

## Lecture 15 - cutting plane

Tuesday, February 26, 2019 14:33

- today:
- more SVM struct properties
- M<sup>2</sup> net dual
- cutting plane alg.

more properties of SVM struct dual:

$$w(\alpha) = \frac{1}{\lambda n} \sum_{i=1}^n \left( \sum_{j \in S_i} \alpha_{ij} y_j \right)$$

$$\text{Let } R_i \triangleq \max_i |\alpha_{ij} y_j| / b$$

$$\bar{R} = \sqrt{\sum_{i=1}^n R_i^2}$$

$$\begin{aligned} \text{Then 1)} \quad \|w^*\|_2 &\leq \frac{1}{\lambda} \sqrt{\sum_{i=1}^n \sum_{j \in S_i} \alpha_{ij}^2} \underbrace{\|\psi_j\|_2}_\text{=1} \\ &\leq \frac{1}{\lambda} \left( \sum_{i=1}^n R_i \right) = \frac{\bar{R}}{\lambda} \end{aligned}$$

$$\text{2) kernel trick: } \langle w, \varphi(x, y) \rangle = \sum_{i=1}^n \sum_{j \in S_i} \alpha_{ij} \underbrace{\langle \psi_j, \varphi(x, y) \rangle}_{k(x^{(i)}, y^{(i)}; x, y)}$$

$$\begin{aligned} &\varphi(x^{(i)}, y^{(i)}) - \varphi(x^{(i)}, \tilde{y}) \\ &\downarrow \\ &k(x^{(i)}, y^{(i)}; x, y) \\ &- k(x^{(i)}, \tilde{y}; x, y) \end{aligned}$$

$$\|w(\alpha)\|^2 \rightarrow \alpha^T K \alpha$$

$$\text{3) suppose scale features } \boxed{\tilde{\psi} = b\psi}$$

3) suppose scale features  $\tilde{\psi} = b\psi$

$$\tilde{H}_i(y; w) = l_i(y) - \langle \tilde{w}, \tilde{\psi}_i(y) \rangle$$

$$\tilde{w}^* = \frac{1}{b} \sum_{i=1}^n \sum_{y \in S} \alpha_i^*(y) \underbrace{\tilde{\psi}_i(y)}_{b\psi_i(y)}$$

Let  $\tilde{\lambda} = b^2 \lambda$

$$\Rightarrow \tilde{w}^*(\alpha^*) = \frac{1}{b} \left[ \sum_{i=1}^n \sum_{y \in S} \alpha_i^*(y) \psi_i(y) \right]$$

$$\text{use } \alpha^* = \alpha^* \Rightarrow \tilde{w}^* = \frac{w^*}{b}$$

$$\Rightarrow \tilde{H}_i(y; \tilde{w}^*) = l_i(y) - \langle w^*, \psi_i(y) \rangle = H_i(y; w^*)$$

i.e.  $\alpha^*$  is really optimal for the new problem with  $\tilde{\psi} \notin \mathcal{X}$  ▷

4) similarly, can show  $\tilde{\lambda} = b \cdot \lambda \Rightarrow \tilde{\lambda} = \frac{\lambda}{b}$  get same sol'n

$M^3$ -net example (dual): (getting compact dual)

$$w(\alpha) = A\alpha = \sum_i A_i \alpha_i \quad \alpha_i \in \Delta(S_i)$$

$$\text{suppose } \psi(y) = \sum_c \psi_c(y_c)$$

$$\begin{aligned} \text{then } A_i \alpha_i &= \sum_y \alpha_i(y) \psi_i(y) = \sum_y \alpha_i(y) \sum_c \psi_{i,c}(y_c) \\ &= \sum_i \sum_c \alpha_{i,c}(y_c) \left[ \sum_i \alpha_i(y) \right] \end{aligned}$$

$$= \sum_c \sum_{\tilde{y}_c} \alpha_{i,c}(\tilde{y}_c) \left[ \sum_{y^*} \alpha_i(y^*) \right]$$

$\tilde{y}_c$  s.t.  
 $\tilde{y}_c = \tilde{y}_0$

$\cong \mu_{i,c}(y_c)$   
"marginal variable"

$$\alpha_i \in \Delta_{\mathcal{Y}_i} \Rightarrow \mu_i \in M_i^0$$

$\downarrow$   
marginal  
polytope

thus  $A_i \alpha_i = \tilde{A}_i \mu_i$  where  $(\tilde{A})_{s,c,c,y_c} = \frac{\alpha_{i,c}(y_c)}{n}$  "marginal variable"

