

IFT6132 reminder: fill survey <http://bit.ly/IFT6132-W24> ASAP!

today:
 • examples of structured prediction
 • structured perception & friends

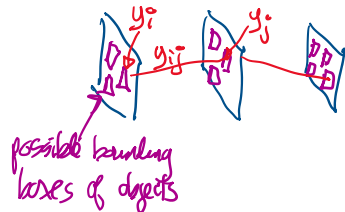
Examples:

I) word alignment (continuation)

here $x = (\underbrace{x_1^E, \dots, x_{L^E}^E}_{\text{English words}}; \underbrace{x_1^F, \dots, x_{L^F}^F}_{\text{French words}})$

$$\mathcal{Y}(x) = \{ y \in \{0,1\}^{L^E \times L^F} : \sum_j y_{ij} \leq 1 \forall i, \sum_i y_{ij} \leq 1 \forall j \}$$

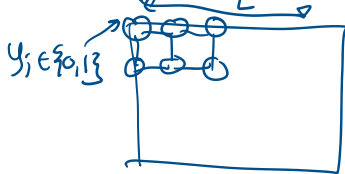
II) multi-object tracking:



$$y_i = \sum_j y_{ij}$$

encoding \rightarrow network flow

III) image segmentation:



$x =$ image of RGB values $L \times L$ pixels

$$\mathcal{Y}(x) = \{0,1\}^{L \times L}$$

\uparrow foreground
 \uparrow background

prediction model $h_w(x)$:

standard: $h_w(x) \triangleq \underset{y \in \mathcal{Y}(x)}{\text{argmax}} \left[\begin{array}{l} s(x,y;w) \text{] compatibility score of } y \text{ for } x \\ -E(x,y;w) \text{] energy fct. of } E \end{array} \right.$

linear model: $s(x,y;w) = \langle w, \underbrace{\varphi(x,y)}_{\text{"joint feature" vector}} \rangle$ $\varphi: X \times \mathcal{Y} \rightarrow \mathbb{R}^d$

word alignment: $\varphi(x,y) = \sum_{i,j} y_{ij} \underbrace{\varphi_{ij}(x_i^E, x_j^F)}_{\substack{\text{features defined} \\ \text{on a pair of English word } x_i^E \\ \text{French word } x_j^F}}$

- string edit distance (x_i^E, x_j^F)
- distance between i & j
- $\{x_i^E, x_j^F\}$ in dictionary etc...

$$s(x,y;w) = \langle w, \varphi(x,y) \rangle = \langle \dots, \langle w, \varphi_{ij}(x_i^E, x_j^F) \rangle, \dots \rangle$$

"French word x_j^F "

$$s(x, y; w) = \langle w, \phi(x, y) \rangle = \sum_{i,j} y_{ij} \langle w, \phi(x_i^F, x_j^F) \rangle$$

("score to match word i to j ")

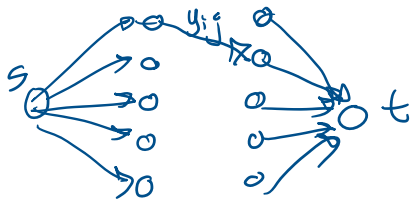
$$h_w(x) = \operatorname{argmax}_{y \in \mathcal{Y}(x)} s(x, y; w)$$

$$\rightarrow \max_y \sum_{i,j} y_{ij} S_{ij}(x)$$

$$\text{s.t. } \begin{cases} y_{ij} \in \{0, 1\} \\ \sum_j y_{ij} \leq 1 \quad \forall i \\ \sum_i y_{ij} \leq 1 \quad \forall j \end{cases}$$

integer LP
but
can be solved exactly
as min linear cost matching problem
eg. Hungarian alg.

or more generally
min cost network flow alg.



[side note: integer program with LP relaxation] (exact)

Learning w ?

