

Lecture 12 - structured SVM

Thursday, February 14, 2019 13:35

today : SVMstruct \rightarrow small opt formulation

saddle point optimization:

necessary stationarity conditions on (w^*, z^*)

$$\min_{w \in W} \max_{z \in Z} \ell(w, z)$$

$$\left\{ \begin{array}{l} \langle \nabla_w \ell(w^*, z^*), w - w^* \rangle \geq 0 \quad \forall w \in W \\ \langle \nabla_z \ell(w^*, z^*), z - z^* \rangle \geq 0 \quad \forall z \in Z \end{array} \right.$$

directional derivative
of ℓ with respect to w
in $w - w^*$ direction is non-neg

$$\text{Let } u \triangleq (w, z)$$

$$F: \mathbb{R}^P \rightarrow \mathbb{R}^P$$

$$F(u) \triangleq \begin{pmatrix} \nabla_w \ell(w, z) \\ -\nabla_z \ell(w, z) \end{pmatrix}$$

$$\boxed{\langle F(u^*), u - u^* \rangle \geq 0 \quad \forall u \in U = W \times Z}$$

"variational inequality"

generalizes:

- minimization
- saddle point
- equilibrium problems
- multi-player games or non zero sum 2-player game

extra gradient algorithm

$$\tilde{u}_{t+1} = u_t - \gamma_t F(u_t)$$

$$u_{t+1} = u_t - \gamma_t F(\tilde{u}_{t+1})$$

\rightarrow solves variational inequality

b) small "complicated" QP formulation (for structural SVM)

$$H_i(w) = \max_{z_i \in Z} g_i(w, z_i) \stackrel{\text{use duality}}{\Leftrightarrow} \min_{v_i \in V_i(w)} \underbrace{g_i(w; v_i)}_{\substack{\text{dual} \\ \text{variables}}} \xrightarrow{\text{convex dual}} \max_{z_i \in Z} g_i(w, z_i)$$

obtain: $\min_{w \in W} \min_{v_i, z_i \in V_i(w)} \frac{\alpha \|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n g_i(w; v_i)$

If \hat{g}_i is jointly convex in $w \setminus v_i$

we get a "tractable" convex min. problem

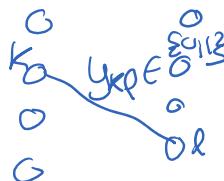
→ can solve with favorite convex min. alg.;

if $d > n$ not too big, use interior point solver

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

examples of $g_i(w; z_i)$

Eng. words Fr. words



I) word alignment:

features on pair (x_K^E, x_e^F)

recall that score $s(x_i y_i | w) = \sum_{k, e} y_{k, e} [w^T \psi(x_k^E, x_e^F)]$

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

Let $y \in \{0,1\}^{L_E \cdot L_F}$

Let matrix F be

$$\begin{bmatrix} \dots & \underbrace{\mathbf{f}_k}_{L_E} & \dots \end{bmatrix} \quad d \times (L_E \cdot L_F)$$

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

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

$$h_w(x^{(i)}) = \arg \max_{y \in \mathcal{Y}_i} s(x^{(i)}, y; w)$$

decoding / inference problem

$$\max w^T F_i y$$

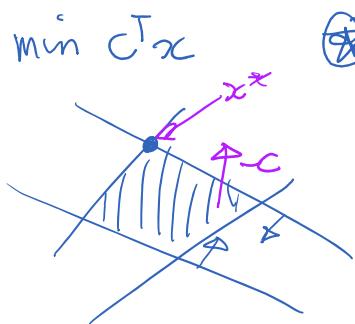
$$\begin{array}{l} y_k \in \{0,1\} \\ y \in M_i \end{array}$$

matching
constraints

$$M_i = \left\{ y \in \mathbb{R}^{L_E \cdot L_F} : \begin{array}{l} \sum y_{k,l} \leq 1, \quad 0 \leq y_{k,l} \leq 1 \\ \sum y_{k,l} \leq 1, \quad l \in \{1, \dots, L_F\} \end{array} \right\}$$

of constraints $L_E + L_F$

"linear integer program"



here, turns out can remove
the integer constraint to
get a "relaxed LP"

give the same
optimal objective value

[actually here, $M_i = \text{convex-hull}(\mathcal{Y}_i)$]

i.e. relaxation is tight

14h36

Reasons that relaxation is tight:

a) write $y \in M_i$ as $A y \leq b$

matrix A here is "totally unimodular"

a) write $y \in M_i$ as $\tilde{A}y \leq \tilde{b}$
 $y \geq 0$ matrix A here is "totally unimodular"
which means that
any subdeterminant of A has value $\begin{cases} +1 \\ -1 \\ 0 \end{cases}$

