

today • FW for SVM struct & FCFW
• BCFW

continue FW for SVM struct

FW on dual is equivalent to batch subgradient update on primal with step-size β_t rel.

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

$$g^{(FW)} = \beta_t \text{ (subgradient)}$$

recall: subgradient method converges $O(\frac{1}{\sqrt{t}})$

when step-size $\beta_t \propto \frac{1}{\sqrt{t}}$

FW method with fixed step-size $\beta_t = \frac{2}{t+2}$

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

* FW gap: (for SVM struct)

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

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

[note: FW gap is not always Lag. gap]
e.g. CRF objective on (dual)

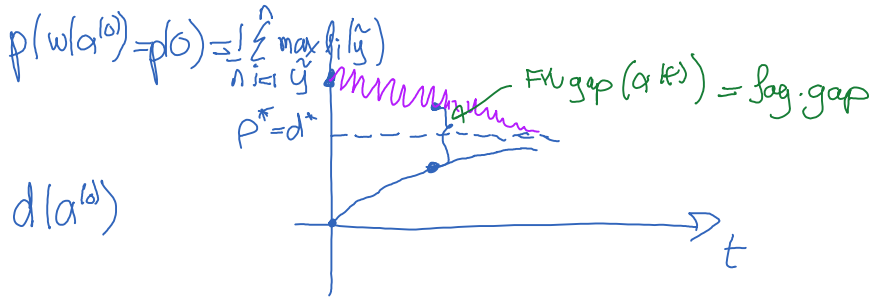
* recall that $\min_{s \leq t} g_s \leq 3 \cdot \frac{2C_f}{t+2}$

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

also (turns out) $g^{FW}(\hat{x}_{WA}^{(t)}) \leq 3 \cdot \frac{2C_f}{t+2}$
weighted avg.

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

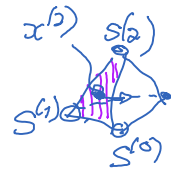
quadratic



FCFW "fully corrective FW" variant

algorithm: re-optimizing f over conv-hull $(\sum_{u=0}^t S^{(u)})$

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



note: could use AFW to do correction step approx.



(see "barrier FW")

[special case: min. norm point alg. MNP] → sequence of affine projections + line search to approximate the correction step

→ state of the art alg. for submodular opt.

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

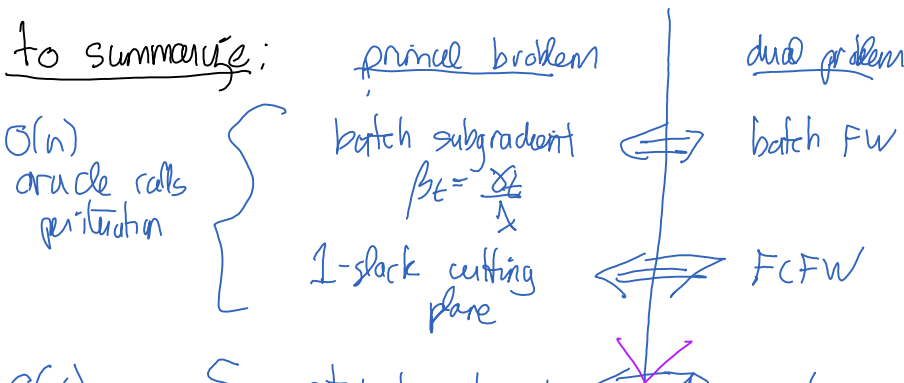
why? every $S^{(t)}$ corresponds to $(y_i^{(t)})_{i=1}^n$
 $\alpha \in \text{conv}(\sum_{u=0}^t S^{(u)}) \iff \alpha = \sum_{u \in T} \tilde{\alpha}_u S^{(u)}$