I) structured perceptron:

- initialize w_0
- repeat for $t=0, \dots$

- sample i_t
- let $\hat{y}_t = h_{w_t}(x^{(i_t)}) = \operatorname{argmax}_{y \in \mathcal{Y}(x^{(i_t)})} \langle w, \phi(x^{(i_t)}, y) \rangle$
- $w_{t+1} = w_t + \eta \left(\underbrace{\langle \phi(x^{(i_t)}, y^{(i_t)}) \rangle}_{\text{stop size} \Rightarrow \text{boost score ground truth}} - \underbrace{\langle \phi(x^{(i_t)}, \hat{y}_t) \rangle}_{\text{penalize \hat{y}_t production}} \right)$

} "decoding oracle"

for stability: output $\hat{w}_T = \frac{1}{T+1} \sum_{t=0}^T w_t$ ← "Polyak averaging"

⊛ structured perceptron can be interpreted as

doing stochastic subgradient method (opt.) on the following non-smooth obj.:

$$\hat{J}(w) = \frac{1}{n} \sum_{i=1}^n J^{\text{percept}}(x^{(i)}, y^{(i)}; w)$$

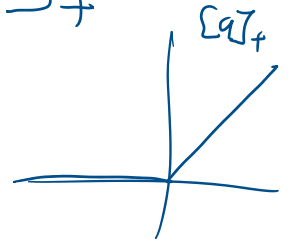
$$J^{\text{percept}}(x, y; w) \triangleq \left[\max_{\tilde{y} \in \mathcal{Y}} \langle w, \phi(x, \tilde{y}) \rangle - \langle w, \phi(x, y) \rangle \right]_+$$

, [9]₄

$$J^*(x, y; w) = \left[\max_{\tilde{y} \in \mathcal{Y}} \langle w, \phi(x; \tilde{y}) \rangle - \langle w, \phi(x; y) \rangle \right]_+$$

where $[a]_+ \triangleq \begin{cases} a & \text{if } a \geq 0 \\ 0 & \text{o.w.} \end{cases}$

(if $y^{(i)} \in \mathcal{Y}$; then this is always ≥ 0 and $[\cdot]_+$ is not needed)



10h32

II) conditional random field

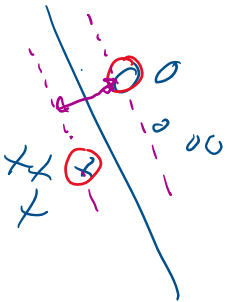
defined $p_w(y|x) \propto \exp(\langle w, \phi(x; y) \rangle) \frac{(\exp(\cdot))}{Z(x)}$

$$h_w(x) = \operatorname{argmax}_{y \in \mathcal{Y}(x)} p_w(y|x) = \operatorname{argmax}_y \langle w, \phi(x; y) \rangle$$

then maximum conditional likelihood on training set to learn \hat{w}

$$\hat{J}^{CRF}(w) = \frac{1}{n} \sum_{i=1}^n \hat{J}^{CRF}(x^{(i)}, y^{(i)}; w) + \underbrace{\lambda \|w\|^2}_{\text{regularizer}}$$

$$\begin{aligned} \hat{J}^{CRF}(x, y; w) &\triangleq -\log p_w(y|x) \\ &= \log \left(\underbrace{\sum_{\tilde{y}} \exp(\langle w, \phi(x; \tilde{y}) \rangle)}_{Z_w(x)} \right) - \langle w, \phi(x; y) \rangle \end{aligned}$$



Issues : $\cdot \ell(y, \tilde{y})$ doesn't appear in it

$\cdot \sum_{\tilde{y} \in \mathcal{Y}} \exp(\langle w, \phi(x; \tilde{y}) \rangle)$ can be difficult

eg. #P-complete for \mathcal{Y} = set of all matchings

III) structured SVM

intuition : want $s(x^{(i)}, y^{(i)}; w) \geq s(x^{(i)}, \tilde{y}; w) + \ell(y^{(i)}, \tilde{y}) \quad \forall \tilde{y} \in \mathcal{Y}_i \triangleq \mathcal{Y}(x^{(i)})$

min $\|w\|^2$ s.t. \uparrow "hard margin structured SVM"

