

- today:
- more SVM struct properties
 - M³-net dual
 - cutting plane alg.
 - FW alg.

more properties of SVM struct dual

certificate of primal or dual suboptimality

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

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

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

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

$$\begin{aligned} & \frac{\lambda \|w(\alpha)\|^2}{2} + \frac{1}{n} \sum_{i=1}^n H_i(w(\alpha)) + \frac{\lambda \|w(\alpha)\|^2}{2} - \frac{1}{n} \sum_{i, \tilde{y}} \alpha_i(\tilde{y}) \ell_i(\tilde{y}) \\ &= \frac{1}{n} \sum_{i=1}^n \left(H_i(w) - \sum_{\tilde{y}} \alpha_i(\tilde{y}) \ell_i(\tilde{y}) \right) + \frac{1}{n} \sum_{i, \tilde{y}} \alpha_i(\tilde{y}) \psi_i(\tilde{y}) \\ & \qquad \qquad \qquad \sum_{\tilde{y}} (-H_i(\tilde{y}; w(\alpha))) \end{aligned}$$

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

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

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

then 1) $\|w\|_2^2 \leq \frac{1}{n} \sum_{i=1}^n \sum_{\tilde{y}} \alpha_i(\tilde{y}) \|\psi_i(\tilde{y})\|_2$

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

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

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

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

$$K(x^{(i)}, y^{(i)}; x, y) = \langle \psi_i(x^{(i)}), \psi_i(y) \rangle$$

$$\|w(\alpha)\|^2 \rightarrow \alpha^T K \alpha \rightarrow k_{i, y; j, y} \triangleq \langle \psi_i(y), \psi_j(y) \rangle$$

1) ... with feature $\tilde{w} \triangleq w(\alpha)$

$$\hookrightarrow K_{i,j}(y) = \langle \psi_i(y), \psi_j(y) \rangle$$

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

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

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

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

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

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

$\Rightarrow \tilde{\alpha}_i$ is really optimal for new problem with $\tilde{\psi}$ & $\tilde{\lambda}$ = 0 as w^* is optimal

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

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

$$w(\alpha) = A\alpha = \sum_{\alpha_i \in \Delta_{|S_i|}} A_i \alpha_i$$

suppose $\Psi(y) = \sum_c \Psi_c(\tilde{y}_c)$

$$\begin{aligned} (\lambda n) A_i \alpha_i &= \sum_{\tilde{y}} \alpha_i(\tilde{y}) \psi_i(\tilde{y}) = \sum_{\tilde{y}} \alpha_i(\tilde{y}) \sum_c \psi_{i,c}(\tilde{y}_c) \\ &= \sum_c \sum_{\tilde{y}_c} \psi_{i,c}(\tilde{y}_c) \left[\sum_{\tilde{y}} \alpha_i(\tilde{y}) \right] \end{aligned}$$

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

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

marginal polytope

thus $A_i \alpha_i = \tilde{A}_i \mu_i$ where $(\tilde{A})_{:,c} = \psi_{i,c}(\tilde{y}_c)$ "marginal variable"

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

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

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

⊗ thus replace

$$\max_{\alpha_i \in \Delta(\mathcal{Y}_i)} \frac{-\lambda \|A\alpha\|^2 + b^T \alpha}{2} \quad \text{with} \quad \boxed{\max_{\mu_i \in M_i} \frac{-\lambda \|\tilde{A}\mu\|^2 + \tilde{b}^T \mu}{2}}$$

→ this is a tractable subproblem if M_i is tractable

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

M³-net paper: sequential minimal optimization
used "structured SMO algorithm"

block-coordinate ascent on
pair of variables on this QP
[similar to "pairwise FW"] (see next classes)

10h 28

constraint generation alg.:

