

# Lecture 1 - scribbles

Friday, January 6, 2017  
14:18

gen. / discriminative continuum

prediction  $\rightarrow$  learn a mapping  $h: X \rightarrow \mathcal{Y}$

$p_w(x, y)$    
  $\nwarrow$  output   
  $\swarrow$  input   
 [generative model on  $X \times \mathcal{Y}$ ]

$p_w(y|x) \rightarrow h_w(x) = \arg \max_{y \in \mathcal{Y}} p_w(y|x)$   
"more discriminative"

$p_w(x, y)$   
 $\downarrow$   
learn  $\hat{w}$   
by ML

$p_w(y|x)$   
 $\downarrow$   
learn  $\hat{w}$   
by MCL

$h_w: X \rightarrow \mathcal{Y}$   
 $\ell(y, y')$   
learn  $\hat{h}_w$  using  
surrogate loss minimization }  
structured SVM

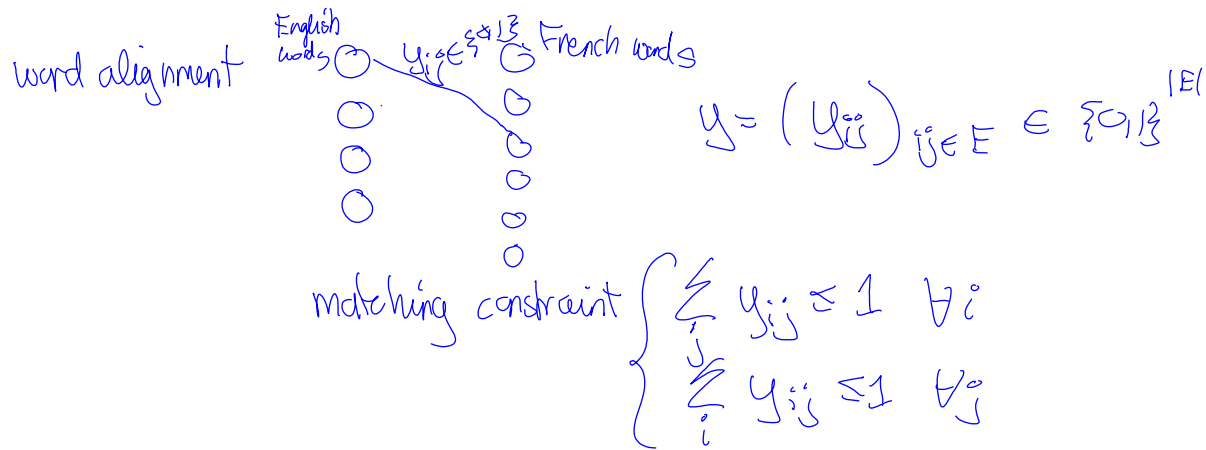
related  
empirical risk  
minimization

$\leftarrow$  more assumptions / loss robust

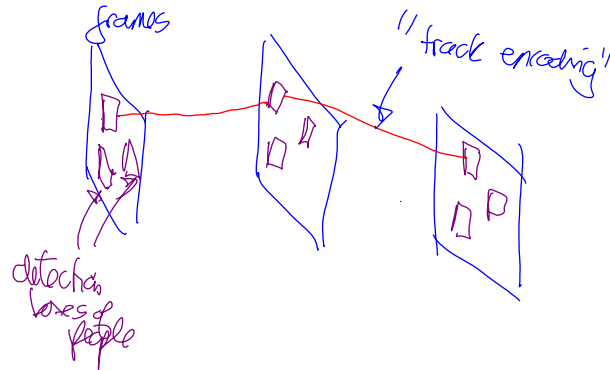
some important aspects of structured prediction:

- 1)  $\mathcal{Y}$  output space is usually exponentially big
- 2) error function  $\ell(y, y')$
- 3) constraints on pieces of  $y$

need an "encoding function" ; e.g.  $y = (y_1, \dots, y_p) \in \mathbb{R}^p$



another example:  
multi-object tracking



simple structured prediction models

1)  $h_w(x) \triangleq \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \underset{\substack{\uparrow \text{'score'}}}{sc(x, y; w)}} \quad ] \text{'compatibility function'}$

