

today: FW convergence & properties  
FW for SVM struct

necessary  
aside: first order optimality condition for constrained opt.

$$\min_{x \in M} f(x)$$

$x^*$  global min

(i.e.  $f$  looks increasing in all feasible directions)

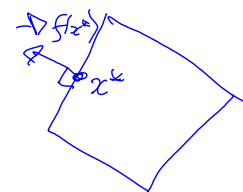
$$\Rightarrow \langle \nabla f(x^*), s - x^* \rangle \geq 0 \quad \forall s \in M$$

$$\Leftrightarrow \left( \min_{s \in M} \langle \nabla f(x^*), s - x^* \rangle \right) \geq 0$$

$$\Leftrightarrow \left( \max_{s \in M} \langle \nabla f(x^*), s - x^* \rangle \right) \leq 0$$

"stationary pt." for const. opt. problem

(sufficient cond. if  $f$  &  $M$  are convex)



FW-gap( $x^*$ )

i.o. want FW-gap( $x^*$ ) = 0  
since gap  $\geq 0$

↓ quantifying "non-stationarity"

see L-J. 2016 annex,

$$\min_{s \in M} \text{gap}(s) \leq O\left(\frac{1}{\sqrt{t}}\right) \quad \text{for FW with line search}$$

and

$f$  L-smooth and  $M$  bounded & convex

"non-convex FW"

but  $f$  not nec. convex

4) affine co-variance of FW

let  $\tilde{M}$  be a new domain

↗ affine transformation  
surjective  
 $\tilde{M} \xrightarrow{A} M$  ie.  $M = A\tilde{M}$

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

$$\downarrow$$

$$x(\tilde{x})$$

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

affine covariance FW: if run FW on  $\tilde{f}$  &  $\tilde{M}$  to get  $\tilde{x}_t$  iterates

$$\begin{array}{ccc} x_t \xleftarrow{A} \tilde{x}_t & \text{then } x_t \triangleq A\tilde{x}_t \text{ corresponds to running FW on } f \text{ \& } M & \text{(modulo tie-breaking)} \\ \text{FW} \uparrow & & \uparrow \text{FW} \\ f \text{ \& } M \xleftarrow{A} \tilde{f} \text{ \& } \tilde{M} \end{array}$$

why? inner product is affine invariant  $\nabla_{\tilde{x}} \tilde{f}(\tilde{x}) = \nabla_x f(A\tilde{x}) = A^T \nabla_x f(x)$

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

$$\langle \tilde{S}, A^T \nabla_x f(x_t) \rangle \quad \leftarrow x_t = A\tilde{x}_t$$

$$\tilde{S}_t = \arg\min_{\tilde{S} \in \tilde{M}} \langle A\tilde{S}, \nabla_x f(x_t) \rangle$$

$$\downarrow$$

$$S = A\tilde{S}$$

$$S_t = \arg\min_{\substack{S \in M \\ \parallel \\ A\tilde{M}}} \langle S, \nabla_x f(x_t) \rangle$$

$$\underline{\underline{S_t = A\tilde{S}_t}}$$

⊗ only first-order method I know with this property (like Newton which is 2nd order)

$\Rightarrow$  we want affine invariant analysis i.e.  $f(x_t) - f(x^*) = f(\tilde{x}_t) - f(\tilde{x}^*)$

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

$x_\gamma = (1-\gamma)x + \gamma s$

potential  
FW step

step-size

worst-case deviation from linear approximation

this depends  
on  $\langle \nabla f(x), s-x \rangle$   
is a fine invariant  
so  $C_f$  is " "

by descent lemma, if  $\nabla f$  is  $L$ -Lipschitz i.e.  $\|\nabla f(x) - \nabla f(y)\|_* \leq L \|x-y\| \quad \forall x, y \in M$

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

$\Rightarrow C_f \leq \sup_{\gamma} \frac{2}{\gamma^2} \left[ \frac{L}{2} \underbrace{\|x_\gamma - x\|^2}_{\|\gamma(s-x)\|^2} \right]$

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

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

affine  
invariant

depends on  $\|\cdot\|$  and thus  $A$

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

$\|d\|_* \triangleq \sup_{\|x\| \leq 1} \langle d, x \rangle$

$\Rightarrow$  generalized CS  $|\langle d, x \rangle| \leq \|d\|_* \|x\|$

$(\|\cdot\|_p)_* = \|\cdot\|_q$  where  $\frac{1}{p} + \frac{1}{q} = 1$

② by def. of  $C_f$ , we get affine invariant version of descent lemma:

$$f(x_8) \leq f(x) + \gamma \langle \nabla f(x), s-x \rangle + \frac{\gamma^2}{2} C_f \quad \forall \gamma \in [0,1] \\ \forall x, s \in M$$

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

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

for FW step with stepsize  $\gamma$ ,

$$(+) \quad \boxed{f(x_8) \leq f(x_t) - \gamma g_t + \frac{\gamma^2}{2} C_f}$$

optimal step-size  
minimizes the RHS

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

$$f(x_{8^*}) \leq f(x_t) - \frac{g_t^2}{2C_f}$$

aff. inv.  $\uparrow$  adaptive step-size

$$\leq f(x_t) - \frac{E_t^2}{2C_f} \quad \text{where } E_t \triangleq f(x_t) - f(x^*) \leq g_t$$

15h37

thm: FW alg. with  $\gamma_t$  chosen either  $\frac{2}{t+2}$  or  $\frac{g_t}{C_f}$  (line search)

