

Lecture 2 - examples

Thursday, January 12, 2023 1:28 PM

Reminder: do not forget to fill survey <http://bit.ly/IFT6132-W23> ASAP if not done!

- today:
- examples of structured prediction
 - structured perceptron { friends }

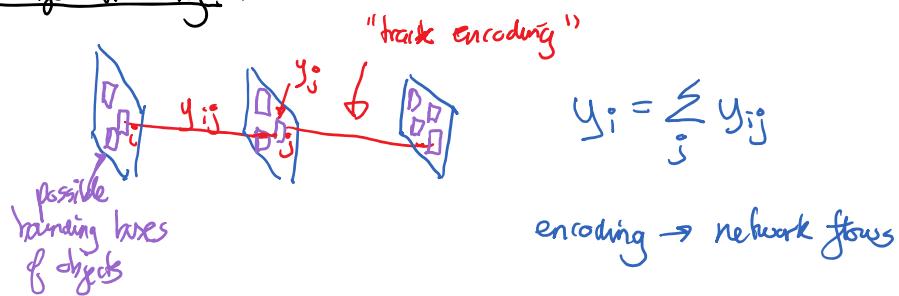
Examples

I) word alignment (continuation)

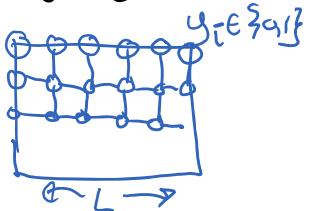
here $x = (\underbrace{x_1^E, x_2^E, \dots, x_{L_E}^E}_{\text{English words}}; \underbrace{x_1^F, \dots, x_{L_F}^F}_{\text{French words}})$

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

II) multi-object tracking:



III) image segmentation



$x = \text{image of RGB values } L \times L \text{ pixels}$

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

↑
foreground
background

prediction model $h_w(x)$

standard: $h_w(x) \triangleq \operatorname{argmax}_{y \in \mathcal{Y}(x)} s(x, y; w) - E(x, y; w)$

$s(x, y; w)$] compatibility score of y for x
 $-E(x, y; w)$] energy fn. E

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

word alignment: $\psi(x, y) = \sum_{i,j} y_{i,j} \psi(x_i^E, x_j^F)$

features defined
on a pair of
English word x_i^E
French " x_j^F

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

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

"score to match i to j"

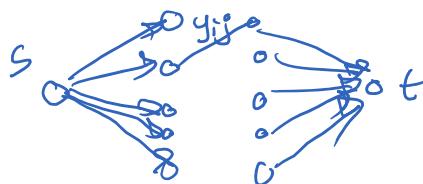
$$h_w(\tau) = \arg \max_{y \in \mathcal{Y}(\tau)} s(x, y; w) \rightsquigarrow \max_y \sum_{i,j} y_{i,j} \tilde{s}_{ij}(x)$$

s.t. $y_{i,j} \in \{0, 1\}$

$\sum_j y_{i,j} \leq 1 \forall i$

$\sum_i y_{i,j} \leq 1 \forall j$

} can be solved
exactly
as min linear cost
matching problem



[↑ denote: integer program
with LP relaxation]

e.g. Hungarian alg
or more generally
min cost network flow alg

14h29

Learning w?

good question about
tradeoff: tractable search
vs. more powerful models

I) structured perceptron:

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

$$\begin{cases} \cdot \text{sample } i_t \\ \cdot \text{let } \hat{y}_t = h_{w_t}(x^{(i_t)}) = \arg \max_{y \in \mathcal{Y}(x^{(i_t)})} \langle w_t, \psi(x^{(i_t)}, y) \rangle \\ \cdot w_{t+1} = w_t + \eta \left(\underbrace{\psi(x^{(i_t)}, y^{(i_t)}) - \psi(x^{(i_t)}, \hat{y}_t)}_{\text{step-size}} \right) \underbrace{\psi(x^{(i_t)}, \hat{y}_t)}_{\text{boost score of ground truth}} \end{cases}$$

penalize incorrect prediction

'decoding oracle'

$$\text{score}_{t+1}(\hat{y}) = \text{score}_t(\hat{y}) + \eta [\langle \psi(y^{i_t}), \psi(\hat{y}) \rangle - \langle \psi(\hat{y}_t), \psi(\hat{y}) \rangle]$$

