

today : • FW } FCFW on SVM struct

• BCFW

continue FW for SVM struct

FW on dual is equivalent to batch subgradient update on primal with step-size relationship

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

$$\gamma^{(FW)} = \lambda \beta^{(\text{subgradient})}$$

recall : subgradient method converges of $O(\frac{1}{\epsilon})$

when step-size $\beta_t \leq \frac{1}{\mu} O(\frac{1}{t})$

FW method with fixed step-size schedule

$$\gamma_t = \frac{2}{t+2}$$

gives $d(\alpha^*) - d(\alpha^{(t)}) \leq \frac{2C\epsilon}{t+2}$

* FW-gap :

$$\begin{aligned} & \langle -\nabla f(\alpha^{(t)}), s^{(t)} - \alpha^{(t)} \rangle \\ &= \frac{1}{n} \sum_{i=1}^n \left[H_i(\hat{y}_i^{(t)}; w^{(t)}) - \sum_{\tilde{y}_i \in \mathcal{Y}_i} \alpha_i^{(t)}(\tilde{y}_i) H_i(\tilde{y}_i; w^{(t)}) \right] \\ &= p(w(\alpha^{(t)})) - d(\alpha^{(t)}) \quad [\text{Lagrangian gap of lecture 15}] \end{aligned}$$

here, FW-gap = Lagrangian gap on $(w(\alpha^{(t)}), \alpha^{(t)})$

[note: FW gap is not always Lagr. gap]
e.g. LP objective (dual)

* recall, that

$$\min_{s \geq 0} g_s^{FW} \leq 3 \cdot \frac{2C\epsilon}{t+2}$$

$$p(w^*) = d(\alpha^*)$$

$$\Rightarrow \text{guarantees on } \underbrace{p(w(\alpha^{(t)})) - p(w^*)}_{\text{primal subopt.}} + \underbrace{d(\alpha^*) - d(\alpha^{(t)})}_{\text{dual subopt.}}$$

$$\text{also, } g^{FW}(\hat{\alpha}_{\text{FW}}^{(t)}) \leq 3 \cdot \frac{2C\epsilon}{t+2}$$

↑
weighted avg.

when $g^{FW}(\cdot)$ is convex [this is case when f is a quadratic]

$$p(w(\alpha^{(t)})) = p(0) = \frac{1}{n} \sum_{i=1}^n \max_{\tilde{y}_i} \ell_i(\tilde{y}_i) \quad p(\alpha) \text{ mmmmmmm } \quad \text{FW-gap}(\alpha^{(t)}) = \text{Lagr-gap}$$

$$p(w^{(t)}) = p(0) = \frac{1}{n} \sum_{i=1}^n \max_{y_i} \ell_i(y_i)$$

$$d(\alpha^{(t)}) = 0$$

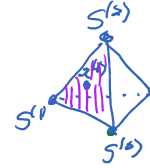
$FW\text{-gap}(\alpha^{(t)}) = p - \text{gap}$

FCFW "felly correctiv fw" variant

algorithm: re-optimize f over conv-hull ($\sum_{i=0}^t s^{(i)}$)

→ think of it as doing "felly line-search" on the "correction polytope"

note: could use AFW to do correction step app.



(see "barrier FW" paper)

[special case: min norm pt. alg. MNP] → sequence of affine projections → line search
→ state of the art alg. for submodular opt.

* turns out that (batch) FCFW on dual sumstruct is equivalent to the constraint generation / cutting plane alg. on the 1-slack formulation (primal)

