

today: structured SVM optimization

other approaches to optimize SVM struct

(UP)
unconstrained primal

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

[Convex
unconstrained
non-smooth]

(PQP)

primal QP
↳ quadratic problem

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

[constrained formulation]

$$\xi_i \geq H_i(w, \tilde{y}) \quad \forall \tilde{y} \in \mathcal{Y}_i \quad \forall i$$

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

i) generic approach to use convexity of loss augmented decoding: [Taskan & al. ICML 2006]

idea: here, we suppose that loss-augmented decoding

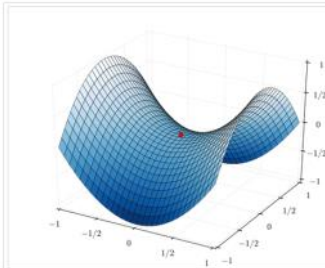
can be expressed as a "compact" maximization problem of a concave f_i

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

$\tilde{y} \in \mathcal{Y}_i$ discrete
 $z \in \mathcal{Z}$ discrete convex set
 discrete variable

where g_i is concave in z
and convex in w

\mathcal{Z} should not depend on w
" " have a tractable description



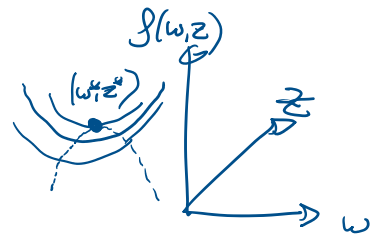
A saddle point (in red) on the graph of $z=x^2-y^2$ (hyperbolic paraboloid)

a) saddle point formulation:

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

$$\min_w \max_z \underbrace{f(w, z)}_{m(w)} \quad \begin{matrix} \text{convex in } w \\ \text{concave in } z \end{matrix}$$

convex-concave saddle pt. problems



(under reg. conditions $\min \max = \max \min$)
→ "saddle point"

$$\forall z f(w^*, z) \leq f(w^*, z^*) \leq f(w, z^*) \quad \forall w$$

$$w^* \in \arg \min_w f(w, z^*)$$

$$z^* \in \arg \max_z f(w^*, z)$$

in general:

$$\min_w \max_z f(w, z) \gg \max_z \min_w f(w, z)$$

$$z^* \in \arg \max_z f(w^*, z)$$

→ might not exist!

↳ but it always exists

for convex-concave + reg. condition
e.g. bounded convex constraints

standard alg. is

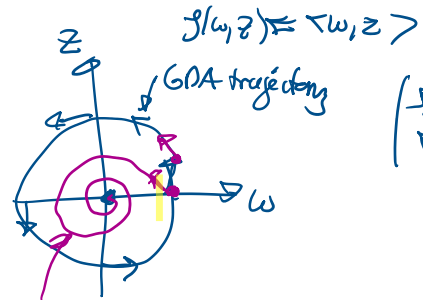
subgradient alg. converges $(\frac{1}{\epsilon})$ for convex-concave games

$$\begin{pmatrix} w_{t+1} \\ z_{t+1} \end{pmatrix} = \begin{pmatrix} w_t \\ z_t \end{pmatrix} + \delta_t \begin{pmatrix} -\nabla_w f(w_t, z_t) \\ \nabla_z f(w_t, z_t) \end{pmatrix}$$

↳ "look ahead step"

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

$$w \quad z \quad \dots \quad z \quad w \quad \dots$$



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

$$\tilde{x}_{t+1} = x_t - \delta_t F(x_t)$$

$$x = \begin{pmatrix} w \\ z \end{pmatrix}$$

$$x_{t+1} = x_t - \delta_t F(\tilde{x}_{t+1})$$

$$F(x) = \begin{pmatrix} \nabla_w f \\ -\nabla_z f \end{pmatrix}$$

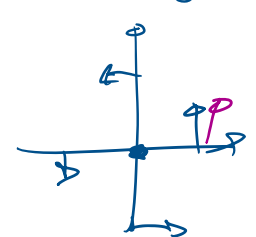
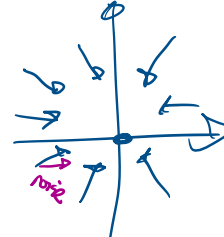
1st order approx. to the implicit method

