

Lecture 15 - cutting plane alg.

Tuesday, March 10, 2020 14:21

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

more properties of SVM struct dual

primal dual gap:

$$\begin{aligned} \text{gap}(\alpha) &= p(w(\alpha), \xi(\alpha)) - d(\alpha) \geq 0 \\ &= \underbrace{p(w(\alpha), \xi(\alpha)) - p(w^*, \xi^*)}_{\text{primal subopt. } w(\alpha)} + \underbrace{p(w^*, \xi^*) - d(\alpha)}_{\text{dual subopt.}} \end{aligned}$$

Certificate of subopt. \otimes

$$\begin{aligned} & \frac{\lambda \|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n k_i(w) + \frac{\lambda \|w\|^2}{2} - \frac{1}{n} \sum_{i=1}^n \sum_{\tilde{y}} \alpha_i(\tilde{y}) l_i(\tilde{y}) \\ &= \frac{1}{n} \sum_{i=1}^n \left[H_i(w) - \sum_{\tilde{y}} \alpha_i(\tilde{y}) k_i(\tilde{y}) \right] + \lambda \langle w, w(\alpha) \rangle \end{aligned}$$

$$\boxed{\text{gap}(\alpha) = \frac{1}{n} \sum_{i=1}^n \left[H_i(w) - \sum_{\tilde{y}} \alpha_i(\tilde{y}) \underbrace{H_i(\tilde{y}; w)}_{k_i(\tilde{y}) - w^T \psi_i(\tilde{y})} \right]}$$

use to get a bound on dual subopt. of α

$$w(\alpha) = \frac{1}{\lambda n} \sum_{i=1}^n \left(\sum_{\tilde{y}} \alpha_i(\tilde{y}) \psi_i(\tilde{y}) \right)$$

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

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

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

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

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

$\phi(x_i, y_i) - \phi(x_i, \tilde{y})$

$k(x_i, y_i; x, y)$

$$\|w(\alpha)\|^2 \rightarrow \alpha^T K \alpha \quad - k(x_i, \tilde{y}_j; x, y)$$

$$\hookrightarrow k_{ij} = k(x_i, \tilde{y}_j) \triangleq \langle \psi_i(x), \psi_j(y) \rangle$$

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

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

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

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

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

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

ie. $\tilde{\alpha}^*$ is really optimal for new problem with $\tilde{\psi}$!

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

M²-net example (dual): (getting small dual)

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

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

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

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

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

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

↑
marginal polytope

$\triangleq \mu_{i,c}(\tilde{y}_c)$
"marginal variables"

thus $A_i \alpha_i = \tilde{A}_i \mu_i$ where $(\tilde{A}_i)_{:,c \in S_i} = \psi_{i,c}(y_c)$

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

similarly, suppose $l(y) = \sum_c l_{i,c}(y_c)$

define $\tilde{b}_{i,c}(y_c) \triangleq \frac{L_{i,c}(y_c)}{n}$
 $\langle b_i, \alpha_i \rangle = \langle \tilde{b}_i, M_i(\alpha_i) \rangle$

⊛ thus we can replace

$$\max_{\alpha_i \in \Delta_{|S_i|}} -\lambda \frac{\|A\alpha\|^2}{2} + b^T \alpha$$

$$\max_{\mu_i \in M_i} -\lambda \frac{\|\tilde{A}_i \mu\|^2}{2} + b_i^T \mu$$

→ this is a tractable size QP
 if M_i is tractable

M³-not paper

used "structured SMO" algorithm

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

∩ C_i is triangulated

then $M_i = L_i$ (local consistency polytope)

15h23

constraint generation alg.:

[Tschantzidis & al. JMLK 2005]

want to solve $\min_{w, \xi} \frac{\lambda \|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n \xi_i$ (P)

(D) $\max -\lambda \|A\alpha\|^2 + b^T \alpha$

n-slack version

s.t. $\xi_i \geq H_i(\tilde{y}_i, w) \quad \forall \tilde{y}_i \in S_i$
 $\xi_i \geq 0$

exp # of constraints $\leq |S_i|$

s.t. $\alpha_i \in \Delta_{|S_i|}$
 # variables

vs,

1-slack version

[ML 2009 paper]

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

s.t. $\xi \geq \frac{1}{n} \sum_{i=1}^n H_i(\tilde{y}_i, w) \quad (\forall \tilde{y}_i \in S_i)_{i=1, \dots, n}$

