

today: energy based methods & surrogate losses
multitasks

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

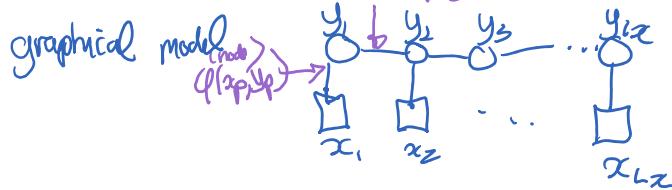
B R A C E

$$\mathcal{S}(x) = \sum_{i=1}^{L_x} x_i$$

$$\Sigma = \{A, \dots, Z\}$$

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

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



$$P_w(y|x) = \frac{1}{Z_w(x)} \exp(\langle w, \ell(x, y) \rangle) = \frac{1}{Z_w(x)} \prod_{c \in C} \pi_{yc}^{N_c(x_c, y_c)}$$

where $C = \{p, p+1\}$

notation $y_C \triangleq (y_i)_{i \in C} \Rightarrow$ can compute argmax $\arg \max_y \langle w, \ell(x, y) \rangle$

dimensions: $16 \cdot 8 \cdot 26$
 \downarrow
of characters

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

node: $\ell^{(\text{node})}(x_p, y_p) = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ \text{vector}(y_p)_{(6 \times 8)} \end{pmatrix} \leftarrow y_p^{\text{in position}} \quad w_a \quad \boxed{a}$

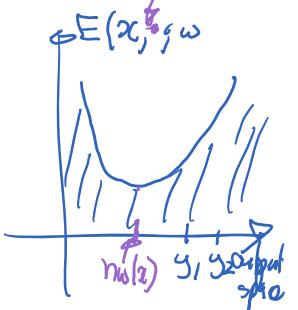
$$\begin{aligned} & \langle w, \ell^{(\text{node})}(x_p, y_p) \rangle \\ &= 0 + 0 + \langle w_{yp}, x_p \rangle + 0 \\ & \quad \text{template for } y_p \end{aligned}$$

edge feature: $\ell(y_p, y_{p+1}) = \begin{pmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \\ \text{if } y_p = y_p, y_{p+1} = y_{p+1} \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{pmatrix} \quad 26^2$

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

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

thing I'm varying



energy based methods : [Secan & al 2006]

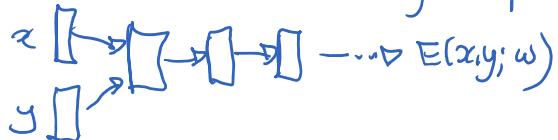
$$\begin{aligned} \text{model: } h_w(x) &= \underset{y \in \mathcal{Y}(x)}{\operatorname{argmin}} E(x, y; w) \quad \text{"energy fn."} \\ &= \underset{y \in \mathcal{Y}(x)}{\operatorname{argmax}} s(x, y; w) \quad \text{"score / compatibility"} \end{aligned}$$

ingredients : 1) what is $E(x, y; w)$?

modelling

$$e.g. s(x, y; w) = \langle w, \ell(x, y) \rangle$$

or $E(x, y; w)$ scalar output of a NN with x, y as input



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

learning

3) how to evaluate $E(x, y; w)$ on a training set? \rightarrow surrogate loss

"quality"

in general: $\tilde{\mathcal{L}}(x^{(i)}, y^{(i)}; E(\cdot, \cdot; w))$

"loss functional"

4) how to minimize $\tilde{\mathcal{L}}(w)$ to learn w ? \rightarrow optimization tricks

flat multiclass case

"flat" (i.e. non structured) setting $h_w(x) = \underset{y}{\operatorname{argmax}} \langle w_y, \ell(x) \rangle$

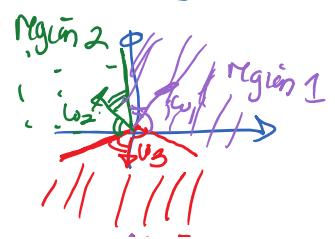
equivalent to $\ell(x, y) = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ \ell(x) \\ 0 \end{pmatrix}_{\in \mathbb{R}^{d \times K}}$

y^{th} position

of classes

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

visually: $|w_y| = 1$



region 3

contrast this flat case
with

structured case

$$\text{e.g. OCR: node feature map } \langle w, \varphi(x, y) \rangle = \sum_p \langle w, \varphi^{(node)}(x_p, y_p) \rangle$$

→ here "sharing" of parameters
between pieces of the joint labels
→ "structure"

aside: in structured prediction, usually absorb "bias" in parameters $\tilde{\varphi}(x)$

$$\langle \tilde{w}, \tilde{\varphi}(x) \rangle = \langle w, \varphi(x) \rangle + b \begin{pmatrix} \varphi(x) \\ 1 \end{pmatrix}$$

$$\tilde{w} = \begin{pmatrix} w \\ b \end{pmatrix}$$

open question: regularizing or not the bias

does it matter in structured pred.?

