

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

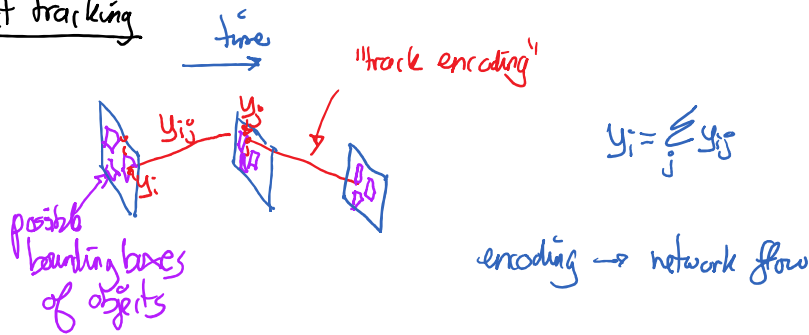
Examples:

I) word alignment (translation)

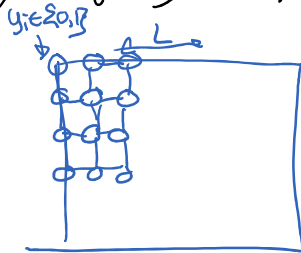
here $x = (\underbrace{x_1^E, \dots, x_{L_E}^E}_{\text{English words}}; \underbrace{x_1^F, \dots, x_{L_F}^F}_{\text{French words}})$

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

II) multi-object tracking



III) image segmentation



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

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

background foreground

prediction model $h_w(x)$

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

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

word alignment: $\phi(x, y) = \sum_{i,j} y_{ij} \underbrace{\phi(x_i^E, x_j^F)}_{\in \mathbb{R}^d}$
 features defined on a pair of English word x_i^E and French x_j^F

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

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

"score to match i to j "

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

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

s.t. $y_{ij} \in \{0, 1\}$
 $\sum_j y_{ij} = 1$
 $\sum_i y_{ij} \leq 1$

can be solved exactly
 linear
 as min cost matching problem
 e.g. Hungarian alg.
 or more generally,
 min cost network flow
 algorithm.



[side note: integer program with tight LP relaxation]

Learning w ?

I) structured perceptron:

- initialize w_0
- repeat for $t=0, \dots$,
 - sample i_t
 - let $\hat{y}_t = h_{w_t}(x^{(i_t)}) = \operatorname{argmax}_{y \in \mathcal{Y}(x^{(i_t)})} \langle w_t, \phi(x^{(i_t)}, y) \rangle$
 - $w_{t+1} = w_t + \eta \left(\underbrace{\phi(x^{(i_t)}, y^{(i_t)})}_{\text{step-size}} - \underbrace{\phi(x^{(i_t)}, \hat{y}_t)}_{\text{penalise prediction score}} \right)$

\Rightarrow boost score of ground truth

"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 optim. on the following non-smooth objective

... stochastic subgradient optim. on the following non-smooth objective

using stochastic subgradient optim. on the following non-smooth objective

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

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

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

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

14h31

II) conditional random fields

define $p_{\omega}(y|x) \propto \exp(\langle \omega, \phi(x, y) \rangle)$

$$h_{\omega}(x) = \arg \max_{y \in \mathcal{Y}} p_{\omega}(y|x) = \arg \max_y \langle \omega, \phi(x, y) \rangle$$

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

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

$$\begin{aligned} g^{\text{CRF}}(x, y; \omega) &\triangleq -\log p_{\omega}(y|x) \\ &= \log \left(\underbrace{\sum_{\tilde{y}} \exp(\langle \omega, \phi(x, \tilde{y}) \rangle)}_{Z_{\omega}(x)} \right) - \langle \omega, \phi(x, y) \rangle \end{aligned}$$

ISSUES: • (y, y') doesn't appear in it

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

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

III) structured SVM

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

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

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

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

soft-margin structured SVM: $R(w)$

GP with exponential # of constraints

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

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

equivalent (non-smooth) formulation:


$$\min_w \frac{\lambda \|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n g^{\text{SVM}}(x^{(i)}, y^{(i)}, w)$$


where $g^{\text{SVM}}(x, y, w) \triangleq \max_{\tilde{y} \in \mathcal{Y}(x)} [\langle w, \phi(x, \tilde{y}) \rangle + l(y, \tilde{y})] - \langle w, \phi(x, y) \rangle$

"loss-augmented decoding"

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

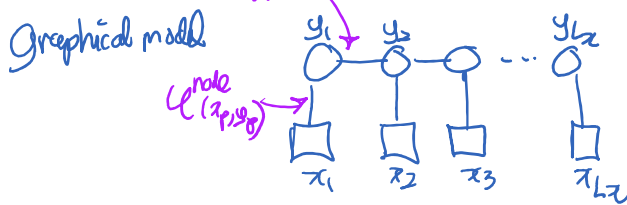
OCR - optical character recognition example

x : sequence of images of characters  $x = (x_1, \dots, x_{L_x})$ $x_p \in \{0, 1\}^{16 \times 8}$

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

in max margin Markov network (M³-net) paper:

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



$$p(y|x) = \frac{1}{Z_w(x)} \exp(\langle w, \phi(x, y) \rangle) = \frac{1}{Z_w(x)} \prod_{C \in \mathcal{C}} \pi_C(\phi_C(x, y_C)) \quad \text{here } \mathcal{C} = \{p, p+1\} \text{ (edges)}$$

notation: $y_c \triangleq (y_i)_{i \in c}$

\Rightarrow can compute $\arg \max_y \langle w, \phi(x, y) \rangle$
using max-product aka Viterbi alg. or max sum

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

16×8

$16 \times 8 \cdot 26$
of characters



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

attempt for character y_p

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

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