

Lecture 17 - FW convergence

Tuesday, March 17, 2020 14:33

- today:
- convergence FW alg.
 - apply SVM struct obj.

curvature constant C_f

$$\text{Curvature constant } C_f \triangleq \sup_{\substack{\delta \in (0,1] \\ x, s \in M}} \frac{1}{\delta^2} [f(x_\delta) - (f(x) + \langle \nabla f(x), x_\delta - x \rangle)]$$

$x_\delta = (1-\delta)x + \delta s$
 ↖ potential FW update

↘ this affine invariant

worst case deviation from linear approximation

⇒ C_f is affine invariant

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

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

$$[\|\nabla f(x) - \nabla f(y)\|_* \leq L \|x-y\| \quad \forall x, y \in M]$$

$$[\langle d, x \rangle \leq \|d\|_* \|x\|]$$

$$\Rightarrow C_f \leq \sup \frac{1}{\delta^2} \left[\frac{L}{2} \|x_\delta - x\|^2 \right]$$

$x_\delta = x + \delta(s-x)$
 $\| \delta(s-x) \|^2$

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

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

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

↖ affine invariant

↖ depends on $\|\cdot\|$

⊕ by def. of C_f , we get an affine invariant version of descent lemma

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

let $x = x_k$ and $s = s_k$, FW corner

FW-gap
+

$$\langle \nabla f(x_t), s_t - x_t \rangle = -g_t$$

for FW step of size γ

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

optimize step-size for bound (RHS)

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

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

↑ this gives you an affine mix adaptive step-size

$$\begin{aligned} \epsilon_t &\stackrel{\Delta}{=} f(x_t) - f(x^*) \leq g_t \\ &\leq f(x_t) - \frac{\epsilon_t^2}{2C_f} \end{aligned}$$

thm: FW alg. with γ_t chosen either as $\frac{2}{t+2}$ (when f is convex) yields $\epsilon_t \leq \frac{2C_f}{t+2}$

↙ $\frac{g_t}{C_f}$ line search

note: non-convex f

$$\min_{s \leq t} g_s \leq O\left(\frac{1}{\sqrt{t}}\right)$$

concave f , $C_f = 0$

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

proof: let $x_\gamma = x_t + \gamma(s_t - x_t)$ + apply (+)

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

by convexity, $g_t \geq \epsilon_t$

$$\underbrace{f(x_{t+1}) - f^*}_{\epsilon_{t+1}} \leq \underbrace{f(x_t) - f^*}_{\epsilon_t} - \gamma \epsilon_t + \frac{\gamma^2}{2} C_f$$

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

* see notes 2017 for a cool OPE trick + induction

here, write force approach to solve recurrence

$$\begin{aligned} \epsilon_{t+1} &\leq (1-\alpha_t) \epsilon_t + \frac{\alpha_t^2}{2} C_f \\ &\leq (1-\alpha_t) \left[(1-\alpha_{t-1}) \epsilon_{t-1} + \frac{\alpha_{t-1}^2}{2} C_f \right] + \frac{\alpha_t^2}{2} C_f \end{aligned}$$

$$\epsilon_{t+1} \leq \prod_{s=0}^t (1-\alpha_s) \epsilon_0 + \frac{C_f}{2} \sum_{s=0}^t \alpha_s^2 \left(\prod_{u=s+1}^t (1-\alpha_u) \right)$$

↑ initial condition
↑ Lipschitz constant

use $(1+\delta) \leq e^\delta \quad \forall \delta$

$(1-\delta) \leq e^{-\delta}$
↑ loose?

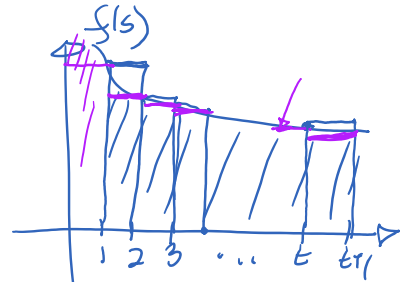
$$\Rightarrow \epsilon_{t+1} \leq \epsilon_0 \exp\left(-\sum_{s=0}^t \alpha_s\right) + \frac{C_f}{2} \sum_{s=0}^t \alpha_s^2 \exp\left(-\sum_{u=s+1}^t \alpha_u\right)$$

$\alpha_s \sim \frac{1}{s} \Rightarrow \sum_{s=0}^t \alpha_s \sim \log(t)$

$\exp\left(-\sum_{s=0}^t \alpha_s\right) \sim \exp(-\log(t)) = \frac{1}{t}$

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

$\sum_{s=0}^t \alpha_s^2 \exp\left(-\sum_{u=s+1}^t \alpha_u\right) \sim O\left(\frac{\log t}{t}\right)$



$$\int_{s=0}^t f(s) ds \geq \sum_{s=1}^t f(s) \geq \int_{s=1}^{t+1} f(s) ds$$

(f decreasing)

⊛ in fact, if use $\alpha_t = \frac{1}{t+1}$, you do get $O\left(\frac{\log t}{t}\right)$ rate

but $\alpha_t = \frac{2}{t+2}$, here our bound says $O\left(\frac{\log t}{t}\right)$; but (tighter) analysis is $O\left(\frac{1}{t}\right)$

see notes in 2017, for $\alpha_t = \frac{\alpha}{t+\alpha}$ ($O\left(\frac{1}{t}\right)$ for $\alpha \geq 2$)

