

IFT6132: reminder: fill survey <http://bit.ly/IFT6132-W21> ASAP!

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

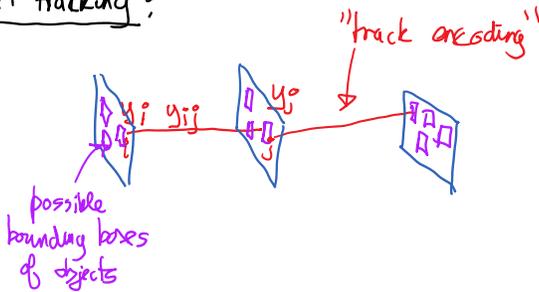
Examples:

I) word alignment (continuation)

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

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

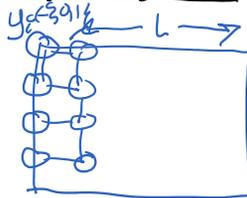
II) multi-object tracking:



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

encoding \rightarrow network flow

III) image segmentation:



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

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

background foreground

prediction model $h_w(x)$:

standard: $h_w(x) \triangleq \arg \max_{y \in Y(x)} s(x, y; w)$] compatibility score of y for x
 $-E(x, y; w)$] energy fct. E

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

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

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

- string edit distance (x_i^E, x_j^F)
- distance between i & j
- $\mathbb{1}\{x_i^E, x_j^F \text{ in dictionary}\}$
- etc.,

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

- match word x_i^E | etc.,

(score to match word i to j)

$$h_w(x) = \text{arg max}_{y \in \mathcal{Y}(x)} S(x, y; w)$$

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

$$\text{s.t. } y_{ij} \in \{0, 1\}$$

$$\sum_j y_{ij} = 1 \quad \forall i$$

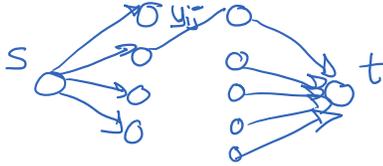
$$\sum_i y_{ij} \leq 1 \quad \forall j$$

can be solved exactly

(is min linear cost matching problem

e.g. Hungarian alg

or more generally min cost network flow alg.



[side note: integer program with LP relaxation]

14/23

Learning w

I) structured perceptron:

- initialize w_0

- repeat for $t=0, \dots$

- sample i_t

$$\bullet \text{ let } \hat{y}_t = h_{w_t}(x^{(i_t)}) = \text{arg max}_{y \in \mathcal{Y}(x^{(i_t)})} \langle w_t, \phi(x^{(i_t)}, y) \rangle$$

"decoding oracle"

$$\bullet w_{t+1} = w_t + \eta \left(\underbrace{\langle \phi(x^{(i_t)}, y^{(i_t)}) \rangle}_{\text{step size}} - \underbrace{\langle \phi(x^{(i_t)}, \hat{y}_t) \rangle}_{\text{penalize prediction}} \right)$$

\Rightarrow boost score ground truth

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.:

$$\hat{\mathcal{L}}(w) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}^{\text{percep.}}(x^{(i)}, y^{(i)}; w)$$

$$\mathcal{L}^{\text{percep.}}(x, u; w) \triangleq \left[\max \langle w, \phi(x, \tilde{u}) \rangle - \langle w, \phi(x, u) \rangle \right]$$

percept. $\mathcal{L}(x, y, w) \triangleq \left[\max_{y \in \mathcal{Y}} \langle w, \phi(x, \tilde{y}) \rangle - \langle w, \phi(x, y) \rangle \right]_+$
 where $[a]_+ \triangleq \begin{cases} a & \text{if } a \geq 0 \\ 0 & \text{otherwise} \end{cases}$

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

II) conditional random field

define $p_w(y|x) \propto \exp(\langle w, \phi(x, y) \rangle)$

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

then maximum conditional likelihood on training set to learn w
 $\mathcal{L}^{CRF}(w) = \underbrace{\frac{1}{n} \sum_{i=1}^n \mathcal{L}^{CRF}(x^{(i)}, y^{(i)}; w)}_{\text{regularizer}} + \lambda \|w\|^2$

$$\mathcal{L}^{CRF}(x, y; w) \triangleq -\log p_w(y|x)$$

$$= \log \left(\underbrace{\sum_{y' \in \mathcal{Y}} \exp(\langle w, \phi(x, y') \rangle)}_{Z_w(x)} \right) - \langle w, \phi(x, y) \rangle$$

issues: • $l(y, y')$ doesn't appear in it

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

e.g. #P-complete for $\mathcal{Y} = \text{set of all matchings}$

III) structural SVM

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

min $\|w\|^2$ s.t. \rightarrow

"hard margin structural SVM"

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

$$R(w) + \frac{1}{n} \sum_{i=1}^n \mathcal{L}^{SUM}(x^{(i)}, y^{(i)}; w) \quad \left\{ \begin{array}{l} \text{min } w, \xi \\ \xi_i + \langle w, \phi(x^{(i)}, y^{(i)}) \rangle \geq \langle w, \phi(x^{(i)}, \tilde{y}) \rangle + l(y^{(i)}, \tilde{y}) \quad \forall \tilde{y} \in \mathcal{Y}_i, \forall i \end{array} \right.$$

soft-margin structural SVM
 GP with an exponential # of constraints

equivalent (non-smooth) formulation:

$$\min_w \frac{\lambda \|w\|^2}{2} + \frac{1}{\lambda} \sum_{i=1}^n \mathcal{J}^{\text{sum}}(x^{(i)}, y^{(i)}; w)$$

"loss augmented decoding"

where $\mathcal{J}^{\text{sum}}(x, y; w) \triangleq \max_{\tilde{y} \in \mathcal{Y}(x)} [\langle w, \phi(x^{(i)}, \tilde{y}) \rangle + \ell(y^{(i)}, \tilde{y})] - \langle w, \phi(x, y) \rangle$

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

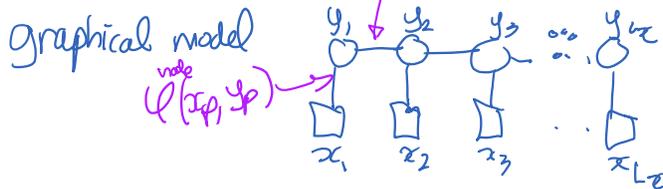
OCR - optical character recognition example:

x : sequence of images of characters  $x = (x_1, \dots, x_{L_x})$ $x_p \in \mathcal{X}_p$ $\mathcal{X}_p \subseteq \mathbb{R}^{16 \times 8}$

y  $y(x) = \sum_{p=1}^{L_x} x_p$ $\Sigma = \{A, \dots, Z\}$

in max-margin Markov network (M^3 -net) paper:

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



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

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

notation $y_c \triangleq (y_i)_{i \in c} \Rightarrow$ can compute $\arg \max_y \langle w, \phi(x, y) \rangle$

using max product alg. aka. Viterbi alg. or max sum

node: $\phi^{(\text{node})}(x_p, y_p) = \begin{pmatrix} 0 \\ 0 \\ \text{vector}(x_p) \\ 0 \\ 0 \end{pmatrix}$ $\leftarrow y_p^{\text{th}}$ position 16×8

$16 \times 8 \times 26$
characters



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

"a template for y_p "

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

\leftarrow (edge) $(1, 1) \dots (1, 1)$ (edge)

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