

Lecture 17 - FW convergence

Tuesday, March 12, 2019 14:19

today: • convergence of FW
• apply FW to sumstruct

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

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

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

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

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

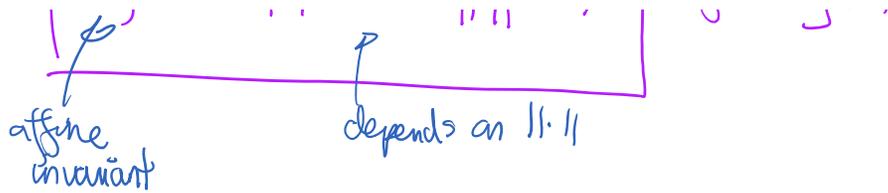
$$\Rightarrow C_f \leq \sup_{\substack{\delta \\ x, s \in M}} \frac{2}{\delta^2} \left[\frac{L}{2} \|x_\delta - x\|^2 \right]$$

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

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

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

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



⊛ by def. of C_f , we get an affine invariant version of descent Lemma

$$f(x_s) \leq f(x) + \gamma \langle \nabla f(x), s-x \rangle + \frac{\gamma^2}{2} C_f \quad \begin{array}{l} \forall \gamma \in [0, 1] \\ \forall x, s \in M \end{array}$$

let $x = x_t$ and $s = s_t$, FW corner

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

for FW-step of size γ

$$(+) \quad f(x_s) \leq f(x_t) - \gamma g_t + \frac{\gamma^2}{2} C_f$$

Optimizing step-size for bound (RHS)

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

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

this gives an affine inv. adaptive step-size

$$\epsilon_t \triangleq f(x_t) - f(x^*) \leq g_t$$

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

thm: FW alg. with α_t chosen either $\frac{2}{t+2}$, $\frac{g_t}{C_f}$ or line search
 (when f is convex)
 yields $E_t \leq \frac{2C_f}{t+2}$

note:

non-convex f

$$\min_{S \subseteq T} g_S \leq O\left(\frac{1}{\sqrt{t}}\right)$$

concave f , $C_f = 0$

$$\min_{S \subseteq T} g_S \leq O\left(\frac{1}{t}\right)$$

proof: let $x_\gamma = x_t + \gamma(x_t - x_t) + \text{apply } (f)$

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

by convexity, $g_t \geq E_t$

$$\underbrace{f(x_{\alpha_t}) - f^*}_{E_{t+1}} \leq \underbrace{f(x_t) - f^*}_{E_t} - \alpha_t E_t + \frac{\alpha_t^2}{2} C_f$$

$$E_{t+1} \leq (1 - \alpha_t) E_t + \frac{\alpha_t^2}{2} C_f$$

* see notes 2017 for a cool GDF trick + induction

here, brute force approach to solve the recurrence

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

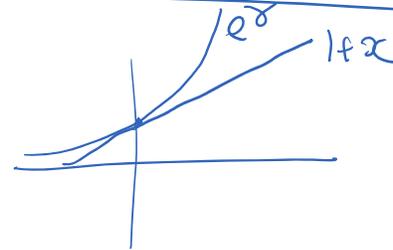
$$\approx (1-\gamma_t) \left[(1-\gamma_{t-1}) \varepsilon_{t-1} + \frac{\gamma_{t-1}^2}{2} C_f \right] + \frac{\gamma_t^2}{2} C_f$$

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

initial condition Lipschitz part

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

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

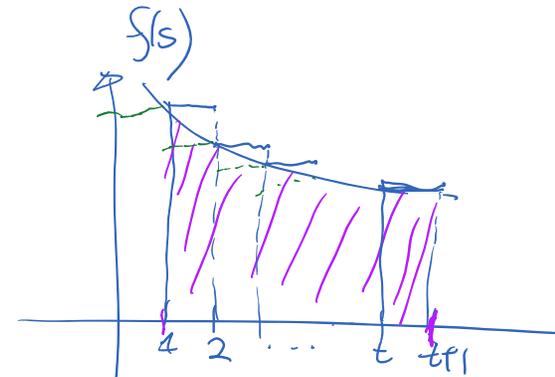


$$\Rightarrow \varepsilon_{t+1} \leq \varepsilon_0 \exp\left(-\sum_{s=0}^t \gamma_s\right) + \frac{C_f}{2} \sum_{s=0}^t \gamma_s^2 \exp\left(-\sum_{u=s+1}^t \gamma_u\right)$$

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

$\exp\left(-\sum_{s=0}^t \gamma_s\right) \approx \exp(-\log t) = O\left(\frac{1}{t}\right)$