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

15h40

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

Linear rate for AFW:

"linear rate constant"

linear rate: $\epsilon_{t+1} \leq (1-\rho) \epsilon_t \leq \epsilon_0 (1-\rho)^t \leq \epsilon_0 \exp(-\rho t)$

(for gradient descent $\rho = \frac{\mu}{L}$)

sublinear rate: $\epsilon_t \leq O\left(\frac{1}{t^{\text{power}}}\right)$

recall for FW (with LS) : $\epsilon_{t+1} \leq \epsilon_t - \frac{Q_t^2}{2C_f}$ [aside $-Q_t^2 \leq -\epsilon_t^2$]

$\epsilon_t (1 - \frac{M_f}{2C_f})$

AFW paper, under some conditions

f is μ -strongly convex

can show $Q_t^2 \geq \frac{M_f}{2} \epsilon_t$

a) FW with LS when x^* is int(M)

b) AFW and M is a polytope

$\Rightarrow \epsilon_{t+1} \leq (1 - \frac{M_f}{4C_f}) \epsilon_t$

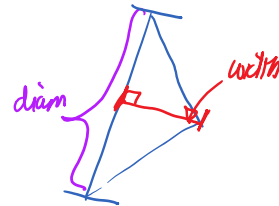
i.e. linear rate with $\rho = M_f/4C_f$

$M_f \rightarrow$ geometric strong convexity constant

$M_f \geq \mu \cdot \text{width}(M)^2$

linear rate $\rho = \frac{M_f}{4C_f} \geq \frac{\mu \cdot \text{width}(M)^2}{4L \cdot \text{diam}(M)^2}$

\downarrow $\frac{1}{K_f}$ \downarrow "condition # of set M"



FW for SVMstruct

dual of SVMstruct: $\min_{\alpha_i \in \Delta_{|S_i|}} \lambda \frac{\|A\alpha\|^2 - b^T \alpha}{2}$

i.e. $M = \prod_{i=1}^n \Delta_{|S_i|}$

$A\alpha = \int_{\Delta_n} \sum_{i=1}^n \sum_y \alpha_i(y) \psi_i(y) = w(\alpha)$

let $\alpha_i^{(0)} = \delta_{y^{(i)}} \Rightarrow w(\alpha^{(0)}) = 0$

FW step:

$x_t = \text{argmin}_{S \in M} \langle S, \nabla f(\alpha_t) \rangle$

$\nabla f(\alpha_t) = \lambda A^T A \alpha_t - b$ $w_t = A \alpha_t$
 $w_t \triangleq w(\alpha_t)$

$(\nabla f(\alpha_t))_{i,y} = \lambda \frac{\psi_i(y)^T w_t}{x_n} - \frac{b_i(y)}{n}$
 $= -\frac{1}{n} H_i(y; w_t)$

$\min_{S \in M} \langle S, \nabla f(\alpha_t) \rangle = \min_{\{S_i \in M_i\}} \sum_i \langle S_i, \nabla_i f(\alpha_t) \rangle$
 $= \sum_i \underbrace{\min_{\{S_i \in M_i\}} \langle S_i, \nabla_i f(\alpha_t) \rangle}_{}$

$$= \sum_i \min_{\{s_i \in M_i\}} \langle s_i, \nabla f(x_t) \rangle$$

$$M_i = \Delta_{\{y_i\}} \quad \min_y \langle \delta y, \nabla f(x_t) \rangle$$

$$\nabla_{y_i} f(x_t) = -\frac{1}{n} H_i(y; w_t)$$

thus $s_t \triangleq (\hat{s}_i)_{i=1}^n$ where $\hat{s}_i = \delta \hat{y}_i(w_t)$ where $\hat{y}_i(w_t) \triangleq \arg \max_{y \in \mathcal{Y}_i} H_i(y; w_t)$

$$\hat{s}_i(y) = \mathbb{1}_{\{y = \hat{y}_i(w_t)\}} \quad \text{[loss-augmented decoding]}$$

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

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

$$w^{(t+1)} = (1-\delta) \frac{A \alpha^{(t)}}{w^{(t)}} + \delta \frac{A s^{(t)}}{\sum_{i=1}^n \psi_i(\hat{y}_i^{(t)})}$$

[here, need to maintain active set $\{ \delta \hat{y}_i^{(t)} \}_{i=1}^n$]

you can choose via analytic LS. on dual objective

$$\gamma_t^* = \arg \min_{\gamma \in [0,1]} f(\alpha^{(t)} + \gamma(s^{(t)} - \alpha^{(t)}))$$

recall primal obj.

$$p(w) = \frac{\lambda \|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n H_i(w)$$

$$p'(w_t) = \lambda w^{(t)} - \frac{1}{n} \sum_{i=1}^n \psi_i(\hat{y}_i^{(t)})$$

$$\max_{\hat{y}} \ell(\hat{y}) - w^T \psi_i(\hat{y})$$

$$w^{(t+1)} = w^{(t)} - \beta p'(w_t)$$

$$= (1-\beta) w^{(t)} + \frac{\beta}{n} \sum_{i=1}^n \psi_i(\hat{y}_i^{(t)})$$

If set $\beta = \frac{\lambda}{\lambda}$

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

FW perspective gives you "adaptive step-size" for batch subgradient