$\xi \geq \frac{1}{n} \sum_{i=1}^n \xi_i$

$\frac{1}{n} \sum_{i=1}^n \xi_i = \langle w, \frac{1}{n} \sum_{i=1}^n M_i(\tilde{y}_i) \rangle$

(D) $\max -\lambda \frac{\|A\alpha\|^2}{2} + b^T \alpha$

$\alpha \in \Delta_{\sum_{i=1}^n |S_i|}$
 $w(\alpha) = \frac{1}{n} \sum_{i=1}^n \alpha_i(\tilde{y}_{i,n})$
 $\sum_{i=1}^n \tilde{y}_{i,n} \in \sum_{i=1}^n S_i$
 $\sum_{i=1}^n \alpha_i(\tilde{y}_i)$

old) instead of $O(dn)$ in n-slack formulation \Rightarrow big memory saving

n-slack SVM sturt alg.:

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 $\hat{y}_i = \arg \max_{\hat{y} \in \mathcal{Y}_i} H_i(\hat{y}; w^{(i)})$ [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, ξ) with these constraints to get $w^{(i+1)}, \xi^{(i+1)}$ [e.g. using CVXopt.]

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

[in 2005, showed that QP step after $O(\frac{1}{\epsilon^2})$ iterations]

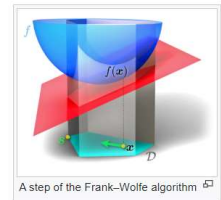
refined later to $O(\frac{1}{\epsilon})$ ← at least for 1-block version [2009]

Frank-Wolfe algorithm

\hookrightarrow for smooth constrained opt. [motivation dual of SVM short $\min_{\alpha_i \in \Delta_{|S|}} \lambda \sum \|A_i\|^2 - b^T \alpha$ in our context]

1940s: simplex alg. to solve LPs

1956: Manguerra Frank & Phil Wolfe
 \rightarrow non-linear opt. by iterated LPs



setup: $\min f(x)$
 s.t. $x \in M$

- f is L -smooth i.e. (∇f is L -Lipschitz) and f is convex
- M is convex and bounded set

FW algorithm

start with $x_0 \in M$

for $t=0, \dots$

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

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

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

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

by convexity
 $f(s) \geq f(x_t) + \langle \nabla f(x_t), s - x_t \rangle \quad \forall s \in M$
 "linear minimization oracle" LMO
 linear approx of f at x_t

min. RHS w.r. to s

stopping criterion

$$x_{t+1} = (1-\alpha_t)x_t + \alpha_t z_t$$

$$\alpha_t \in [0,1] \quad (\text{convex combo})$$

$$= x_t + \alpha_t \underbrace{(z_t - x_t)}_{d_t}$$

end
output x_t

step-size choice: $\alpha_t = \begin{cases} \text{universal choice } \frac{2}{t+2} \\ \text{line search } \alpha_t = \arg \min_{\delta \in [0,1]} f(x_t + \delta(z_t - x_t)) \\ \text{adaptive: } \left[\frac{\alpha_t}{\|z_t - x_t\|^2} \text{ or } \frac{\alpha_t}{\|g_t\|} \right] \text{ truncate at 1} \end{cases}$

offline unconstrained const.

⊛ big motivation for FW

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

properties:

$$1) f(x_t) - \underbrace{\min_{x \in M} f(x)}_{f^*} \leq O\left(\frac{1}{t}\right)$$

$$2) \text{FW-gap } g_t \geq f(x_t) - f^* \rightarrow \text{certificate of subopt.}$$

$$\min_{s \in t} g_s \leq O\left(\frac{1}{t}\right) \quad [\text{i.e. we will stop in } O\left(\frac{1}{\epsilon}\right) \text{ iterations?}]$$

$$3) x_t = p_0^+ x_0 + \sum_{u=1}^t p_u^+ s_{u-1}$$

where $\sum_{u=0}^t p_u^+ = 1$
 $p_u^+ \geq 0$

$\rightarrow x_t$ has "sparse" expansion in terms of the FW-comers $\{s_u\}_{u=1}^{t-1}$

⊛⊛ "sparse method"
 \rightarrow popular in ML

\rightarrow see later how to apply on dual of SVM struct

4) FW is affine covariant (like Newton's method)

5) there is $\Omega\left(\frac{1}{t}\right)$ lower bound for FW-like methods for $t \leq d$