w_t) = A w_t - \frac{1}{n} \sum_{i=1}^n y_i(\tilde{y}_i^{(t)})$$

$$\begin{aligned} w^{(t+1)} &= w^{(t)} - \beta p'(w_t) \\ &= (1 - \beta \lambda) w^{(t)} + \beta \sum_{i=1}^n y_i(\tilde{y}_i^{(t)}) \end{aligned}$$

$$\text{if we set } \boxed{\beta = \frac{\lambda}{\lambda}}$$

Then batch subgradient step on primal
is equivalent
to a batch FW step on dual
with $\beta = \frac{\lambda}{\lambda}$ step-size relationship

$$w^{(t)} = A \alpha^{(t)}$$

FW perspective gives you "adaptive step-size"

for batch subgradient method: