

Lecture 15 - cutting plane alg.

Tuesday, March 14, 2023 2:28 PM

- today:
- more SVM struct properties
 - M^2 -net dual
 - cutting plane alg.
 - FW alg.

more properties of SVM struct dual

$$\text{primal-dual gap} \quad p(w) - d(\alpha) \geq 0 \quad \forall w \text{ or feasible}$$

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

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

primal subpt. dual subpt.

certificate of
primal or dual
suboptimality

$$\text{gap}(\alpha) = p(w(\alpha), s(\alpha)) - d(\alpha)$$

$$= \frac{1}{n} \sum_{i=1}^n \left(\|w\|^2 + \sum_j H_i(w(\alpha)) \right) + \frac{1}{n} \|w\|^2 - \sum_{i,j} \alpha_i y_j \langle w, y_j \rangle$$

$$= \frac{1}{n} \sum_{i=1}^n \left(H_i(w(\alpha)) - \sum_j \alpha_i y_j \langle w, y_j \rangle \right) + \frac{1}{n} \langle w, w(\alpha) \rangle$$

$$\text{gap}(\alpha) = \frac{1}{n} \sum_{i=1}^n \left[H_i(w(\alpha)) - \sum_j \alpha_i y_j H_i(y_j; w(\alpha)) \right]$$

$$\max_{\tilde{y}} H_i(\tilde{y}; w)$$

use this to get a
bound dual subpt. of α

$$w(\alpha) = \frac{1}{n} \sum_{i,j} \alpha_i y_j \gamma_i y_j$$

then 1) $\|w^*\|_2 \leq \frac{1}{n} \sum_{i,j} \alpha_i y_j \|\gamma_i y_j\|_2$

$$\text{let } R_i \triangleq \max_{\tilde{y}} \|\gamma_i(\tilde{y})\|_2$$

$$\bar{R} = \frac{1}{n} \sum_i R_i$$

$$\leq \frac{1}{n} \left(\sum_i R_i \right) = \bar{R}$$

$$\varphi(x^{(i)}, y^{(i)}) - \varphi(x^{(i)}, \tilde{y})$$

2) kernel trick : $\langle w(\alpha), \varphi(x, y) \rangle = \frac{1}{n} \sum_{i,j} \alpha_i y_j \underbrace{\langle \gamma_i y_j, \varphi(x, y) \rangle}_{k(x^{(i)}, y^{(i)}, x, y) - k(x^{(i)}, y)}$

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

$$K_{i,j} \triangleq \langle \gamma_i y_j, \gamma_j y_j \rangle$$

$$\Rightarrow k_{ij}(\tilde{y}; \tilde{y}) \triangleq \langle \psi_i(\tilde{y}), \psi_j(\tilde{y}) \rangle$$

3) Suppose scale features $\tilde{\psi} \triangleq b\psi$

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

$$\tilde{w}(\tilde{\alpha}) = \frac{1}{n} \sum_{i=1}^n \tilde{\alpha}_i \tilde{\psi}_i(\tilde{y})$$

$$\text{Let } \tilde{\lambda} \triangleq b^2 \lambda \quad \tilde{w}(\tilde{\alpha}) = \frac{1}{b^2 n} \sum_{i=1}^n \tilde{\alpha}_i \tilde{\psi}_i(\tilde{y}) \psi_i(y)$$

$$\text{If you use } \tilde{\alpha}_i^* \triangleq \alpha_i^* \Rightarrow \tilde{w}(\tilde{\alpha}) = \frac{w(\alpha)}{b}$$

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

$\Rightarrow \tilde{\alpha}_i^*$ is really optimal for new problem with $\tilde{\psi}$ & $\tilde{\lambda}$

4) similarly, can show $\tilde{b} = b\lambda \Rightarrow \tilde{\lambda} = \frac{\lambda}{b}$, get same solution

M³-net example (dual) : (getting small dual)

$$w(\alpha) = A\alpha = \sum_i A_i \alpha_i \quad \text{suppose } \psi(y) = \sum_c \psi_c(y_c)$$

$$\alpha_i \in \Delta_{\text{pr}_i}$$

$$(A\alpha)_i = \sum_{\tilde{y}} \alpha_i(\tilde{y}) \psi_i(\tilde{y}) = \sum_{\tilde{y}} \alpha_i(\tilde{y}) \sum_c \psi_{i,c}(\tilde{y}_c)$$

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

$$\sum_{\tilde{y}} \alpha_i(\tilde{y}) = \tilde{\alpha}_i$$

$$\tilde{\alpha}_i \in \Delta_{\text{pr}_i}$$

$$\alpha_i \in \Delta_{\text{pr}_i} \Rightarrow \mu_i \in M_i$$

marginal polytope

$$\text{thus } A_i \alpha_i = \tilde{A}_i \mu_i \text{ where}$$

"marginal variable"

$$(\tilde{A})_{i,c}(c, \tilde{y}_c) = \frac{\psi_{i,c}(\tilde{y}_c)}{\lambda_n}$$