\Rightarrow that if b has integer entries

then all vertices of $\{y : \tilde{A}y \leq \tilde{b}, y \geq 0\}$
have integer coordinates

\Rightarrow relaxation is tight for any linear costs

Idea: $\tilde{A}\tilde{y} \leq \tilde{b}$, a corner of this polytope is obtained by solving

$$\tilde{A}_{I_1} \tilde{y} = \tilde{b}_I \text{ for } \tilde{A}_{I_1} \text{ invertible}$$

$$|I_1| = \dim(y)$$

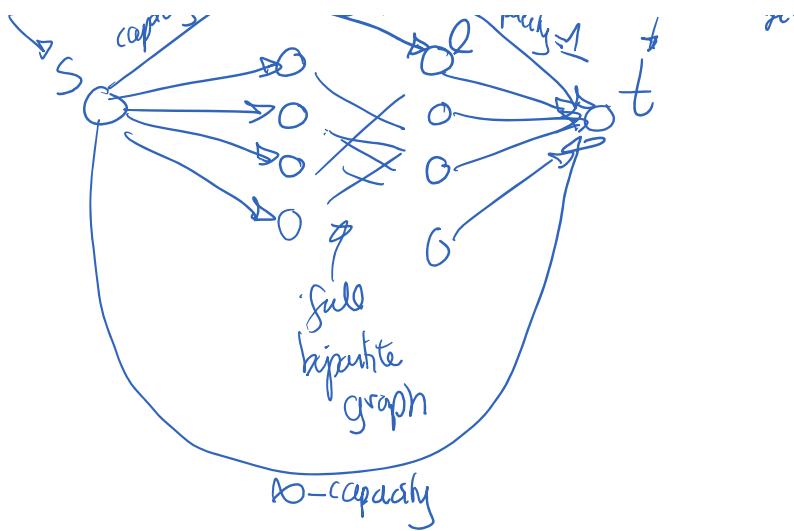
$$\tilde{y} = \tilde{A}_{I_1}^{-1} \tilde{b}_I$$

↑ Cramer's rule : ratio of subdetn.
 \Rightarrow integers etc...

b) equivalent to min cost network flow problem with integer capacities

[which has integer solutions by min-cut/max-flow thm.]





Conclusion: can write decoding as $\max_{z \in M^0} w^T F_i z$

what about loss?

Hammming Loss example:

$$\begin{aligned}
 l(y, \tilde{y}) &= \sum_{k,l} \mathbb{I}\{y_{k,l} \neq \tilde{y}_{k,l}\} \\
 &= \sum_{k,l} (y_{k,l} - \tilde{y}_{k,l})^2 \\
 &\quad \underbrace{\left(y_{k,l}^2 - 2y_{k,l}\tilde{y}_{k,l} + \tilde{y}_{k,l}^2 \right)}_{\text{Hamm Loss}}
 \end{aligned}$$

$$l_i(\tilde{y}) = a_i + \underbrace{(1 - 2y^{(i)})^T}_{c_i^0} \tilde{y}$$

cool trick: $y^2 = y$

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

Done now + 1

$$\underbrace{c_i}_{\tilde{c}_i}$$

Poss augmented
inference becomes:

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

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

$$= \max_{z \in M_i} \underbrace{(F_i^T w + c_i)^T z + a_i - w^T F_i y^{(i)}}_{\triangleq g_i(w; z)}$$