[Tsochantzidis & el. JMLR 2005]

want to solve $\frac{\lambda \|w\|^2}{2} + \sum_{i=1}^n \xi_i$ (P) $\max_{\alpha_i \in \Delta(\mathcal{Y}_i)} \frac{-\lambda \|A\alpha\|^2 + b^T \alpha}{2}$ (D)

st. $\xi_i \geq H_i(y_i; w) \forall y_i \in \mathcal{Y}_i$ } exp # of constraints $\sum_i |\mathcal{Y}_i|$ $\alpha_i \in \Delta(\mathcal{Y}_i)$ # variables

$\xi_i \geq 0$

n-slack version

vs. 1-slack version
[ML 2001 paper]

(P) $\min_{w, \xi} \frac{\lambda \|w\|^2}{2} + \xi$
 $\xi \geq \frac{1}{n} \sum_{i=1}^n H_i(y_i; w) \quad (y_i \in \mathcal{Y}_i)_{i=1, \dots, n}$
of constraints $\prod_i |\mathcal{Y}_i|$

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

$\frac{1}{n} \sum_{i=1}^n H(y_i) - \frac{1}{n} \langle w, \sum_{i=1}^n \gamma_i(y_i) \rangle$
(d) to store

$w(\alpha) = \frac{1}{n} \sum_{i=1}^n \alpha_i \gamma_i(y_i) \left(\sum_{i=1}^n \alpha_i \gamma_i(y_i) \right)$
 $\forall \alpha, y_i \in \mathcal{Y}_i$

instead of $O(d \cdot n)$ storage in 1-slack formulation \Rightarrow big memory saving

n-slack SVM struct alg.: cutting 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 $\hat{y}_i = \underset{y \in \mathcal{F}_i}{\text{argmax}} H_i(y; w^{(i)})$ [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 $w^{(i+1)}, \xi^{(i+1)}$

stop when primal-dual gap $\leq \epsilon$ [e.g. CVXopt]

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

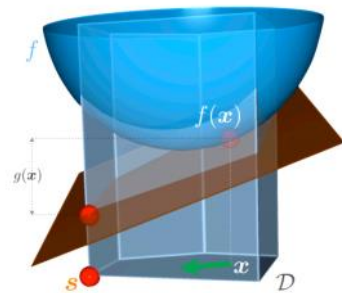
refined later [2009] to $O(\frac{1}{\epsilon})$ for 1-slack formulation

Frank-Wolfe algorithm

\hookrightarrow for smooth constrained opt. [motivation dual of SVM struct $\min_{x \in A} \frac{1}{2} \|Ax\|^2 - b^T x$ in con context]

1940s: simplex alg. to solve LPs

1956: Margerite Frank & Phil Wolfe \rightarrow non-linear opt. by iterating LPs



setup $\min f(x)$ s.t. $x \in M$
 • f is L -smooth
 • M is convex and bounded set

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

"linear minimization oracle" LMO

by convexity

$f(s) \geq f(x_t) + \langle \nabla f(x_t), s - x_t \rangle$ $\forall s \in M$
 $f(z^*) \geq$ linear app. of f at x_t

FW algorithm

start with $x_0 \in M$
 for $t=0, \dots$, compute $s_t = \underset{s \in M}{\text{argmin}} \langle s, \nabla f(x_t) \rangle$

stopping criterion

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

$x_t = (1-\alpha)x_t + \alpha s_t$ [convex combo]

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

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

end
output x_t

step size choice: $\alpha_t = \begin{cases} \text{unwind} & \frac{2}{t+2} \\ \text{line search} & \alpha_t = \underset{\alpha \in [0,1]}{\text{argmin}} f(x_t + \alpha(s_t - x_t)) \end{cases}$

⊗ big motivation for FW

↳ LMO is often much cheaper than projections

and cheap for many structured M appearing in ML

adaptive choice: $\left[\frac{g_t}{L \|s_t - x_t\|^2} \text{ or } \frac{g_t}{C_S} \right]$
 truncate at 1
 affine invariant constant