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

$$x^* = x^* - \delta F(x^*)$$

stoch. min

vs stochastic games



subgradient to structural SVM [Taskan et al. JMLR 2006]

10/33

b) small 'complicated' QP formulation (for SVM struct)

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

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

using strong duality

dual variables

obtain $\min_w \min_{v_i \in V_i(w)} \frac{\|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n \tilde{g}_i(w; v_i)$

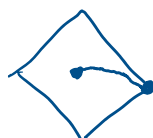
assume $v_i \in V_i(w)$
→ convex constraint in w

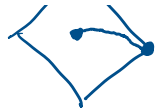
if \tilde{g}_i is jointly convex in w & v_i
we get a "tractable" convex min.

→ can solve with favorite convex min alg.

if d & n not too big, use interior point solver

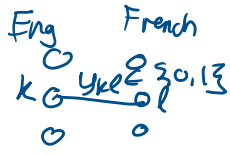
e.g. Mosek, Cplex (commercial)
CVXPopt (free, python)





Examples $Q_i(w; z_i)$

I) word alignment:



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

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

let $y \in \{0, 1\}^{L^E \cdot L^F}$

let matrix F be $d \times L^E \cdot L^F$

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

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



decoding: $hw(x^{(i)}) = \underset{y^i \in \mathcal{Y}_i}{\text{argmax}} s(x^{(i)}, y^i; w) \rightarrow \max_{\substack{y_{kl} \in \{0, 1\} \\ y \in M_i}} w^T F_i y$ Linear integer program

$M_i = \{ y \in \mathbb{R}^{L^E \cdot L^F} : \begin{cases} \sum_k y_{kl} \leq 1 \\ \sum_l y_{kl} \leq 1 \\ 0 \leq y_{kl} \leq 1 \end{cases} \}$

of constraints $(L^E + L^F)$

↑ matching constraints

$M_i = \{ z : \begin{cases} A z \leq 1 \\ z \geq 0 \end{cases} \}$

$A_i = \begin{pmatrix} y_{11} & y_{12} & y_{21} & y_{22} \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ \hline & & & I \end{pmatrix}$

⊛ here, turns out that can remove integer constraint to get a "relaxed LP" give the same opt. obj. value ie. relaxation is tight

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

$(2L^E + L^E) \times L^E$ matrix

why is relaxation tight?

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 $\begin{cases} \pm 1 \\ 0 \end{cases}$

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

\Rightarrow relaxation is tight for any linear cost

→ relaxation is tight for any linear cost

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

$$\tilde{A}_I z = \tilde{b}_I \text{ for } \tilde{A}_I: \text{ is invertible}$$

$$|I| = \dim(z)$$

$$\tilde{z}^* = \tilde{A}_I^{-1} b_I$$

use Cramer's rule ratio of subdeterminant

→ get integer values //

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

what about loss?

Hamming loss example

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

code trick:

$$y^2 = y \text{ when } y \in \{0, 1\}$$

$$= \sum_{k \in R} (y_{k \ell} - \tilde{y}_{k \ell})^2$$

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

$$l_i(\tilde{y}) = \sum_{k \in R} a_i + (1 - 2y_{k \ell}^{(i)})^T \tilde{y}$$

loss augmented decoding

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

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

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

LP duality:

$$\max_{\substack{Az \leq b \\ z \geq 0}} \tilde{c}^T z = \min_{\substack{A^T v \geq \tilde{c} \\ v \geq 0}} b^T v$$

$$\max_{z \in Z_i} g_i(w, z) = \min_{v \in V_i(w)} \hat{g}_i(w, z)$$

$$\text{here } \hat{c}_i \leq F_i^T w + c_i$$

$$z \geq 0$$

$$v \geq 0$$

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

$$A_i \text{ is } 2L \times L^2$$

(note $b_i = -1$)
here

Sum struct obj.
becomes (small QP)

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

$$\text{s.t. } A_i^T v_i \geq F_i^T w + c_i \quad \text{"small complicated" QP}$$
$$v_i \geq 0$$

compare with
saddle pt.
formulation

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

$$z_i \in M_i$$

→ simpler constraints