↪ # of columns is  $\sum_c |\mathcal{Y}_c|$

Similarly, suppose  $\ell_i(y) = \sum_c \ell_{i,c}(y_c)$

define  $\tilde{b}_{i,c}(y_c) \triangleq \frac{\ell_{i,c}(y_c)}{n}$

$\langle b_i, \alpha_i \rangle = \langle \tilde{b}_i, \mu_i(\alpha_i) \rangle$

④ Thus we can replace

$$\max_{\alpha_i \in \Delta_{\mathcal{Y}_i}} \rightarrow \frac{\|\tilde{A}\alpha\|^2}{2} + \tilde{b}^\top \alpha$$

with 
$$\max_{\mu_i \in M_i^0} -\lambda \|\tilde{A}\mu\|^2 + \tilde{b}^\top \mu$$

→ this is a tractable size GP

if  $M_i$  is tractable

M3-net paper

used "structured SMO" algorithm

In order to make it tractable...

if  $G_i$  was triangulated;

then  $M_i = 1 \Rightarrow$  (partial constraint) and  $L = 1$

used structured SMO algorithm

block-coordinate ascent using  
pair of variables on this QP

[similar to "pairwise FW"]

15h 15

If  $\mathcal{G}_i$  was triangulated,

then  $M_i = L_i$  (local consistency polytope)

$$\max_{\mu \in L_i} \underbrace{\mu_i}_{\rightarrow \text{tractable QP?}}$$

### constraint generation alg.

[Tsochantarakis et al. JMLR 2005]

$$\min_{w, \xi} \frac{\Delta \|w\|^2}{2} - \left[ \sum_{i=1}^n \xi_i \right] \quad (\text{P})$$

n-slack  
version

$$\text{s.t. } \begin{cases} \xi_i \geq H_i(y_i; w) & \forall y \in \mathcal{Y}_i \\ \xi_i \geq 0 \end{cases}$$

$$\max_{\alpha} \underbrace{-\frac{\Delta}{2} \|A\alpha\|^2}_{(\text{D})} + b^T \alpha$$

s.t.  $\alpha_i \in \Delta_i$

$\rightarrow$  # constraints  $\rightarrow$  # variables

Vs.

1-slack version

[ML 2009 paper]

$$\min_{w, \xi} \frac{\Delta \|w\|^2 + \xi}{2} \quad (\text{P})$$

$$\text{s.t. } \begin{cases} \xi \geq \sum_{i=1}^n H_i(y_i; w) & (\forall y_i \in \mathcal{S}_i) \\ \xi \geq 0 \end{cases}$$

$\sum_{i=1}^n H_i(y_i)$  constraints

$$\left[ \sum_{i=1}^n H_i(y_i) + \langle w, \sum_{i=1}^n y_i \rangle \right]$$

$$\max_{\alpha} \underbrace{-\frac{\Delta}{2} \|A\alpha\|^2}_{(\text{D})} + \langle b, \alpha \rangle$$

$$\alpha \in \Delta \left( \bigcap_{i=1}^n \mathcal{S}_i \right)$$

$$w(\alpha) = \sum_{y_i \in \mathcal{Y}_i} \alpha_i H_i(y_i) \left( \sum_{j=1}^n \gamma_j H_j(y_i) \right)$$

$$\left[ \sum_{i=1}^n \ell_i(\tilde{y}_i) + \langle w, \sum_i \xi_i \tilde{y}_i \rangle \right]$$

(Old)  
instead of  $O(d \cdot n)$   
in m-slack formulation

$\rightarrow$  big memory saving

### n-slack SVMstruct algorithm:

iterate solving QP with more and more constraints

1) start with no constraint on  $w \Rightarrow w^{(0)} = 0$   
 $\xi^{(0)} = 0$

2) repeat: for each  $i$ , find  $\hat{y}_i = \underset{y \in \mathcal{S}_i}{\operatorname{argmax}} H_i(y; w^{(t)})$  [loss-augmented decoding]

• add  $\xi_i \geq H_i(\hat{y}_i; w)$  constraint to QP (if not already there)

$\hookrightarrow$  then resolve QP( $w, \xi$ ) with these constraints  
 b get  $w^{(t+1)}, \xi^{(t+1)}$

