

Lecture 15 - cutting plane alg.

Thursday, March 11, 2021 13:27

- Today:
- more SVM struct properties
 - M²-net dual
 - cutting plane alg.
 - FW alg.

more properties of SVM struct dual

primal-dual gap

$$p(w) - d(\alpha) \geq 0 \quad \forall w, \alpha \text{ feasible}$$

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

Certificate of primal or dual suboptimality

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

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

$$= \frac{\lambda \|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n H_i(w(\alpha)) + \frac{\Delta}{2} \|w\|^2 - \frac{1}{n} \sum_{i=1}^n \alpha_i \langle \tilde{y}_i, \psi_i(\tilde{y}) \rangle$$

$$= \frac{1}{n} \sum_{i=1}^n \left(H_i(w) - \sum_{\tilde{y}} \alpha_i \langle \tilde{y}_i, \psi_i(\tilde{y}) \rangle \right) + \lambda \langle w, w(\alpha) \rangle$$

$$\quad \quad \quad \frac{1}{n} \sum_{i=1}^n \alpha_i \langle \tilde{y}_i, \psi_i(\tilde{y}) \rangle$$

$$\text{gap}(\alpha) = \frac{1}{n} \sum_{i=1}^n \left[H_i(w(\alpha)) - \sum_{\tilde{y}} \alpha_i \langle \tilde{y}_i, \psi_i(\tilde{y}; w(\alpha)) \rangle \right]$$

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

← use to get a bound on dual subopt. of α

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

Let $R_i \triangleq \max_{\tilde{y}} \|\psi_i(\tilde{y})\|_2$

then 1) $\|w^*\|_2 \leq \frac{1}{n} \sum_{i=1}^n \sum_{\tilde{y}} \alpha_i \|\psi_i(\tilde{y})\|_2$

$$R \triangleq \frac{1}{n} \sum_{i=1}^n R_i$$

$$\leq \frac{1}{n} \left(\sum_{i=1}^n R_i \right) = R$$

2) kernel trick & $\langle w(\alpha), \phi(x, y) \rangle = \frac{1}{n} \sum_{i=1}^n \sum_{\tilde{y}} \alpha_i \langle \tilde{y}_i, \psi_i(\tilde{y}) \rangle, \phi(x, y)$

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

$$\hookrightarrow k_{ij} = \langle \psi_i(y), \psi_j(y) \rangle = K(x^{(i)}, y^{(i)}; x, y) - K(x^{(i)}, y^{(i)}; x^{(i)}, y^{(i)})$$

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

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

$$\tilde{w}(\tilde{\alpha}) = \frac{1}{\tilde{\lambda}} \frac{1}{n} \sum_i \sum_{\tilde{y}} \tilde{\alpha}_i(\tilde{y}) \frac{\psi_i(\tilde{y})}{b \psi_i(\tilde{y})}$$

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

$$\tilde{w}(\tilde{\alpha}) = \frac{1}{b^2} \left[\frac{1}{\lambda} \frac{1}{n} \sum_i \sum_{\tilde{y}} \tilde{\alpha}_i(\tilde{y}) \psi_i(\tilde{y}) \right]$$

if you use $\tilde{\alpha}_i \triangleq \alpha_i^*$

$$\Rightarrow \tilde{w}(\tilde{\alpha}) = \frac{w^*}{b}$$

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

$\Rightarrow \tilde{\alpha}_i$ 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)

$$w(\alpha) = A\alpha = \sum_i A_i \alpha_i$$

$\alpha_i \in \Delta_{|S_i|}$

suppose $\psi(y) = \sum_c \psi_c(y_c)$

$$(An) A_i \alpha_i = \sum_{\tilde{y}} \alpha_i(\tilde{y}) \psi_i(\tilde{y}) = \sum_{\tilde{y}} \alpha_i(\tilde{y}) \sum_c \psi_{i,c}(\tilde{y}_c)$$

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

$\left[\sum_{\tilde{y}} \alpha_i(\tilde{y}) \right]$
 s.t.
 $\tilde{y}_c = \tilde{y}_c$
 $\cong M_{i,c}(\tilde{y}_c)$
 marginal variables

$$\alpha_i \in \Delta_{|S_i|} \Rightarrow \mu_i \in M_i$$

marginal polytope

thus $A_i \alpha_i = \tilde{A}_i \mu_i$ where $(\tilde{A})_{i,c} = \frac{\psi_{i,c}(\tilde{y}_c)}{\lambda_i}$

\hookrightarrow # of columns is $\sum_c |S_c|$

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

define $\tilde{b}_{i,c}(\tilde{y}_c) \triangleq \frac{l_{i,c}(\tilde{y}_c)}{\lambda_i} \Rightarrow \langle b_i, \alpha_i \rangle = \langle \tilde{b}_i, \mu_i(\alpha_i) \rangle$

⊛ thus replace

$$\max_{\alpha_i \in \Delta_{|S_i|}} -\frac{\lambda}{2} \|A\alpha\|^2 + b^T \alpha \quad \text{with} \quad \max_{\mu_i \in M_i} -\frac{\lambda}{2} \|\tilde{A}\mu\|^2 + \tilde{b}^T \mu$$

