

# Lecture 17 - scribbles - FW for SVMstruct

Tuesday, March 20, 2018 14:27

- today:
- FW for SVMstruct
  - BCFW
  - variance reduced SGD [SAG]

## FW for SVMstruct:

FW on dual of SVMstruct objective:

$$\min_{\alpha: \alpha_i \in \Delta_{\mathcal{Y}_i}} \underbrace{\frac{\lambda \|A\alpha\|^2}{2} - b^T \alpha}_{f(\alpha) \triangleq -d(\alpha)}$$

$$M = \bigtimes_{i=1}^n \Delta_{\mathcal{Y}_i}$$

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

$$\begin{aligned} (\nabla f(\alpha))_{i,y} &= -\frac{1}{n} H(y; w(\alpha)) \\ &= (\lambda A^T A \alpha - b)_{i,y} \end{aligned}$$

$$w(\alpha) = A\alpha$$

$$b_i(y) = \frac{1}{n} \ell_i(y)$$

$$A_{i,y} = \frac{\gamma_i(y)}{\lambda n}$$

$$\text{FW step: } \min_{S \in M} \langle S, \nabla f(\alpha_t) \rangle$$

$$= \sum_{i=1}^n \min_{S_i \in \Delta_{\mathcal{Y}_i}} \langle S_i, \nabla_i f(\alpha_t) \rangle$$

$$S_i^{(t)} = \delta \hat{y}_i^{(t)} \text{ where } \hat{y}_i^{(t)} = \underset{\tilde{y} \in \mathcal{Y}_i}{\text{argmax}} H_i(\tilde{y}; w^{(t)}) \leftarrow \text{loss-augmented decoding}$$

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

A ↗

$$w^{(t+1)} = (1-\delta) \underbrace{A \alpha^{(t)}}_{w^{(t)}} + \delta \underbrace{A S^{(t)}}_{\frac{1}{\lambda n} \sum_{i=1}^n \gamma_i(\hat{y}_i^{(t)})}$$

recall primal:

$$p(w) = \frac{\lambda \|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n H_i(w) \quad \nabla \max_{\tilde{y} \in \mathcal{Y}_i} \ell_i(\tilde{y}) - w^T \gamma_i(\tilde{y})$$

$$\partial p = \lambda w^{(t)} - \frac{1}{n} \sum_{i=1}^n \gamma_i(\hat{y}_i^{(t)})$$

⊛ advantage of FW: adaptive step-size  
using analytic line search

$$\partial p = \lambda w^{(t)} - \frac{1}{n} \sum_{i=1}^n \gamma_i \tilde{y}_i^{(t)}$$

using analytic line search

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

$$w^{(t+1)} = w^{(t)} - \beta \frac{\partial p}{\partial w} = (1 - \lambda \beta) w^{(t)} + \beta \sum_{i=1}^n \gamma_i \tilde{y}_i^{(t)}$$

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

then batch subgradient step on primal is equivalent to batch FW step on dual

recall: Subgradient method convergence result  $O(\frac{1}{t})$

needed a step size  $\beta_t \leq \frac{1}{\mu} O(\frac{1}{t})$

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

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

$$\text{here, FW-gap } \langle -\nabla f(\alpha^{(t)}), S^{(t)} - \alpha^{(t)} \rangle = \frac{1}{n} \sum_{i=1}^n \left[ H(\tilde{y}_i^{(t)}; w^{(t)}) - \sum_{\tilde{y} \in \mathcal{Y}_i} \alpha_i^{(t)}(\tilde{y}) H(\tilde{y}; w^{(t)}) \right]$$

$$= p(w(\alpha^{(t)})) - d(\alpha^{(t)})$$

(Bregman gap of lecture 14?)

\* note one can show

$$\cdot \text{ that } \min_{S \in \mathcal{Y}} g_S \leq 3 \cdot \frac{2C_S}{t+2}$$

FW gap

$$\cdot \text{ also have } g^{\text{FW}}(\hat{x}_{\text{FW}}^{(t)}) \leq 3 \cdot \frac{2C_S}{t+2}$$

↑  
weighted average

when  $g(\cdot)$  is convex [this is case for  $\mathcal{S}$  being quadratic]