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



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

$$\sum_{s=1}^t \frac{1}{s} = 1 + \sum_{s=2}^t \frac{1}{s} \leq 1 + \int_{s=1}^t \frac{1}{s} ds$$

$$\sum_{s=1}^t \frac{\gamma_s^2}{s} \exp(-\sum_{r=1}^s \gamma_r) \approx O\left(\frac{\log t}{t}\right)$$

$\frac{1}{s^2} \approx \frac{\gamma}{t}$

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

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

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

see notes in 2012, for $\gamma_t = \frac{\alpha}{t+\alpha}$ ($O\left(\frac{1}{t}\right)$ for $\alpha \geq 2$)

15h36

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

Linear rate for AFW:

"linear rate constant"

linear rate: $\varepsilon_{t+1} \leq (1-\rho) \varepsilon_t \leq \varepsilon_0 (1-\rho)^t \leq \varepsilon_0 \exp(-\rho t)$
 (for gradient descent $\rho = \frac{\mu}{L}$)

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

recall for FW (with LS): $\varepsilon_{t+1} \leq \varepsilon_t - \frac{g_t^2}{2\zeta_t}$

[aside $-g_t^2 \leq \varepsilon_t^2$]
 $\varepsilon_t \left(1 - \frac{\varepsilon_t}{2\zeta_t}\right)$

AFW paper, under some conditions

can show $\sigma_c^2 \geq \frac{\mu_g}{2} \epsilon_t$

$$\Rightarrow \epsilon_{t+1} \leq \left(1 - \frac{\mu_g}{4C_g}\right) \epsilon_t$$

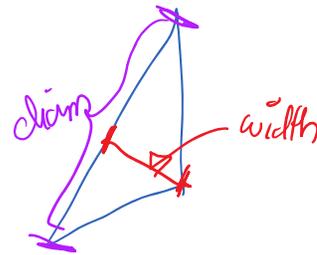
ie. linear rate with $\rho = \frac{\mu_g}{4C_g}$

$\mu_g \rightarrow$ geometric strong convexity constant

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

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

\downarrow \downarrow
 K_g condition # of set M



f is μ -strongly convex

a) FW with L-S, when $x^* \in \text{int}(M)$

b) AFW and M is a polytope

FW for SVM struct

dual SVM struct: $\min_{\alpha_i \in \Delta(1,1)} \frac{\Delta \|A\alpha\|^2}{2} - b^T \alpha$

ie. $M = \prod_{i=1}^n \Delta(1,1)$

$$A\alpha = \int \sum_{i=1}^n \sum_{\tilde{y}} \alpha_i(\tilde{y}) \Psi_i(\tilde{y}) = w(\alpha)$$

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

Elman et al.

$$x_{t+1} = \alpha y^{(t)} \quad \rightarrow \quad w^{(t+1)} = \alpha$$

FW step:

$$s_t = \arg \min_{s \in M} \langle s, \nabla f(x_t) \rangle$$

$$\nabla f(x_t) = \lambda \underbrace{A^T A x_t - b}_{w(x_t)} \quad w_t = A x_t$$

$$\begin{aligned} (\nabla f(x_t))_{i,y} &= \lambda \frac{\varphi_i(y)}{\lambda n} w_t - \frac{r_i(y)}{n} \\ &= \frac{1}{n} H_i(y; w_t) \end{aligned}$$

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

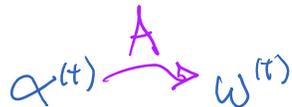
$$= \sum_i \underbrace{\min_{s_i \in M_i} \langle s_i, \nabla_i f(x_t) \rangle}_{M_i = \Delta(r_i)}$$

$$M_i = \Delta(r_i)$$

$$\min_y \langle \delta_y, \nabla_i f(x_t) \rangle$$

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

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



[loss-augmented decoding] ~~⊗~~ ~~⊗~~

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

$$\alpha^{(t+1)} = (1-\lambda)\alpha^{(t)} + \lambda \hat{\alpha}^{(t)}$$

↑ how need to maintain active set

$$\alpha_c^{(t+1)} = (1-\gamma)\alpha_c^{(t)} + \gamma \hat{S}_i^{(t)}$$

[here, need to maintain active sets:
 $\{S_i^{(t)}\}_{i=1}^n$]

$$w^{(t+1)} = (1-\gamma) \underbrace{Aq^{(t)}}_{w^{(t)}} + \gamma \underbrace{As^{(t)}}_{\sum_{i=1}^n \lambda_i \hat{y}_i^{(t)}}$$

you can choose via analytic LS or dual objective

recall primal obj:

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

$$p(w) = \frac{\lambda}{2} \|w\|^2 + \sum_{i=1}^n \max(0, \lambda_i (y_i - w^T \hat{y}_i))$$

(batch)
 Subgradient method update

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

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

if set $\beta = \frac{\gamma}{\lambda}$

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

FW perspective gives you an "adaptive step size" batch subgradient method