\rightarrow this is a tractable size GP if M_i is tractable

→ this is a tractable size QP
if M_i is tractable

M^3 not paper:
used "structured SVM algorithm"

if C_i is triangulated
then $M_i = L_i$ (local consistency)
polytope

block-coordinate ascent using
pair of variables on this QP
[similar to "pairwise FW"]

14h 29

constraint generation algo:

[Tsochantzakis & al. JMLR 2005]

want to solve $\min_{w, \xi} \frac{\lambda \|w\|^2}{2} + \frac{1}{n} \sum_i \xi_i$ (P)
s.t. $\xi_i \geq H_i(\tilde{y}_i; w) \forall \tilde{y}_i \in \mathcal{Y}_i$
 $\xi_i \geq 0$

$\left. \begin{array}{l} \text{exp \#} \\ \text{of constraints} \\ \sum_i |\mathcal{Y}_i| \end{array} \right\} \xrightarrow{\text{\# variables}}$

$\max -\frac{\lambda \|Aa\|^2}{2} + b^T a$
s.t. $a_i \in \Delta_i$

n-slack version

vs.
1-slack version
[ML 2009 paper]

$\min_{w, \xi} \frac{\lambda \|w\|^2}{2} + \xi$ (P)
s.t. $\xi \geq \frac{1}{n} \sum_{i=1}^n H_i(\tilde{y}_i; w) \quad (\forall \tilde{y}_i \in \mathcal{Y}_i)_{i=1, \dots, n}$

of constraints $\rightarrow \prod_i |\mathcal{Y}_i|$

$\left(\sum_{i=1}^n k_i \langle \tilde{y}_i \rangle \right) \leq w, \left(\sum_{i=1}^n \psi_i(\tilde{y}_i) \right)$
O(d) to store

(D)
 $\max -\frac{\lambda \|Aa\|^2}{2} + b^T a$
 $a \in \Delta \left(\prod_{i=1}^n \mathcal{Y}_i \right)$

$w^T a = \frac{1}{n} \sum_{i=1}^n \alpha(\tilde{y}_i; n) \left(\sum_{i=1}^n \psi_i(\tilde{y}_i) \right)$
 $\forall n \tilde{y}_i: n \in \mathcal{Y}_i$

instead of O(d·n) storage in n-slack formulation \Rightarrow big memory saving

n-slack SVM struct algo:

working plane / constraint generation

iterate solving QP with more & more constraints

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

2) repeat: for each i , find $\hat{y}_i^* = \arg \max_{\tilde{y}_i \in \mathcal{Y}_i} H_i(\tilde{y}_i; w^{(t)})$ [loss-augmented decoding]
• add $\xi_i \geq H_i(\hat{y}_i; w)$ constraint to QP (is not always true)

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

stop when primal-dual gap $\leq \epsilon$ } $O(n)$ time

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

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

Frank-Wolfe algorithm

↳ for smooth constrained opt. [motivation on an contest dual of SVM sheet $\min_{\alpha_i \in \Delta, \sum \alpha_i = 1} \frac{\sum \alpha_i \|A_i\|^2}{2} - b^T \alpha$]

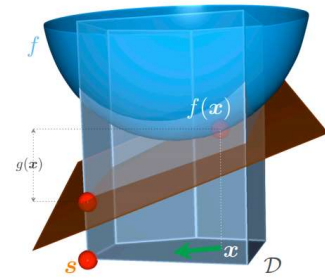
1940s: simplex alg. to solve LPs

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

Setup: $\min f(x)$
 $\text{st. } x \in M$

• f is L -smooth

• M is convex and bounded set



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

FW algorithm

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

FW corner

$$s_t = \arg \min_{s \in M} \langle s, \nabla f(x_t) \rangle$$

"linear minimization oracle" LMO

min RHS w.r.t s

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

linear app. of f at x_t

[let $g_t \triangleq \langle s_t - x_t, -\nabla f(x_t) \rangle$ FW gap if $g_t \leq \epsilon$; output x_t]

stopping criterion

$$x_{t+1} = (1 - \gamma_t) x_t + \gamma_t s_t \quad \gamma_t \in [0, 1] \text{ (convex combo)}$$

$$= x_t + \underbrace{\gamma_t}_{\alpha_t} (s_t - x_t)$$

end
 output x_t

step size choice: $\gamma_t = \begin{cases} \text{universal} & \frac{2}{t+2} \\ \text{line search} & \gamma_t = \arg \min_{\gamma \in [0, 1]} f(x_t + \gamma(s_t - x_t)) \end{cases}$

⊗ big motivation for FW

is LMO is often much cheaper than projections.

adaptive choice: $\left[\frac{g_t}{L \|s_t - x_t\|^2} \text{ or } \frac{g_t}{\epsilon} \right]$ truncated

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

adaptive choice : $\left[\begin{array}{c} g_t \\ \text{truncated} \\ \text{at } 1 \end{array} \right. \text{ or } \left. \begin{array}{c} g_t \\ \frac{g_t}{C} \end{array} \right]$
 affine
 invariant const.