

today: structured SVM optimization

other approaches to optimize SVM struct

(UP)
unconstrained primal

$$\min_w \frac{\lambda \|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n H_i(w)$$

[unconstrained]
non-smooth

(PQP)
primal QP
→ quadratic program

$$\min_{w, \xi_i} \frac{\lambda \|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n \xi_i$$

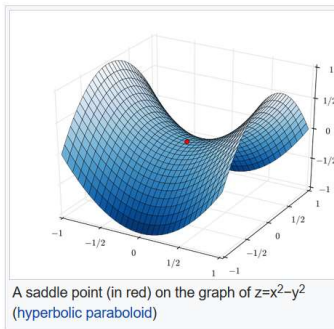
[constrained formulation]

$$\xi_i \geq H_i(w; y_i) \quad \forall y_i \in \mathcal{Y}_i; \xi_i$$

(smooth) convex QP
with exp. # of
linear constraints

1) generic approach to use convexity of loss-augmented decoding: [Taskan & al. ICLR 2005]

idea: here, we suppose that loss-augmented decoding can be expressed as a "compact" maximization problem of a concave fct.



$$i.e. H_i(w) = \max_{\tilde{y} \in \mathcal{Y}_i} \ell(\tilde{y}) - \langle w, \psi(\tilde{y}) \rangle = \max_{z \in \mathcal{Z}} g_i(w; z)$$

↑ discrete

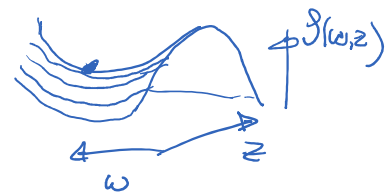
\mathcal{Z}_i ↑ convex set
where g_i is concave in z
and convex in w

\mathcal{Z} : • should not depend on w
• $\forall i$ have tractable description

a) saddle point formulation:

$$\min_w \max_{z_i \in \mathcal{Z}_i} \frac{\lambda \|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n g_i(w; z_i)$$

$\min_w \max_z \mathcal{L}(w, z)$
convex in w
concave in z
convex-concave
saddle point problem



(under reg. conditions $\min_w \max_z = \max_z \min_w \rightarrow$ "saddle point")

In general:

$$\min_w \max_z \geq \max_z \min_w$$

$$\forall z \quad \mathcal{L}(w^*, z) \leq \mathcal{L}(w^*, z^*) \leq \mathcal{L}(w, z^*) \quad \forall w$$

critical dependence
 $w^* \in \arg \min_w \mathcal{L}(w, z^*)$
 $z^* \in \arg \max_z \mathcal{L}(w^*, z)$

dependence $\rightarrow z^* \in \arg \max_z g(w^*, z)$

\Rightarrow might not exist? [but always exist for convex-concave + reg. conditions] e.g. 1-coercive

Standard alg:

\rightarrow converges $O(\frac{1}{t})$ for convex-concave game

lookahead step

extragradient algorithm:

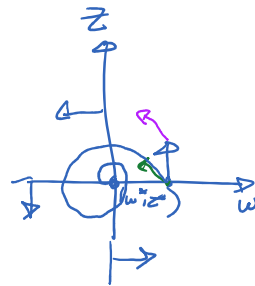
$$\begin{pmatrix} \tilde{w}_{t+1} \\ \tilde{z}_{t+1} \end{pmatrix} = \begin{pmatrix} w_t \\ z_t \end{pmatrix} + \alpha_t \begin{pmatrix} -\nabla_w g(w_t, z_t) \\ \nabla_z g(w_t, z_t) \end{pmatrix}$$

$$\begin{pmatrix} w_{t+1} \\ z_{t+1} \end{pmatrix} = \begin{pmatrix} w_t \\ z_t \end{pmatrix} + \alpha_t \begin{pmatrix} -\nabla_w g(\tilde{w}_{t+1}, \tilde{z}_{t+1}) \\ \nabla_z g(\tilde{w}_{t+1}, \tilde{z}_{t+1}) \end{pmatrix}$$

$$\begin{aligned} \tilde{x}_{t+1} &= x_t - \alpha_t F(x_t) \\ x_{t+1} &= x_t - \alpha_t F(\tilde{x}_{t+1}) \end{aligned} \quad x = \begin{pmatrix} w \\ z \end{pmatrix}$$

1st-order approx. to implicit method:

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



$$g(w, z) = \langle w, z \rangle$$

$$\begin{pmatrix} -\nabla_w g \\ \nabla_z g \end{pmatrix} = \begin{pmatrix} -z \\ w \end{pmatrix}$$

applied to structured SVM [Taskan et al. ICLR 2006]

b) small "complicated" Cp formulation (for structured SVM)

$$H_i(w) = \max_{z_i \in Z_i} g_i(w; z_i) = \min_{v_i \in V_i(w)} \tilde{g}_i(w; v_i)$$

use strong duality

\rightarrow convex dual of $\max_{z_i \in Z_i} g_i(w; z_i)$
dual variables

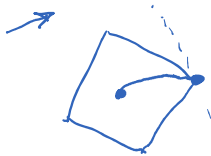
obtain: $\min_{w \in W} \min_{\substack{v_i \in V_i(w) \\ v_i}} \frac{\Delta \|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n \tilde{G}_i(w; v_i)$

if \tilde{G}_i is jointly convex in w, v_i we get a "tractable" convex min problem

\rightarrow can solve with favorite convex min alg

if d is not too big, use interior point solvers

e.g. Mosek, Cplex (commercial)
CVXopt (free python)



14th 35

examples of $g_i(w; z_i)$

I) word alignment:

Eng French
 0 0 3 0 1 0
 0 0 4 0 0
 0 0

\rightarrow features on min (F. r.p.)