(binary SVM : $y_i \in \{-1, +1\}$ $h_w(x) = \operatorname{sgn}(\langle w, \phi(x) \rangle)$)

$$y_i \langle w, \phi(x^{(i)}) \rangle \geq 1$$

hinge loss : $[1 - y_i \langle w, \phi(x^{(i)}) \rangle]_+$

soft-margin structured SVM : $R(w) \uparrow$ $\frac{1}{n} \sum_{i=1}^n \hat{J}^{svm}(x^{(i)}, y^{(i)}; w)$

soft-margin structured SVM

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

$$\xi_i + \langle w, \phi(x^{(i)}, y^{(i)}) \rangle \geq \langle w, \phi(x^{(i)}, \tilde{y}) \rangle + \ell(y^{(i)}, \tilde{y}) \quad \forall y \in \mathcal{Y}; \forall i$$

QP with an exponential # of constraints

equivalent (non-smooth) formulation:

$$\min_w \frac{\lambda \|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n f^{SVM}(x^{(i)}, y^{(i)}; w) \quad \xi_i \geq 0$$

where $f^{SVM}(x, y; w) = \left[\max_{\tilde{y} \in \mathcal{Y}(x)} \left[\langle w, \phi(x, \tilde{y}) \rangle + \ell(y, \tilde{y}) \right] - \langle w, \phi(x, y) \rangle \right]^+$

"structured hinge loss" (suppose that $y \in \mathcal{Y}(x)$)

not needed

OCR - optical character recognition example

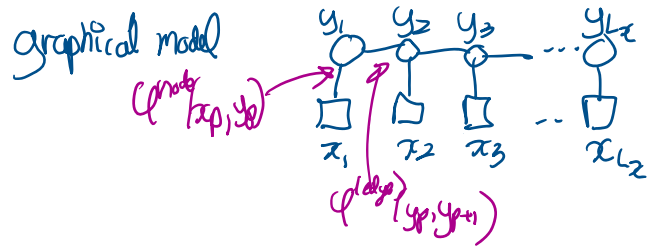
x : sequence of images of characters 
 y : characters \rightarrow B R A C E

$$x = (x_1, \dots, x_{L_x}) \quad x_p \in \{0, 1\}^{16 \times 8}$$

$$\mathcal{Y}(x) = \sum^{L_x} \mathcal{Z} \quad \mathcal{Z} = \{A, B, \dots, Z\}$$

in max-margin Markov network (M2-net) paper 8

$$\langle w, \phi(x, y) \rangle = \sum_{p=1}^{L_x} \langle w^{(node)}, \phi^{(node)}(x_p, y_p) \rangle + \sum_{p=1}^{L_x-1} \langle w^{(edge)}, \phi^{(edge)}(y_p, y_{p+1}) \rangle$$



$$p_w(y|x) = \frac{1}{Z_w(x)} \exp(\langle w, \phi(x, y) \rangle) = \frac{1}{Z_w(x)} \prod_{c \in \mathcal{C}} \psi_c(x, y_c)$$

where $\mathcal{C} = \{p, p+1\}$ (edges)

(chain structure) notation: $y_C \triangleq (y_i)_{i \in C}$

\Rightarrow can compute $\arg \max_{y \in \mathcal{Y}(x)} \langle w, \phi(x, y) \rangle$ using max product / sum alg. aka. Viterbi alg.

node: $\phi^{(node)}(x_p, y_p) = \begin{pmatrix} 0 \\ 0 \\ \text{vector}(x_p) \end{pmatrix}$ $\begin{matrix} y_p^m \text{ position} \\ \leftarrow (6 \times 8) \end{matrix}$ $\langle w, \phi(x_p, y_p) \rangle$

now: $\psi(x_p, y_p) \rightarrow \begin{pmatrix} 0 \\ \text{vector}(x_p) \\ 0 \\ 0 \end{pmatrix}$

\swarrow JP position
 \leftarrow 16×8
 $16 \times 8 \times 26$

$$\langle w, \psi(x_p, y_p) \rangle = 0 + 0 + 0 \dots \langle w_{y_p, x_p} \rangle + 0$$

\uparrow
 a template for y_p

w_a a
 w_b b

edge feature $\phi(y_p, y_{p+1}) = \begin{pmatrix} \downarrow (y_p, y_{p+1}) \\ \mathbb{1}\{y_p = y_p, y_{p+1} = y_{p+1}\} \end{pmatrix} 26^2$

$$\langle w^{(edge)}, \phi(y_p, y_{p+1}) \rangle = w_{y_p, y_{p+1}}^{(edge)}$$