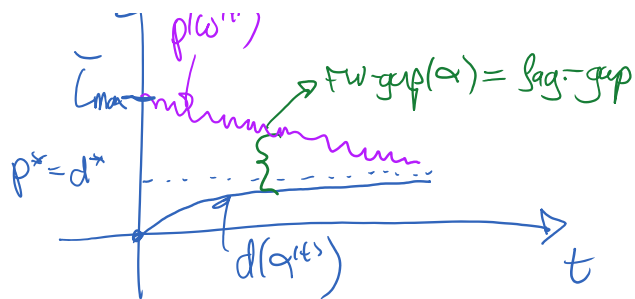
$$p(\alpha) = \frac{1}{n} \sum_{i=1}^n \max_{\tilde{y}} \ell_i(\tilde{y}) = \bar{L}_{\max}$$

$$\bar{L}_{\max} \uparrow p(w^{(t)})$$

$$\text{FW-gap}(\alpha) = \text{Breg-gap}$$

$$p(0) = \downarrow \sum_{i=1}^n \max_{\tilde{y}} (l_i | \tilde{y}) = L_{\max}$$

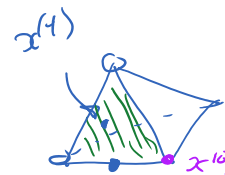
$$d(\alpha^{(0)}) = 0$$



FCFW "fully correct FW" variant:

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

[special case: min norm point (MNP) algorithm]  
→ state of art for submodular optimization



turns out that

(batch) FCFW on dual of sumstruct

is equivalent to the constraint generation / cutting plane approach for 1-stock formulation

why?

every  $S^{(t)}$  corresponds  $\{\hat{y}_i^{(t)}\}_{i=1}^n$

$$\alpha \in \text{conv}(\{S^{(u)}\}_{u \leq t}) \quad \alpha = \sum_{u \leq t} \tilde{\alpha}_u S^{(u)}$$

$$w = A\alpha = \sum_{u \leq t} \tilde{\alpha}_u \downarrow_{A_n} \left( \sum_{i=1}^n \gamma_i | \hat{y}_i^{(u)} \right)$$

to summarize

$O(n)$  oracle calls  
per iteration

primal problem

batch subgradient  
 $\beta_t = \tilde{\alpha}_t$

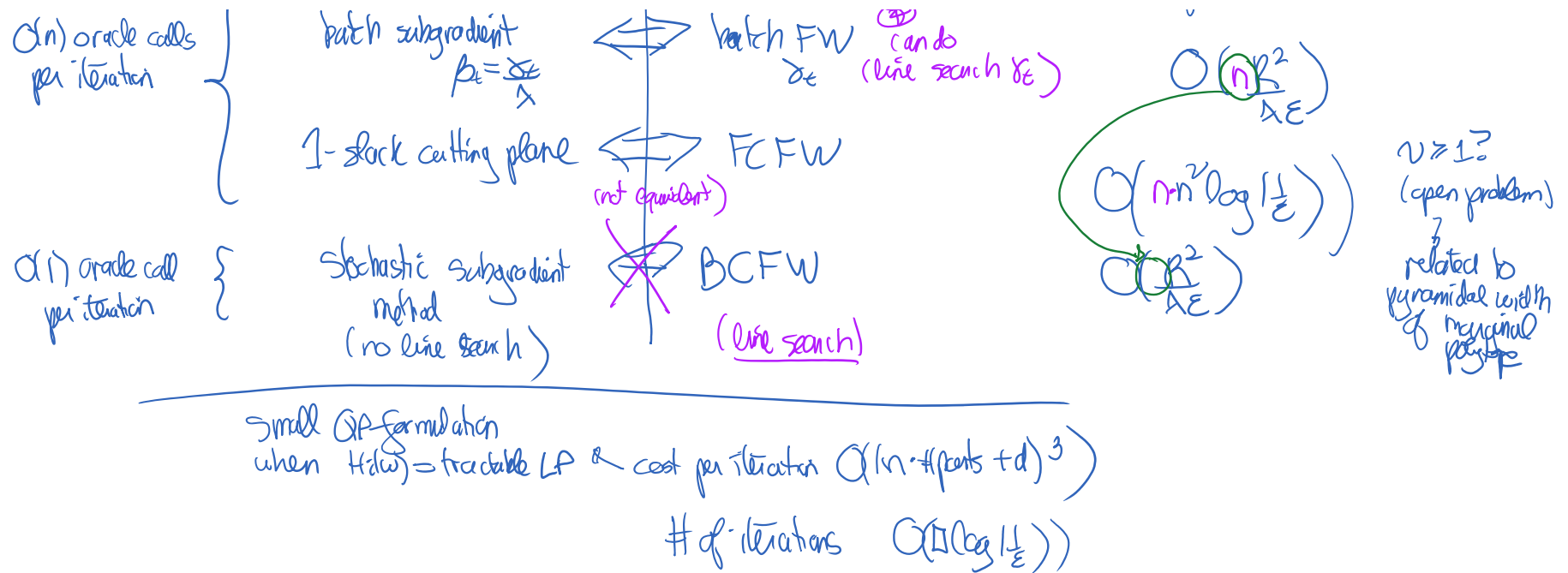


dual problem

batch FW  
 $\tilde{\alpha}_t$  (can do like search  $\tilde{\alpha}_t$ )

complexity in #  
of oracle calls

$$O(n^2)$$



## Block-coordinate optimization

[wright reference 2016 survey?]