15n 45'

Surrogate losses

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

I) perceptron loss [Collins et al 2002 EMNLP]

$$J_{\text{percep}}(x, y; w) = \left[\max_{\hat{y} \in \mathcal{Y}(x)} s(x, \hat{y}; w) - s(x, y; w) \right]_+$$

↑ score of ground truth
not needed if assume $y \in \mathcal{Y}(x)$
by using $\hat{y} = y$

$$s(x, y; w) = \langle w, \varphi(x, y) \rangle$$

$$\max_{\hat{y}} \langle w, (\varphi(x, \hat{y}) - \varphi(x, y)) \rangle - \psi(\hat{y}) \geq 0$$

observations: 1) degenerate solution to $\hat{J}(w)$ with $w=0$ or constant $s(x, y)$
2) overaged perceptron alg. :

- Observations:
- 1) degenerate solution to $\mathcal{J}(\mathbf{w})$ with $\mathbf{w} = \mathbf{0}$ w/ constant rate over y
 - 2) averaged perceptron alg.:

- amounts to running constant step size stochastic subgradient method on $\mathcal{J}(\mathbf{w})$

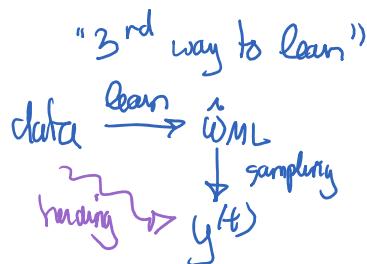
- output $\hat{\mathbf{w}}_T = \frac{1}{T+1} \sum_{t=0}^T \mathbf{w}_t$ (Polyak avg.)

\hookrightarrow will converge to $\mathbf{w} = \mathbf{0}$ when data is not separable

Comments: 1) Cöllin's paper \rightarrow he gives error bound

and generalization error guarantees for perceptron

2) (aside) connection with the "hawking" alg. by Welling & al. [ICML 2012]



II) log-loss (CRF) (probabilistic interpretation)

suppose $p(y|x; \mathbf{w}) \propto \exp(\beta s(x, y; \mathbf{w}))$

β
"inverse temperature" parameter

Boltzmann dist

$$\beta = \frac{1}{k_B \cdot T}$$

temperature

MCL \rightarrow log-loss

$$\begin{aligned}
 \mathcal{L}(x, y; \mathbf{w}) &= -\frac{1}{\beta} \log p(y|x; \mathbf{w}) = -\frac{1}{\beta} \log \left(\frac{\exp(\beta s(x, y))}{\sum_y \exp(\beta s(x, y))} \right) \\
 &\quad \text{rescaling} \\
 &= -\frac{1}{\beta} \log \left(\sum_y \exp(\beta s(x, y)) \right) - s(x, \hat{y}) \quad \text{partition fct.} \\
 &\quad \text{"log-sum-exp" } \rightarrow \text{soft max" why?} \\
 &\quad \text{let } \hat{y} = \arg \max_y s(x, y)
 \end{aligned}$$

$$\text{let } \hat{y} = \underset{y}{\operatorname{argmax}} s(y)$$

$$\frac{1}{\beta} \log \left(\exp(\beta s(\hat{y})) \left[\sum_y \exp(\beta(s(\hat{y}) - s(\hat{y})) \right] \right)$$

$$= \frac{1}{\beta} s(\hat{y}) + \frac{1}{\beta} \log \left(\sum_y \exp(\beta(s(\hat{y}) - s(\hat{y})) \right)$$

$\beta \rightarrow \infty$ (i.e. zero temperature limit)

$$\frac{1}{\beta} \log \left(\sum_y \exp(\beta \cdot s(y)) \right) \xrightarrow{\beta \rightarrow \infty} \underset{\hat{y}}{\operatorname{max}} s(\hat{y})$$

note:
In deep learning
they call "softmax"

$$\left(\frac{\exp(s_y)}{\sum_y \exp(s_y)} \right)_{y \in \mathcal{Y}}$$

I call this
"soft argmax"

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