\hat{y}_t

y^{i_t}

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

⊕ Structured perceptron can be interpreted as

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

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

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

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

II) conditional random field

define $p_\omega(y|x)$ or $\exp(\langle \omega, \varphi(x, y) \rangle)$

$$h_\omega(x) = \underset{y \in \mathcal{Y}(x)}{\operatorname{argmax}} p_\omega(y|x) = \underset{y}{\operatorname{argmax}} \langle \omega, \varphi(x, y) \rangle$$

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

$$\hat{S}^{\text{CRF}}(\omega) = \frac{1}{n} \sum_{i=1}^n S^{\text{CRF}}(x^{(i)}, y^{(i)}; \omega) + \lambda \underbrace{\| \omega \|_2^2}_{\text{regularizer}}$$

$$S^{\text{CRF}}(x, y; \omega) \triangleq -\log p_\omega(y|x)$$

$$= \log \underbrace{\left(\sum_{\tilde{y}} \exp(\langle \omega, \varphi(x, \tilde{y}) \rangle) \right)}_{Z(\omega|x)} - \langle \omega, \varphi(x, y) \rangle$$

Issues:

- $\varphi(y, y)$ doesn't appear in it

- $\sum_{\tilde{y} \in \mathcal{Y}} \exp(\langle \omega, \varphi(x, \tilde{y}) \rangle)$ is often intractable

e.g. #complete problem for \mathcal{Y} = set of all matchings!

III) structured sum

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

$\min \| \omega \|_2^2$ st. \rightarrow "hard margin structured sum"

binary S/M: $y \in \{-1, +1\}$, $h_\omega(x) = \operatorname{sgn}(\langle \omega, \varphi(x) \rangle)$

$$\sum_i \langle \omega, \varphi(x^{(i)}) \rangle \geq 1$$

$$\boxed{\begin{aligned} & L(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}) = \sum_{i=1}^n \ell(\mathbf{x}^{(i)}, y^{(i)}; \mathbf{w}) \\ & y_i \langle \mathbf{w}, \varphi(\mathbf{x}^{(i)}) \rangle \geq 1 \end{aligned}}$$

soft-margin structured SVM: $\mathcal{R}(\mathbf{w})$

$$\text{QP with an exponential of constraints} \quad \left\{ \begin{array}{l} \min_{\mathbf{w}, \xi} \frac{\lambda \|\mathbf{w}\|^2}{2} + \frac{1}{n} \sum_{i=1}^n \xi_i \\ \xi_i + \langle \mathbf{w}, \varphi(\mathbf{x}^{(i)}, y^{(i)}) \rangle \geq \langle \mathbf{w}, \varphi(\mathbf{x}^{(i)}, \tilde{y}) \rangle + l(y^{(i)}, \tilde{y}) \quad \forall \tilde{y} \in \mathcal{Y}_i, \forall i \end{array} \right.$$

equivalent (non-smooth) formulation:

$$\min_{\mathbf{w}} \frac{\lambda \|\mathbf{w}\|^2}{2} + \frac{1}{n} \sum_{i=1}^n g_{\text{SVM}}(\mathbf{x}^{(i)}, y^{(i)}; \mathbf{w})$$

where $\boxed{g_{\text{SVM}}(\mathbf{x}, y; \mathbf{w}) \triangleq \max_{\tilde{y} \in \mathcal{Y}(\mathbf{x})} [\langle \mathbf{w}, \varphi(\mathbf{x}, \tilde{y}) \rangle + l(y, \tilde{y})] - \langle \mathbf{w}, \varphi(\mathbf{x}, y) \rangle}$

loss augmented decoding

"structured hinge loss"

(suppose that $y \in \mathcal{Y}(x)$)