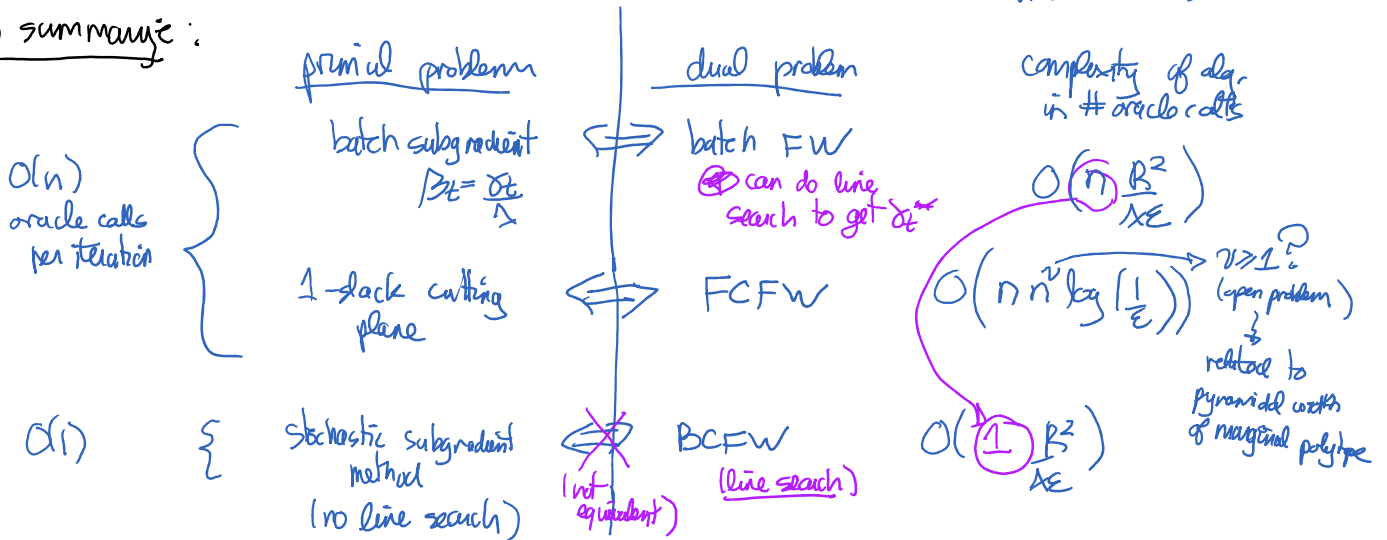
why? every $s^{(t)}$ corresponds $(\hat{y}_i^{(t)})_{i=1}^n$
 $\alpha \in \text{conv}(\{s^{(u)}\}_{u \leq t}) \quad \alpha = \sum_{u \leq t} \tilde{\alpha}_u s^{(u)}$

g strongly convex
 $g(Ax) + b^T x$

"gene. strongly convex"

$$w = Ax = \sum_{u \leq t} \tilde{\alpha}_u \left(\frac{1}{\Delta_n} \sum_{i=1}^n \mathbb{1}_{\{i \in u\}} \hat{y}_i^{(u)} \right)$$

to summarize:



14h51

Block-coordinate optimization

"huge scale" optimization \rightarrow [Nesterov 2010-2012]

proposed: randomized block-coordinate projected gradient method

setup: $\min f(x)$
 s.t. $x \in \prod_{i=1}^n M_i$
 $x = (\underbrace{x_1, \dots, x_n}_{\text{block}})$

algorithm: pick at random i_t
 then let $x_{i_t}^{(t+1)} = \text{Proj}_{M_{i_t}}(x_{i_t}^{(t)} - \frac{1}{L_{i_t}} \nabla_{i_t} f(x^{(t)}))$
 $x_j^{(t+1)} = x_j^{(t)} \quad \forall j \neq i_t$
Lipschitz const L_{i_t}

Only update block i_t at iteration t

Nesterov showed (uniform sampling) $\mathbb{E} f(x^{(t)}) - f^* \leq \frac{2}{t+4} [\sum_{i=1}^n L_i] \|x_0 - x^*\|^2$

(for convex f with L_i -lips. gradient per block)

Block-coordinate FW (BCFW): idea: do a FW step on block i

alg.: for $t=0, \dots$
 pick i unif. at random from 1 to n
 let $s_i^{(t)} = \text{argmin}_{s_i \in M_i} \langle s_i, \nabla_{i_t} f(x^{(t)}) \rangle$ [FW corner for block i]
 $\begin{cases} x_{i_t}^{(t+1)} = x_{i_t}^{(t)}(1-\delta_t) + \delta_t s_i^{(t)} \\ x_j^{(t+1)} = x_j^{(t)} \quad \forall j \neq i_t \end{cases}$
 $x_t = \text{line search } \text{argmin}_{\delta \in [0,1]} f(x^{(t)} + \delta(s_{[i]}^{(t)} - x_{[i]}^{(t)}))$
using predictive structure
 $\frac{2n}{t+2n}$ # of blocks

an important property: FW-gap = $\max_{s \in M} \langle -\nabla f(x), s-x \rangle = \sum_i \max_{s_i \in M_i} \langle -\nabla_{i_t} f(x), s_i - x_{i_t} \rangle$

$\triangleq \sum_i g_i(x)$ "block FW gap"

$g^{FW}(x) = \sum_i g_i^{FW}(x)$

as before, can show that $g_i(x) \geq f(x) - \min_{y \in M} f(y)$
 s.t. $y_j = x_j \quad \forall j \neq i_t$

\rightarrow motivated "gap sampling" variant

[Oskain & al. ICLR 2016]