Solving QP in primal for cutting plane

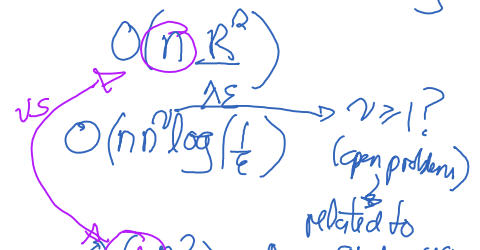
$$W = A\alpha = \sum_{u \in T} \tilde{\alpha}_u \left(\sum_{k=1}^n w_k(y_i^{(u)}) \right)$$

correction step on dual for FCFW

to summarize:



complexity of alg in # oracle calls → loss augmented decoding



$O(1)$ plane
 { stochastic subgradient method (no line search) }
 BCFW (not equivalent) (line search)
 (\mathbb{R}^2) (open problem) related to pyramidal width of marginal polytope

15h35

block-coordinate optimization

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

proposed: randomized block-coordinate projected gradient method

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

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

[only updating block i_t at iteration t]

Nesterov showed (uniform sampling)

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

(for convex f with L_i -Lipschitz gradient per block)

Block-coordinate FW (BCFW) : idea is to do a FW step on block i_t

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]
 $\begin{cases} x_{i_t}^{(t+1)} = x_{i_t}^{(t)} (1-\delta_t) + \delta_t s_{i_t}^{(t)} \\ x_j^{(t+1)} = x_j^{(t)} \quad \forall j \neq i_t \end{cases}$

when $s_{[i]} = \begin{pmatrix} 0 \\ \vdots \\ s_i \\ \vdots \\ 0 \end{pmatrix}$ block i

$\delta_t = \begin{cases} \text{line search } \text{argmin}_{\delta \in [0,1]} f(x^{(t)} + \delta (s_{[i]}^{(t)} - x_{[i]}^{(t)})) \\ \frac{2}{t+2n} \end{cases}$ # of blocks

* an important property: FW-gap = $\max_{S \in \mathcal{M}} \langle -\nabla f(x), S-x \rangle \stackrel{\text{product using structure}}{=} \max_{S_i \in \mathcal{M}_i} \langle -\nabla f(x), S_i - x_i \rangle$

$$g^{FW}(x) = \sum_{i=1}^n g_i^{FW}(x)$$

↳ motivated "gap sampling" variant

as before, you can show that $g_i(x) \geq f(x) - \min_y f(y)$
 $\forall S_i, y_i = \arg \min_{y \in \mathcal{M}_i} f(y)$

[Osofin et al. ICML 2016]

Convergence of BCFW

$$C_f^{(i)} \leq L_i \text{diam}(\mathcal{M}_i)^2$$

$$C_f^{(i)} \stackrel{\text{product using structure}}{=} \sum_{j=1}^n C_f^{(j)}$$

$$\underbrace{\mathbb{E}f(x^{(t)}) - f^*}_{\triangleq \epsilon_t} \leq \frac{2n}{t+2n} [C_f^{(i)} + f(x^{(t)}) - f^*] \text{ for } \delta_t = \frac{2n}{t+2n}$$

if you use line search $\epsilon_t \leq \frac{2n C_f^{(i)}}{t-t_0+2n}$ for $t \geq t_0$

$t_0 \leq n \log \frac{2(f(x^{(0)}) - f^*)}{\epsilon}$
 time to ensure that $\epsilon_t \leq C_f^{(i)}$

batch FW

$$\epsilon_t \leq \frac{2C_f}{t+2}$$

one can show that $C_f^{(i)} \leq C_f$ for quadratic functions

BCFW

$$\epsilon_t \leq \frac{2n C_f^{(i)}}{t-t_0+2n}$$

$$C_f \text{ vs. } n C_f^{(i)}$$

but BCFW is n times cheaper than batch FW complexity # of oracle calls FW $\rightarrow O(n C_f / \epsilon)$

\Rightarrow BCFW is "never" slower than batch FW (when not using parallelization)

but for SUM-struct you have $n C_f^{(i)}$ is actually n times smaller than C_f (?)

$$\text{BCFW } O\left(\frac{n C_f^{(i)}}{\epsilon}\right)$$

extensions to BCFW:

- non-uniform sampling eg- $\begin{cases} g_i(x) \\ C_f^{(i)} \end{cases}$

ICML 2016

-
- Non-uniform sampling e.g. $\begin{cases} g_i(x) \\ f^{(i)} \end{cases}$
 - Using away-step, etc. to get "linear convergence"
- } ICML 2016