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

proof: let  $x_8 = x_t + \gamma(s_t - x_t)$  + apply (+)

$$\begin{aligned}
 f(x_t) &\leq f(x_t) - \gamma g_t + \frac{\gamma^2}{2} C_f & \alpha_t \geq \varepsilon_t \text{ by convexity} \\
 &\leq f(x_t) - \gamma \varepsilon_t + \frac{\gamma^2}{2} C_f \\
 \underbrace{f(x_t) - f^*}_{\varepsilon_t} &\leq \underbrace{f(x_t) - f^*}_{\varepsilon_t} - \gamma \varepsilon_t + \frac{\gamma^2}{2} C_f
 \end{aligned}$$

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

- see notes (not given for cool ODE trick + induction)

here, brute force approach to solve the recurrence :

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

$\leq \dots$

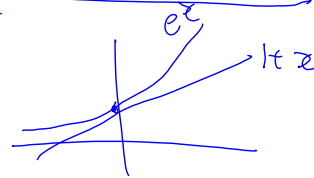
$$\boxed{\varepsilon_{t+1} \leq \prod_{s=0}^t (1 - \gamma_s) \varepsilon_0 + \frac{C_f}{2} \sum_{s=0}^t \gamma_s^2 \left( \prod_{u=s+1}^t (1 - \gamma_u) \right)}$$

initial condition
Lipschitz part

use  $(1 - \gamma) \leq e^{-\gamma} \quad \forall \gamma$

$(1 - \gamma) \leq e^{-\gamma}$

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



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

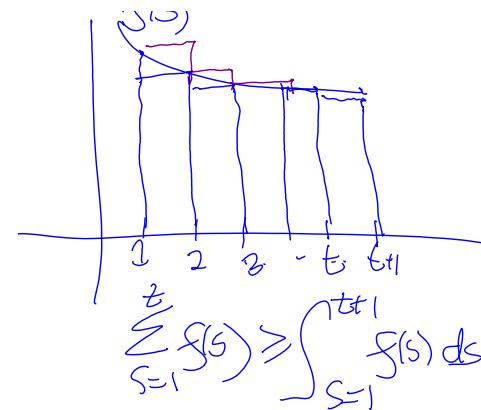
$$\gamma_s \sim \frac{1}{s} \Rightarrow \sum_{s=0}^t \gamma_s \sim \log(t)$$

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

$$\exp\left(-\sum_{u=s+1}^t \gamma_u\right) \sim \exp(-\log t/s) \sim O\left(\frac{s}{t}\right) \Rightarrow \sum_{s=0}^t \gamma_s^2 \exp\left(-\sum_{u=s+1}^t \gamma_u\right) \sim \sum_{s=0}^t \frac{1}{s} \cdot \frac{s}{t} \sim O\left(\frac{\log t}{t}\right)$$

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

see notes last year for  $\gamma_s = \frac{\alpha}{t+\alpha}$  ( $O\left(\frac{1}{t}\right)$  for  $\alpha \geq 2$ )



\* Linear rate for AFW:

linear rate:  $\varepsilon_{t+1} \leq (1-p) \varepsilon_t \leq \varepsilon_0 (1-p)^t \leq \varepsilon_0 \exp(-pt)$  (linear rate constant)

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

recall for FW:  $\varepsilon_{t+1} \leq (1-\gamma_t) \varepsilon_t$   $\gamma_t \approx \frac{g_t}{C_t} \approx \frac{\varepsilon_t}{C_t}$

AFW paper: get  $\varepsilon_{t+1} \leq \varepsilon_t - \gamma g_t + \frac{\gamma^2}{2} C_t$   $\gamma^* = \frac{g_t}{C_t}$

$$\varepsilon_{t+1} \leq \varepsilon_t - \frac{g_t^2}{2C_t}$$

under some conditions,  
can show that  $\alpha \varepsilon_t^2 \geq$

"growth condition"

$\varepsilon_t$  is extremely small

$$\varepsilon_{t+1} \leq \varepsilon_t - \frac{g_t^2}{2L_f}$$

under some conditions, can show that  $g_t^2 \geq \frac{\mu_f}{2} \varepsilon_t$   $\rightarrow$   $f$  is  $\mu$ -strongly convex

$$\Rightarrow \varepsilon_{t+1} \leq \left(1 - \frac{\mu_f}{2L_f}\right) \varepsilon_t$$

ie. linear rate with  $\rho = \frac{\mu_f}{2L_f}$

a) FW with L.S. when  $x^* \in \text{int}(M)$   
 b) AFW and  $M$  is polytope

Application of FW to structured SVM:

(dual)  $\min_{\alpha \in \Delta(\mathcal{S})} \frac{\lambda \|A\alpha\|^2 - b^T \alpha}{2}$

ie.  $M = \sum_{i=1}^n \Delta(\mathcal{S}_i)$

$$A\alpha = \frac{1}{\lambda n} \sum_i \sum_y \alpha(y) \varphi_i(y) = w(\alpha)$$

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

FW step:

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

$$\nabla f(\alpha_t) = \lambda A^T \frac{w_t}{\lambda n} - b$$

$$w_t = A\alpha_t$$

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

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

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

$$M_i = \Delta(\mathcal{S}_i)$$

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

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

$$1. \quad \sim \quad 1 \sim n$$

$$V(y|x) = \frac{1}{n} \sum H(y; w_t)$$

Thus  $S_t = (\hat{s}_i)_{i=1}^n$

$$\hat{s}_i = \delta_{\hat{y}_i(w_t)} \text{ where } \hat{y}_i(w_t) = \underset{\tilde{y} \in \mathcal{Y}_i}{\operatorname{argmax}} H(\tilde{y}; w_t) \quad [\text{loss-augmented decoding}]$$

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