\hookrightarrow # of columns is $\sum_c |\mathcal{S}_c|$

$$\text{similarly, suppose } l_i(\tilde{y}) = \sum_c l_{i,c}(\tilde{y}_c)$$

$$\text{define } \tilde{b}_{i,c}(\tilde{y}_c) \stackrel{\Delta}{=} \underbrace{l_{i,c}(\tilde{y}_c)}_n \Rightarrow \langle b_i, \alpha_i \rangle = \langle \tilde{b}_i, \mu_i \rangle$$

thus replace

$$\max_{\alpha \in \Delta_{\mathcal{M}_i}} -\frac{1}{2} \|A\alpha\|^2 + b^T \alpha \quad \text{with}$$

$$\max_{\mu \in \mathcal{M}_i} -\frac{1}{2} \|\tilde{A}\mu\|^2 + \tilde{b}^T \mu$$

→ this is a tractable size QP

if \mathcal{M}_i are tractable

if G_i is triangulated
then $M_i = L_i$ (local consistency polygons)
⇒ tractable

M^3 -net paper:

used "structured SMO algorithm"

block-coordinate ascent using pairs of variables on this QP

[similar to "pairwise FW"]

14h27

constraint generation alg.:

[Tschauder et al. JMLR 2005]

$$\text{want to solve } \min_{w, \xi} \frac{\|w\|^2}{2} + \frac{1}{n} \sum_i \xi_i \quad (\text{P}) \quad \max_{\alpha} -\frac{1}{2} \|A\alpha\|^2 + b^T \alpha \quad (\text{D})$$

D-slack version

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

constraints → $\prod_i |\mathcal{Y}_i|$
variables

$$\sum_i \xi_i$$

vs.
1-slack version

[ML 2009 paper]

$$\min_{w, \xi} \frac{\|w\|^2}{2} + \xi \quad (\text{P}) \quad \max_{\alpha} -\frac{1}{2} \|A\alpha\|^2 + b^T \alpha \quad (\text{D})$$

$$\text{s.t. } \xi \geq \frac{1}{n} \sum_{i=1}^n H_i(\tilde{y}_i; w) \quad (\forall \tilde{y}_i \in \mathcal{Y}_i)$$

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

$$\frac{1}{n} \sum_i \xi_i \geq \frac{1}{n} \sum_i \langle w, \sum_{j=1}^n \psi_j(\tilde{y}_i) \rangle \rightarrow \langle w, \sum_{i=1}^n \xi_i \psi_i(\tilde{y}_i) \rangle$$

of constraints → $\prod_i |\mathcal{Y}_i|$
(d) to store

$$w = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \psi_j(\tilde{y}_i) \quad \sum_i \alpha_i = 1$$

$\sum_{i=1}^d \|\hat{y}_i\|_1$ instead of $O(d)$ to store

Instead of $O(d \cdot n)$ storage \Rightarrow big memory saving
in n -slack formulation

$$\sum_{i=1}^n \|\hat{y}_i\|_1$$

n -slack SUMstruct alg.: cutting plane / constraint generation

Iterate solving QP with more & more constraints

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

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

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

↳ then resolve QP(w, g) with these constraints to get $w^{(t+1)}, g^{(t+1)}$ [e.g. CVXopt.]

stop when primal-dual gap $\leq \epsilon$

[in 2005 paper, show that alg. stop after $O(\frac{1}{\epsilon^2})$ iteration]

refined later [2009] to $O(\frac{1}{\epsilon})$ for 1-slack formulation

Frank-Wolfe algorithm

↳ for smooth constrained opt. [motivation in our context dual of SUMstruct min $\min_{\alpha_i \in \Delta(M)} \sum_i \|f_i\|_\infty - b_i \alpha_i$]

1940s: Simplex alg. to solve LPs

1956: Marguerite Frank & Phil Wolfe
→ non-linear opt. by iterating LPs

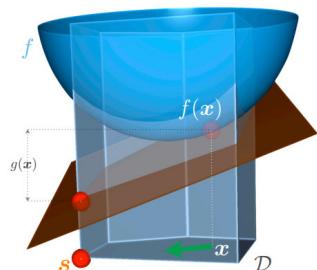
setup $\min f(x)$ • f is L -smooth

s.t. $x \in M$

• M is convex and bounded set

and assume we can solve efficiently $\min_{s \in M} \langle s, d \rangle$ for any d

$\min_{s \in M} \langle s, d \rangle$ by convexity



FW algorithm

Start with $x_0 \in M$
for $t = 0, \dots,$

FW corner

linear
minimization
oracle
(LMO)

$f(s) \geq f(x_t) + \langle \nabla f(x_t), s - x_t \rangle \quad \forall s \in M$
min \downarrow linear app. of f at x_t

for $t = 0, \dots$

\downarrow FW convex

\downarrow linear app. of f at x_t

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

\downarrow min KHS w.r.t. s

minimum finding oracle (LMO)

[let $g_t \stackrel{\Delta}{=} \langle s_t - x_t, -\nabla f(x_t) \rangle$ FW gap if $g_t \leq \epsilon$; output x_t]

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

$$= x_t + \gamma_t (s_t - x_t)$$

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

end
output x_t

step size choice:

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

big motivation for FW

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

adaptive choice:

$$\underbrace{\frac{g_t}{\|s_t - x_t\|^2}}_{\text{truncate at 1}} \text{ or } \frac{g_t}{C_f}$$

affine invariant const.