Stop when primal-dual gap  $\leq \varepsilon$  } note: this takes  $O(n)$

[in 2005, showed that alg. stops after  $O(\frac{1}{\varepsilon})$  iterations]

refined later to  $O(\frac{1}{\varepsilon})$  at least for 1-slack version

### Frank-Wolfe algorithm

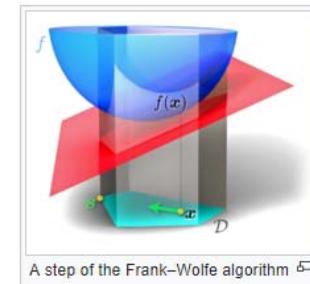
$\hookrightarrow$  for smooth constrained optimization

$$\left[ \begin{array}{c} \text{motivation} \\ \text{in an context} \end{array} \quad \text{dual of SVMstruct} \quad \min_{\alpha \in \Delta^{n+1}} \frac{\|A\alpha\|^2}{2} - b^T \alpha \right]$$

1940s: simplex algorithm to solve LPs

1956: Marguerite Frank & Phil Wolfe

→ non-linear opt. by iterating LPs



$$\text{Setup: } \min_x f(x)$$

st.  $x \in M$

•  $f$  is  $L$ -smooth i.e.  $\nabla f$  is  $L$ -Lipschitz and  $f$  is convex

•  $M$  is convex and bounded set

(more generally can also get rates for  $f$  only convex e.g.  $\nabla f$  is Hölder cont.)

by convexity

$$f(s) \geq f(x_t) + \langle \nabla f(x_t), s - x_t \rangle \quad \forall s \in M$$

FW algorithm:

start with  $x_0 \in M$

for  $t = 0, \dots$

FW convex

"linear minimization oracle"

LMO

minimizing RHS

linear approx. of  $f$  at  $x_t$

$$\text{compute } s_t = \underset{s \in M}{\operatorname{argmin}} \langle s, \nabla f(x_t) \rangle$$

stepping criterion

$$\left[ \text{let } \alpha_t \triangleq \langle s_t - x_t, -\nabla f(x_t) \rangle \text{ FW gap} \quad \text{if } \alpha_t \leq \xi, \text{ output } x_t \right]$$

$$x_{t+1} = (1 - \gamma_t)x_t + \gamma_t s_t \quad \text{step-size}$$

$\gamma_t \in [0, 1]$  (convex combo)

$$= x_t + \gamma_t \frac{(s_t - x_t)}{\|s_t\|}$$

end

output  $x_{t+1}$

end  
output  $x_{t+1}$

$\leftarrow$

$$\text{step size choice : } \gamma_t = \begin{cases} \text{universal choice } \frac{2}{L^2} \\ \text{line search } \gamma_t = \underset{\gamma \in [0, 1]}{\operatorname{argmin}} f(x_t + \gamma(s_t - x_t)) \end{cases}$$

adaptive :  $\frac{g_t}{L \|s_t\|_2^2}$  or  $\frac{g_t}{C_S}$  truncated at 1

Affine invariant constant

big motivation for FW

is LMO is often much cheaper  
than projections  
and cheap for many sets  $M$  appearing in ML

properties: 1)  $f(x_t) - \min_{\substack{x \in M \\ \in g^*}} f(x) \leq O\left(\frac{1}{t}\right)$

2) FW-gap  $g_t \geq f(x_t) - f^*$   $\rightarrow$  certificate of optimality

$$\min_{s \in S} g_s \leq O\left(\frac{1}{t}\right)$$

[i.e. we will stop in  $O\left(\frac{1}{\epsilon}\right)$  iterations?]

3)  $x_t = p_0^t x_0 + \sum_{u=1}^t p_u^t s_{u-1}$   $\rightarrow x_t$  has "sparse" expansion in terms of the FW-corners  $\{s_u\}_{u=1}^{2^k-1}$   
 where  $\sum_{u=0}^t p_u^t = 1$   $p_u^t \geq 0$   $\circlearrowleft$  "Sparse method"  
 $\rightarrow$  popular in ML

[we'll see can run FW on structured dual assuming can compute LMO]  $\downarrow$

(which is loss-<sup>ed</sup>-segmented )  
decoding here

- 4) FW is affine covariant (like Newton's method)
- 5) there is a  $\mathcal{Q}(f_t)$  lower bound for FW-like methods for  $t \leq d$