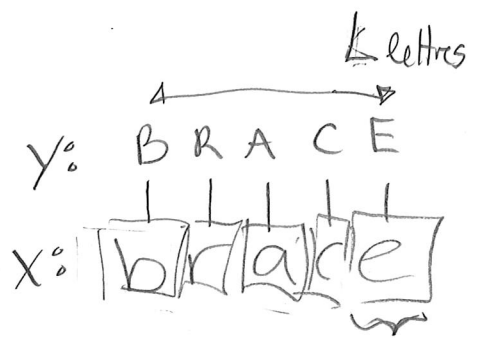


prédiction structurée

motivation: OCR



$f: X \rightarrow Y$  objets structurés  $K$  lettre possible

$|Y| = K^L \rightarrow$  # exponentiel de possibilités

approche multiclasse: one-vs-rest } impossible?  
régression logistique multiclasse

approche indépendante  $\rightarrow$  ignore les corrélations entre décision

$Y = (y_1, \dots, y_L)$   
 $p(Y|X)$  à corrélation?

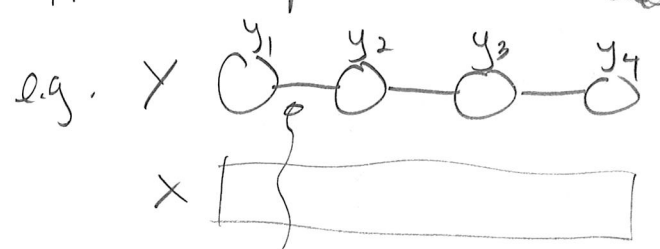
prédiction structurée  $\rightarrow$  décision jointe

problème:  $p(Y|X)$   
pour  $X$  fixé;  $K^L - 1$  paramètres en général (!)  
[# augacho très rapidement...]

exploiter structure pour < statistique (trop de paramètres)  
& relié  
computational (max  $p(Y|X)$  NP-dur en général)  
yes

outil (cours MVA): modèle graphique

suppose des indépendances conditionnelles entre les variables:



chaîne

factorise avec  $\prod_{t=1}^{L-1} \psi_t(y_t, y_{t+1}, X)$  (tableau  $y_t, y_{t+1}$ )

#paramètres  $K^2$  potentiel

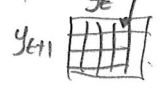
donc passe de  $\approx K^L$  à  $L K^2$  paramètres

clique  $(y_t, y_{t+1})$

graphe encafé les indépendances conditionnelles

donc passe de

$\approx K^L$  à  $L K^2$  paramètres



dans cas MVA, montre graphe  $\Leftrightarrow$  factorisation  $\Leftrightarrow$  indépendances conditionnelles

ici  $y_{t+2} \perp\!\!\!\perp y_t \mid y_{t+1}$

ie.  $p(y_{t+2}, y_t \mid y_{t+1}, x) = p(y_{t+2} \mid y_{t+1}, x) \cdot p(y_t \mid y_{t+1}, x)$

pourquoi factorisation aide?

$$\text{score}(x, y) \triangleq \log p(y \mid x) = \sum_t \underbrace{\log \psi_t(y_t, y_{t+1}, x)}_{\text{score}(y_t, y_{t+1}, x)} - \log Z(x)$$

$$\max_{y \in \mathcal{Y}} [\text{score}(x, y)] = \max_{y_1} \left[ \max_{y_2} \left[ \dots + \max_{y_{k-1}} \left[ \max_{y_k} \left[ \text{score}(y_{k-1}, y_k) \right] \right] \right] \right]$$

$m(y_{k-1})$

$m(y_{k-2})$

$m(y_2)$  } "messages" dans graphe

temps  $\rightarrow O(L K^2)$  [vs.  $O(K^L)$ !]

programmation dynamique pour trouver argmax

algorithme de Viterbi

par similitude:

[anneau  $(\max, +)$  vs.  $(+, \times)$ ]

$$Z(x) = \sum_{y_1} \left[ \sum_{y_2} \left[ \sum_{y_3} \left[ \dots \sum_{y_L} \psi_L(y_{L-1}, y_L) \right] \right] \right]$$

$m(y_{L-1})$

"sum-product algorithm"

$$ax + ay + bx + by = (a+b)(x+y)$$

donc manière "tractable" de faire regression logistique structurée (2)

CRF = Conditional Random Field

$p(y|x)$  is random field

special case  
of graphical model

$$p(y|x) \propto \exp(\text{score}(x, y))$$

$$\text{score}(x, y) = y^T w^T x$$

pour  $y \in \{0, 1\}$

done regression logistique

$\dots w^T x \rightarrow$  multiclass

$$Z(x) = \sum_y \exp(\text{score}(x, y))$$

exemple de  
paramétrisation:

$$\text{score}(w, x, y) \triangleq \sum_t (w_{\text{node}}^T \psi(x_t, y_t))$$

$$+ \sum_t (w_{\text{edge}}^T \psi(y_t, y_{t+1}))$$

$$s(x, y_{t+1}, x)$$

$$= w^T \psi(x, y)$$

pour OCR :

$$\psi(y_t, y_{t+1}) = \mathbb{1}\{y_t = 'a', y_{t+1} = 'c'\}$$

feature "indicatrice"

dimension  $K^2$

done possibilité d'apprendre  
les probabilités  $p(y_t, y_{t+1})$

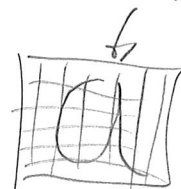
marginales

[corrélation adjacente  
de lettres]

$$\psi(x_t, y_t) = \begin{pmatrix} \vdots \\ \mathbb{1}\{y_t = 'a'\} x_t \\ \vdots \end{pmatrix}$$

$Kd$  paramètres

$w_{\text{node}}^{'a'}$



régression logistique régularisée :

$$\min_w \lambda \|w\|^2 + \frac{1}{n} \sum_i [-\log p(y_i|x_i)] = J_\lambda(w)$$

$$\frac{\exp(w^T \psi(x_i, y_i))}{Z(x)}$$

$$\lambda \|w\|^2 + \frac{1}{n} \sum_i [\log(Z(x)) - w^T \psi(x_i, y_i)]$$

$$\nabla J_\lambda(w) = \lambda \vec{w} + \frac{1}{n} \sum_i \nabla_w (\log(Z(w)) - \psi(x_i, y_i))$$

$$\log\left[\sum_y \exp(w^T \psi(x_i, y))\right]$$

$$\frac{\sum_y \exp(w^T \psi(x_i, y)) \psi(x_i, y)}{\sum_y \exp(w^T \psi(x_i, y))}$$

$$p(y|x)$$

$$= \mathbb{E}_{p(y|x)} \psi(x_i, y) - \psi(x_i, y_i)$$

peut calculer avec sum-product

[sans régularisation,  $\nabla J_\lambda = 0 \Rightarrow \frac{1}{n} \sum_i [\mathbb{E}_{p(y|x_i)} \psi(x_i, y) - \psi(x_i, y_i)] = 0$  "feature moment matching"]

voir slides Ben Taskar

autres exemples :

- reconnaissance vocale
- traduction automatique
- coding région detection in DNA
- etc?
- alignement de mots

autres sujets

$p(y)$  apprendre  
continue

- "unsupervised learning"
- dimensionality reduction etc...

CS + STAT