I) word alignment:



features on pair (x_k^E, x_l^F)

recall that score $s(x, y; w) = \sum_{k,l} y_{k,l} [w^T \phi(x_k^E, x_l^F)]$

let $y \in \{0, 1\}^{L_E \times L_F}$

let matrix F be $\begin{bmatrix} \dots & \phi(x_k^E, x_l^F) & \dots \end{bmatrix}$ $d \times (L_E \times L_F)$

$s(x, y; w) = w^T F y$

$s(x^{(i)}, \tilde{y}; w) = w^T F_i \tilde{y}$

decoding: $h_w(x^{(i)}) = \arg \max_{\tilde{y} \in \mathcal{Y}_i} s(x^{(i)}, \tilde{y}; w)$ \rightarrow $\max_{\substack{y_{k,l} \in \{0,1\} \\ y \in M_i}} w^T F_i y$ linear integer program

$M_i = \left\{ \sum_{k,l} y_{k,l} e_{k,l}^{(i)} \mid \sum_{k,l} y_{k,l} \leq 1, \sum_{k,l} y_{k,l} \leq 1 \right\}$
 # of constraints $\sum_{k,l} y_{k,l} \leq 1$
 $L_E + L_F$ constraints

matching constraints

$A: \begin{pmatrix} y_{11} & y_{12} & y_{22} & y_{21} \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{pmatrix}$

$M_i = \{z \mid Az \leq \mathbf{1}, z \geq 0\}$

⊗ here, turns out can remove the integer constraint to get a "relaxed LP" give the same opt. objective value

ie. relaxation is tight

Actually here, $M_i = \text{conv-hull}(\mathcal{Y}_i)$

reasons that relaxation is tight

a) write $z \in M_i$ as $Az \leq b, z \geq 0$

matrix A here is "totally unimodular"

which means any subdeterminant of A has value $\begin{cases} +1 \\ -1 \\ 0 \end{cases}$

\Rightarrow that if b has integer entries then all vertices of $\{z \mid Az \leq b, z \geq 0\}$

have integer coordinates

\Rightarrow relaxation is tight for any linear cost

idea: $\tilde{A} \tilde{z} \leq \tilde{b}$, a corner of this obtained by solving $\tilde{A}_I \tilde{z} = \tilde{b}_I$ for \tilde{A}_I : invertible

$$|I| = \dim(z) \quad \tilde{z} = A_i^{-1} b_i$$

↑ Cramer's rule: ratio of subdeter.
⇒ integer

Conclusion: can write decoding as $\max_{z \in M_i} w^T F_i z$

What about loss?

Hamming loss example

$$l(y, \tilde{y}) = \sum_{k,l} \mathbb{1}\{y_{k,l} \neq \tilde{y}_{k,l}\}$$

$$= \sum_{k,l} (y_{k,l} - \tilde{y}_{k,l})^2$$

$$(y_{k,l}^2 - 2y_{k,l}\tilde{y}_{k,l} + \tilde{y}_{k,l}^2) = (y_{k,l}^2 + (1-2y_{k,l})\tilde{y}_{k,l})$$

= $\tilde{y}_{k,l}$

$$l_i(\tilde{y}) = a_i + \underbrace{(1-2y^{(i)})}_{c_i}^T \tilde{y}$$

cool trick: $y^2 = y$

when $y \in \{0,1\}$

loss-augmented decoding:

$$\max_{\substack{\tilde{y} \in \{0,1\} \\ \tilde{y} \in M_i}} \underbrace{a_i + c_i^T \tilde{y}}_{l_i(\tilde{y})} - \frac{(w^T F_i y^{(i)} - w^T F_i \tilde{y})}{w^T \phi_i(\tilde{y})}$$

$$= a_i - w^T F_i y^{(i)} + \max_{\substack{\tilde{y} \in \{0,1\} \\ \tilde{y} \in M_i}} (F_i^T w + c_i)^T \tilde{y}$$

$$= \max_{z \in M_i} \underbrace{(F_i^T w + c_i)}_{\tilde{c}_i} z + a_i - w^T F_i y^{(i)}$$

$\triangleq g_i(w; z)$

LP duality:

$$\max_{\substack{A_i z \leq b_i \\ z \geq 0}} \tilde{c}_i^T z = \min_{\substack{A_i^T v \geq \tilde{c}_i \\ v \geq 0}} b_i^T v$$

$$\max_{z \in M_i} g_i(w; z) = \min_{v \in V_i(w)} \tilde{g}_i(w; v)$$

here $\tilde{c}_i \triangleq F_i^T w + c_i$

A_i is $2L \times L^2$

SVM struct objective becomes

$$\min_w \min_{\{v_i\}_{i=1}^n} \frac{1}{2} \lambda \|w\|^2 + \frac{1}{n} \sum_{i=1}^n [a_i - w^T F_i y_i] + b_i^T v_i$$

st. $A_i^T v_i \geq F_i^T w + c_i$ "small complicated QP"
 $v_i \geq 0$

$$A^T \begin{pmatrix} w \\ v \end{pmatrix} \leq b$$

compare with saddle pt. formulation

$$\min_w \max_{\{z_i\}_{i=1}^n} \frac{1}{2} \lambda \|w\|^2 + \frac{1}{n} \sum_{i=1}^n [a_i - w^T F_i y_i] + [(F_i^T w + c_i)^T z_i]$$

$z_i \in M_i$
 \uparrow simplex constraints