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

proposed: randomized block-coordinate projected gradient method

$$\nabla = \begin{pmatrix} \nabla_1 \\ \vdots \\ \nabla_i \\ \vdots \\ \nabla_n \end{pmatrix}$$

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

alg: pick  $i$  at random

$$\begin{aligned} \text{then let } & \begin{cases} x_i^{(t+1)} = \text{Proj}_{M_i} \left( x_i^{(t)} - \frac{1}{L_i} \nabla_i f(x^{(t)}) \right) \\ x_j^{(t+1)} = x_j^{(t)} \quad \forall j \neq i \end{cases} \end{aligned}$$

Lipschitz constant of  $\nabla_i f(x)$

[only update block  $i$  at iteration  $t$ ]

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

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

alg: for  $t=0, \dots$

pick  $i$  at random

let  $S_i^{(t)} \triangleq \arg \min_{S_i \in M_i} \langle S_i, \nabla_i f(x^{(t)}) \rangle$  [FW corner for block  $i$ ]

$$S[i] = \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

$$\begin{cases} x_i^{(t+1)} = x_i^{(t)}(1-\gamma_t) + \gamma_t S_i^{(t)} \\ x_j^{(t+1)} = x_j^{(t)} \text{ for } j \neq i \end{cases}$$

$$\gamma_t = \begin{cases} \text{line search:} \\ \arg \min_{\gamma \in [0,1]} f(x^{(t)} + \gamma(S[i] - x^{(t)})) \end{cases}$$

$2n$  # of blocks  
 $t+2n$

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

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

$g_i(x)$  "block FW gap"

↳ motivated "gap sampling"  
variant [Oscin & al. JICML 2016]

as before, can show  $g_i(x) \geq f(x) - \min_{y: y_j=x_j \forall j \neq i} f(y)$

convergence result for BCFW:

$$C_f^{(i)} \leq L_i \text{diam}(M_i)^2$$

$$C_f^{\otimes} \triangleq \sum_{i=1}^n C_f^{(i)}$$

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

$$\text{if you use line search: } \epsilon_t \leq \frac{2n C_f^{\oplus}}{t - t_0 + 2n} \quad \text{for } t \geq t_0$$

$$\text{where } t_0 \leq n \log \frac{2(f(x^{(t)}) - f(z^*))}{C_f^{\oplus}}$$

↑  
time to ensure  $\epsilon_t \leq C_f^{\oplus}$

one can show that  $C_f^{\oplus} \leq C_f$  for quadratic functions

$$\text{compare batch FW} \rightarrow \frac{2C_f}{t+2}$$

vs.

$$\text{BCFW with line search} \rightarrow \frac{2n C_f^{\oplus}}{t - t_0 + 2n}$$

but BCFW is  $n$  times cheaper than batch FW

$\Rightarrow$  BCFW is "never" slower than batch FW

ie. # of oracle calls  
to reach  $\epsilon$ -error

$$\text{FW: } O\left(\frac{n C_f}{\epsilon}\right)$$

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

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

• using away step, etc... to "get"  $\epsilon$ -error convergence

} ICMC 2016

Application of BCFW for SVM struct:

- getting  $\hat{y}_i$  is one loss-augmented decoding call for example  $i$

- $w = Ax = \sum_i A_i x_i$   
 $\triangleq w_i$

↑ need to store those in memory

→ you update  $\begin{cases} w_i^{(t+1)} = w_i^{(t)} + \delta (w_i^{(t)} - w_i^{(t)}) \\ w_j^{(t+1)} = w_j^{(t)} \quad \forall j \neq i \end{cases}$

for SVM struct,  $C_f^{(i)} \leq \frac{4R^2}{\lambda n^2} \Rightarrow C_f^{(i)} = \frac{4R^2}{\lambda n} \approx \frac{C_f}{n}$

vs.  $C_f \leq \frac{4R^2}{\lambda}$

uses affine invariance crucially

ie. BCFW is "n times" faster than batch FW for SVM struct? 28

$$C_f \leq L_{1,1} \cdot \text{diam}(M)^2$$

if use  $\ell_2$ :  $\|\cdot\|_2$  get bad bound  $\text{diam}(M)^2 = 2n$

Lipschitz constant in  $\ell_2$ -norm = largest e-value of Hessian

$$\lambda ATA = \frac{1}{\lambda n^2} \langle \psi_i(y), \psi_j(y) \rangle_{(y_i, y_j)}$$

say e.g.  $\langle \psi_i(y), \psi_j(y) \rangle \approx 1$  for lots of output

→ get <sup>largest</sup> e-value scale with dimension of matrix

$$(\mathbb{1} \mathbb{1}^T) \mathbb{1} = \mathbb{1} \cdot \text{dim}$$

→ here, dim is exponential  
 ⇒ useless bound?