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

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

$$A_i \text{ is } (2L+L^2) \times L^2$$

SVM shurt objective becomes

$$\boxed{\begin{array}{ll} \min_w & \min_{\substack{\tilde{y} \in \mathcal{S}_{0,1}^n \\ A_i^T v_i \geq F_i^T w + c_i \\ v_i \geq 0}} \\ & \frac{\|w\|^2}{2} + \sum_{i=1}^n [a_i - w^T F_i y^{(i)}] + b^T v_i \end{array}}$$

"small complicated QP"

compare with saddle pt.
formulation

$$\boxed{\begin{array}{ll} \min_w & \max_{\substack{\tilde{y} \in \mathcal{S}_{0,1} \\ Z_i \in M_i}} \\ & \frac{\|w\|^2}{2} + \sum_{i=1}^n [a_i - w^T F_i y^{(i)}] + [(F_i^T w + c_i)^T Z_i] \end{array}}$$

simpler constraints



⊕ note: primal error vs. dual error not the same ☺

II) M³net example (RF score)

* graph $G = (V, E)$ $y_p, p \in V$ $y_C \stackrel{def}{=} (y_p)_{p \in C}$
 set of cliques ℓ

MRF model

$$\begin{aligned} p(y|x; w) &= \frac{1}{Z(x; w)} \prod_{C \in \ell} \pi_C(y_C; x; w) \\ &= \frac{1}{Z} \exp \left(\sum_C \underbrace{\log \pi_C(y_C; x; w)}_{S(x, y; w)} \right) \\ &\quad w^T \ell_C(y_C; x) \end{aligned}$$

Example: OCR
 → chain graph on y

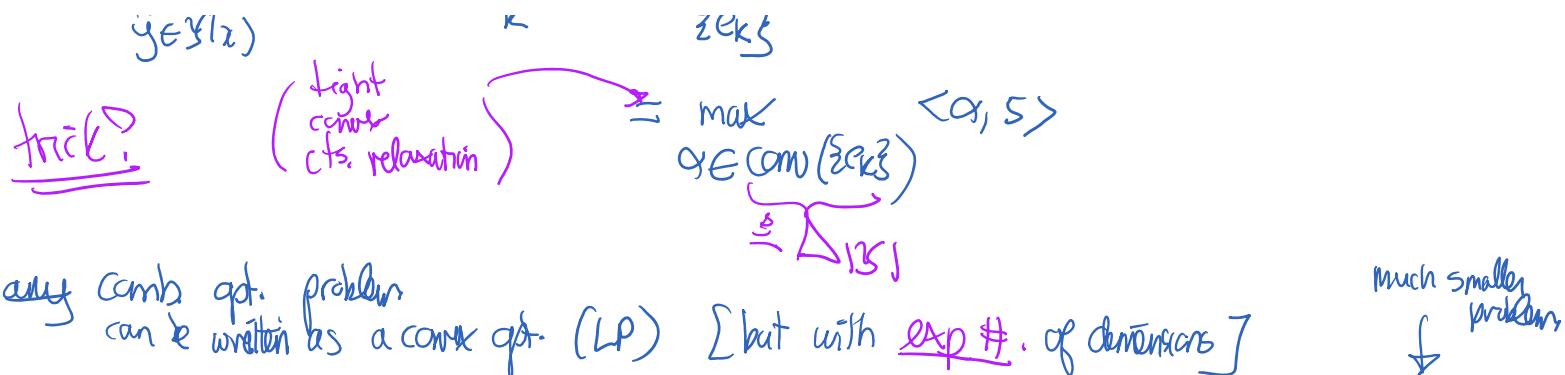
$$\text{inference: } \max_{\tilde{y} \in Y} \sum_{C \in \ell} w^T \ell_C(y_C; x)$$

goal: let's rewrite as a LP

* fix x, w , let s be a vector of size $|Y(x)|$ of scores $e_k = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \end{pmatrix}$ k^{th} position

$$\max_{\tilde{y} \in Y(x)} s(x, \tilde{y}; w) = \max_K s_K = \max_{\{e_k\}} \langle e_k, s \rangle$$

↑ right ↗ max / min ↙ \sim



④ insight of M³-net paper → use MRF structure to transform $\max_{\alpha \in \Delta_{12}} \langle \alpha, s \rangle = \max_{M \in \mathcal{M}_c} \sum_{\mu \in M} \langle \mu, s \rangle$