

today:
 • convergence of FW alg,
 • apply SVM struct alg.

Curvature constant C_f

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

this is affine invariant
 \Downarrow
 C_f is affine invariant

↑
 potential FW step update

→ worst case deviation from linear approximation

* by descent lemma, if ∇f is L -Lipschitz $\rightsquigarrow \left[\|\nabla f(x) - \nabla f(y)\|_* \leq L \|x - y\| \right]$

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

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

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

\searrow $x + \gamma(s - x)$

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

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

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

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

↑ affine invariant ↓ depends on $\|\cdot\|$

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

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

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

$\gamma^* = \min\{g_t, 1\}$

$$\gamma = \min\left\{\frac{g_t}{c_f}, 1\right\}$$

this gives you affine invariant adaptive step-size

$$f(x_{t+1}) \leq f(x_t) - \frac{g_t^2}{2c_f} \quad [\text{when } \frac{g_t}{c_f} \leq 1]$$

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

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

Thm.: FW alg. with γ_t chosen either as $\frac{2}{t+2}$ (when f is convex) or $\frac{g_t}{c_f}$ (line search) yields $\epsilon_t \leq \frac{2c_f}{t+2}$

- note!
- non-convex f
 - $\min_{S \subseteq T} g_S \leq O\left(\frac{1}{\sqrt{t}}\right)$
 - f is concave, $c_f = 0$?
 - $\min_{S \subseteq T} g_S \leq O\left(\frac{1}{t}\right)$

proof: let $x_\gamma = x_t + \gamma(x_t - x^*)$ } 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, $g_t \geq \epsilon_t$

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

$-g_t \leq -\epsilon_t$

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

* see notes 2017 for a cool CDF trick + induction proof

here, brute force approach to solve recurrence

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

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

initial condition Lipschitz constant

initial condition Lipschitz constant

use $(1+\gamma) \leq e^\gamma \quad \forall \gamma$
 $(1-\gamma) \leq \exp(-\gamma)$

base?

$$\Rightarrow E_{t+1} \leq E_0 \exp\left(-\sum_{s=0}^t \gamma_s\right) + \frac{C_\gamma}{2} \sum_{s=0}^t \gamma_s^2 \exp\left(-\sum_{u=s+1}^t \gamma_u\right)$$

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

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

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

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

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

$$\frac{1}{t} \sum_{s=0}^t \frac{1}{s} \rightarrow O\left(\frac{\log t}{t}\right)$$

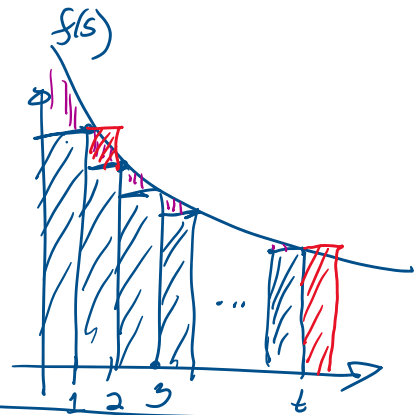
* in fact, if use $\gamma_t = \frac{1}{t+1}$, you do get $O\left(\frac{\log t}{t}\right)$ rate

but if $\gamma_t = \frac{2}{t+2}$, have bound says $O\left(\frac{\log t}{t}\right)$;

but a (tighter) direct analysis $O\left(\frac{1}{t}\right)$

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

[telescoping product]



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

f is decreasing.

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

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

10h30

linear rate for AFW:

(lemmings): linear rate: $E_{t+1} \leq (1-p) E_t \leq (1-p)^t E_0 \leq \epsilon \cdot \exp(-pt)$

(for gradient descent $\beta = \frac{\mu}{L} = \frac{1}{\kappa}$)
 condition #

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

⊗ recall for FW with LS: $\epsilon_{t+1} \leq \epsilon_t - \frac{g_t^2}{2c_g}$

[aside: $-g_t^2 \leq -\epsilon_t^2$]

$\Rightarrow \epsilon_t (1 - \frac{\epsilon_t}{2c_g})$
 $\rightarrow 0$

\Rightarrow this is why you don't rec. get linear rate

AFW paper, under some conditions

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

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

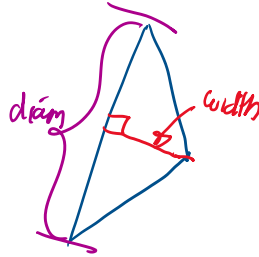
i.e. linear rate with $\rho = \frac{\mu_g}{4c_g}$

$\mu_g \rightarrow$ "geometric strong convexity constant"

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

f is μ -strongly convex
 a) FW with LS when x^* is int(M)
 b) AFW and M is a polytope

linear rate: $\rho = \frac{\mu_g}{4c_g} \Rightarrow \frac{\mu \cdot \text{Pwidth}(M)^2}{4L \cdot \text{diam}(M)^2}$
 \downarrow K_g \downarrow "condition # of set M"



FW for SVM struct dual

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

(i.o. $M = \prod_{i=1}^n \Delta_{|S_i|}$)

$A\alpha = \frac{1}{\lambda_n} \sum_{i=1}^n \sum_{\tilde{y}} \alpha_i(\tilde{y}) \psi_i(\tilde{y}) \triangleq w(\alpha)$

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

FW step:

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

let $w_t \triangleq w(\alpha_t) = A\alpha_t$

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

$(\nabla f(\alpha_t))_{i,y} = \lambda \sum_{\tilde{y}} \psi_i(\tilde{y})^T w_t - l_i(y)$
 $= \frac{1}{n} [l_i(y) - w_t^T \psi_i(y)]$

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

$H_i(y_i; w_t)$

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

(block structure)

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

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

$$\text{thus } z_t \triangleq \left(\hat{s}_i \right)_{i=1}^n$$

$$\text{where } \hat{s}_i \triangleq \delta \hat{y}_i(w_t)$$

$$\text{where } \hat{y}_i(w_t) = \underset{y \in Y_i}{\text{argmax}} H_i(y; w_t) \quad \text{[loss-augmented decoding]}$$

$$\alpha_t = \sum_u \alpha_u^+ s_u$$

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

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

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

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

$$\alpha^{(t+1)} = (1-\delta) \alpha^{(t)} + \delta S^{(t)}$$

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

you can choose via analytic L.S. on dual obj.

recall primal obj:

$$\varphi(w) = \frac{\lambda}{2} \|w\|^2 + \left(\frac{1}{n} \sum_{i=1}^n H_i(w) \right) \rightarrow \max_y \left(\varphi(\hat{y}) - w^T \psi(\hat{y}) \right)$$

$$p^*(w_t) = \lambda w_t - \frac{1}{n} \sum_{i=1}^n \psi_i(\hat{y}_i^{(t)})$$

$$w^{(t+1)} = w^{(t)} - \beta p^*(w^{(t)})$$

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

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

then batch subgradient step on primal is equivalent

to a batch FW step on dual with $\beta = \frac{\delta}{\lambda}$ step-size relationship

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

\otimes FW perspective gives you "adaptive step-size" for batch subgradient method \triangleright