

## Lecture 2 - scribbles - EBM

Friday, January 26, 2018

13:35

today: - OCR example  
- energy based models

OCR - optical character recognition example

$x$ : sequence of  
images of characters



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

$y$ :

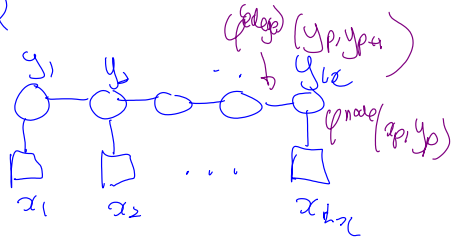
B R A C E

$$\gamma(x) = \sum_{\mathbf{y}} \quad \Sigma = \{A, B, \dots, Z\}$$

in Max-Margin Markov network (M<sup>3</sup>-net) paper:

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

graphical model



$$p(y|x) = \frac{1}{Z_w(x)} \exp(\langle w, \psi(x, y) \rangle) = \frac{1}{Z_w(x)} \prod_{C \in \mathcal{C}} \psi_C(x, y_C)$$

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

$$\text{here } \mathcal{C} = \{p, p+1\} \text{ (edge)}$$

$$\Rightarrow \text{can compute } \arg \max_y \langle w, \psi(x, y) \rangle$$

using max-product alg or  
viterbi alg or max-sum

feature function:

node:  $\ell(x_p, y_p) = \begin{pmatrix} \vdots \\ \text{vector}(x_p) \\ \vdots \end{pmatrix}$   $\leftarrow y_p$  in position

$16 \cdot 8 \cdot 26$   
# of characters

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

a template for character  $y_p$

$w_a$  

$w_b$  

edge feature:

$$\ell(y_p, y_{p+1}) = \begin{pmatrix} \vdots \\ f(y_p, y_{p+1}) \\ \vdots \end{pmatrix} \downarrow 26^2$$

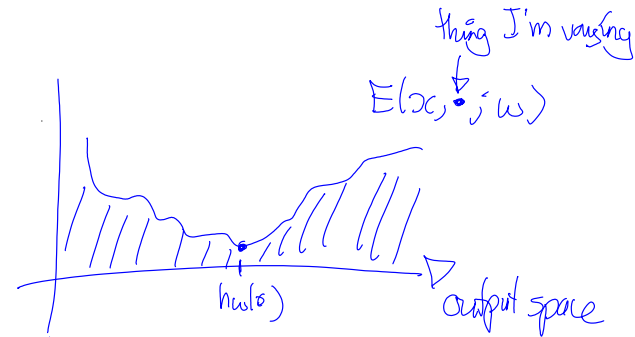
$$\mathbb{1}\{y_p = y_p, y_{p+1} = y_{p+1}\}$$

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

energy based methods: (LeCun et al. 2006)

model:  $h_w(x) = \underset{y \in \mathcal{Y}(x)}{\text{argmin}} E(x, y; w)$  "energy fct."

$= \underset{y \in \mathcal{Y}(x)}{\text{argmax}} S(x, y; w)$  "score / compatibility fct."



ingredients; modeling { 1) what is  $E(x, y; w)$ ? e.g.  $S(x, y; w) = \langle w, \ell(x, y) \rangle$

2) how do you compute  $\underset{y \in \mathcal{Y}(x)}{\text{argmin}} E(x, y; w)$ ? "inference" / "decoding"

learning { 3) how to evaluate  $E(x, y; w)$  on training set?  $\rightarrow$  surrogate loss  $\hat{J}(w)$   
in general:  $\hat{J}(x^{(i)}, y^{(i)}, E(\cdot, \cdot; \cdot))$  "loss functional"  
4) how to minimize  $\hat{J}(w)$  to learn  $w$ ?  $\rightarrow$  optimization tricks

4) how to minimize  $\hat{J}(w)$  to learn  $w$ ?  $\rightarrow$  optimization tricks

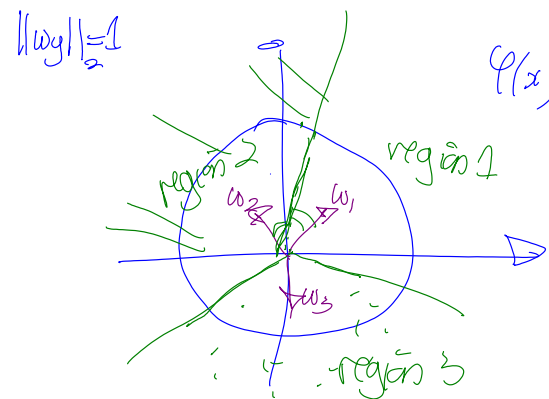
flat multiclass case:

"flat" setting  $h_w(x) = \arg \max_y \langle w_y, \phi(x) \rangle$

equivalent to:

$$\ell(x, y) = \begin{pmatrix} 0 \\ 0 \\ \phi(x) \\ 0 \\ 0 \end{pmatrix} \in \mathbb{R}^{d+2}$$

# of classes  
 $\leftarrow y^{\text{th}}$  position



surrogate losses:

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

$$h_w(x) \triangleq \arg \max_{y \in \mathcal{Y}} \langle w, \phi(x, y) \rangle$$

I) perception loss [Collins 2002 EMNLP]

$$\hat{J}(x, y; w) = \max_{\tilde{y} \in \mathcal{Y}(x)} S(x, \tilde{y}; w) - \overbrace{S(x, y; w)}^{\text{score of ground truth}}$$

$$S(x, y; w) = \langle w, \ell(x, y) \rangle$$

$$\max_{\tilde{y}} \langle w, \ell(x, \tilde{y}) - \ell(x, y) \rangle \geq 0$$

$\leftarrow$  by using  $\tilde{y} = y$

observations: 1) degenerate solution  $w=0$  or constant score over  $y$

does not converge in general

2) averaged perceptron algorithm:

- run constant step-size stochastic subgradient method on  $\hat{S}(w)$

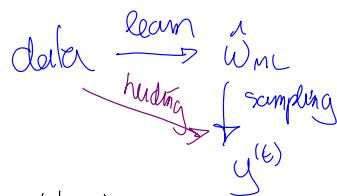
- output  $\hat{w}_T = \frac{1}{T} \sum_{t=0}^T w_t$  (Polyak averaging)  $\rightarrow$  will converge to  $w^* = 0$  (?)

comments: 1) Collins paper  $\rightarrow$  he gives error bound and generalization error guarantees for perceptron

2) connection with the "hiding" algorithm by Welling & al

"2nd way to learn"

see ICML 2012



II) log-loss (probabilistic interpretation)

suppose  $p(y|x; w) \propto \exp(\beta S(x, y; w))$   
 $\beta$  inverse temperature parameter

Boltzmann dist. in physics  
 $(\beta = \frac{1}{k_B T})$

MCL  $\rightarrow$  log-loss:

$$\mathcal{L}(x, y; w) = \underbrace{-\log p(y|x; w)}_{\text{rescaling}} = -\frac{1}{\beta} \log \left[ \frac{\exp(\beta S(x, y; w))}{\sum_y \exp(\beta S(x, \tilde{y}; w))} \right]$$

$Z_\beta(x; w)$  partition function

$$= \frac{1}{\beta} \log \left( \sum_y \exp(\beta S(x, \tilde{y}; w)) \right) - S(x, y; w)$$

$UN / \left( \sum_y \exp(S_y) \right)$

$$= \frac{1}{\beta} \log \left( \sum_{\tilde{y}} \exp(\beta s(\tilde{y})) \right) - s(y)$$

"log-sum-exp"  $\rightarrow$  "soft-max"  $\hat{y} = \underset{\tilde{y}}{\operatorname{argmax}} s(\tilde{y})$

why?  $\frac{1}{\beta} \log \left[ \exp(\beta s(\hat{y})) \left[ \sum_{\tilde{y}} \exp(\beta (s(\tilde{y}) - s(\hat{y}))) \right] \right]$

$\leq |\mathcal{Y}|$

$$= s(\hat{y}) + \frac{1}{\beta} \log \left( \underbrace{\sum_{\tilde{y}} \exp(\beta (s(\tilde{y}) - s(\hat{y})))}_{\leq |\mathcal{Y}|} \right)$$

$\beta \rightarrow \infty$  (ie, zero-temperature limit)  $\frac{1}{\beta} \log \left( \sum_{\tilde{y}} \exp(\beta s(\tilde{y})) \right) \xrightarrow{\beta \rightarrow \infty} \max_{\tilde{y}} s(\tilde{y})$

$\lim_{\beta \rightarrow \infty} \log \text{loss}(\beta) \leadsto \text{perceptron loss?}$

UN  $\left( \frac{\exp(s_y)}{\sum_{\tilde{y}} \exp(s_{\tilde{y}})} \right)_y$   
 $\rightarrow$  they call it soft-max  
 I call it soft-argmax