$-E(x, y; w) \quad ] \text{energy function}$

example:  $sc(x, y; w) = \langle w, \underbrace{\phi(x, y)}_{\substack{\text{joint feature vector} \\ \in \mathbb{R}^d}} \rangle$

word alignment example: recall  $y = (y_{ij}) \in \{0,1\}^{|E|}$

word alignment example: recall  $y = (y_{ij}) \in \{0,1\}^{|E|}$

$$\phi(x, y) = \sum_{i,j} y_{ij} \underbrace{\psi(x_{ij})}_{\text{features defined on English word } i \text{ and French } j}$$

(string edit distance  $(i, j)$   
 $\mathbb{1}_{\{(i,j) \text{ are in dictionary}\}}$   
 distance index  $i \neq j$ )

here, you  
 can do:  
 $\arg \max_{y \in \text{matchings}} \phi(x, y)$   
 using min-cost network  
 algorithms

$$\langle w, \phi(x, y) \rangle = \sum_{i,j} y_{ij} \underbrace{\langle w, \psi(x_{ij}) \rangle}_{\text{represents how much we want to align } i \text{ \& } j}$$

another example from a undirected graphical model

$$p_w(y|x) = \frac{1}{Z_w(x)} \prod_c \psi_c(y_c, x)$$

$$= \exp\left(\sum_c \underbrace{\log \psi_c(y_c, x)}_{\langle w, \psi_c(y_c, x) \rangle} - A(w|x)\right)$$

sequence model  $\stackrel{\text{bg}}{p_w(y|x)} = \sum_t \langle w, \phi(y_t, y_{t+1}; x) \rangle + \text{const.}$



$\arg \max$  here  $\rightarrow$  Viterbi algorithm

II) how do we learn  $w$ , given some training data  $(x_i, y_i)_{i=1}^n$

simplest algorithm:

Structured perceptron:

go over training set:

• sample  $i_t$  [predict on  $x_{i_t}$ ]

• let  $\hat{y}_t = h_{w_t}(x_{i_t}) = \arg \max_{y \in \mathcal{Y}} \langle w, \phi(x_{i_t}, y) \rangle$

$$w_{t+1} = w_t + \underbrace{\eta}_{\text{step-size}} \underbrace{(\phi(x_{i_t}, y_{i_t}) - \phi(x_{i_t}, \hat{y}_t))}_{\text{boost ground truth score} - \text{pending prediction}}$$

for stability:

$$\text{output } \hat{w} = \frac{1}{T} \sum_{t=1}^T w_t \quad \leftarrow \text{"Polyak averaging"}$$

this can be interpreted as doing stochastic subgradient optimization on the following non-smooth objective:

$$L(w) = \frac{1}{n} \sum_{i=1}^n \left[ \max_{y \in \mathcal{Y}} \langle w, \phi(x_i, y) \rangle - \langle w, \phi(x_i, y_i) \rangle \right]_+$$

if  $y_i \in \mathcal{Y}$ , then always positive and  $\therefore [ \cdot ]_+$  is not needed

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

Conditional random fields (CRF)

define  $p_w(y|x) \propto \exp(\langle w, \phi(x,y) \rangle)$

then do MCL on the training set

caveat: need (implicitly) to be able to sum over  $\mathcal{Y}$  i.e.  $\sum_{y \in \mathcal{Y}} \exp(\langle w, \phi(x,y) \rangle)$

#P-complete for set of matchings!

Structured SVM:

intuition: want  $\langle w, \phi(x_i, y_i) \rangle \geq \langle w, \phi(x_i, y) \rangle + \ell(y_i, y) \quad \forall y \in \mathcal{Y}_i$

min  $\|w\|^2$

s.t.

(hard margin structured SVM)

$$\textcircled{I} \min_{w, \xi} \left[ \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \xi_i \right] \quad \left[ \begin{array}{l} \text{QP} \\ \text{with exponential} \\ \text{\# of constraints} \end{array} \right]$$

$$\xi_i + \langle w, \phi(x_i, y_i) \rangle \geq \langle w, \phi(x_i, y) \rangle + \ell(y_i, y) \quad \forall y \in \mathcal{Y}_i, \forall i$$

equivalently

$$\xi_i + \langle w, \phi(x_i, y_i) \rangle \geq \max_{y \in \mathcal{Y}_i} \{ \langle w, \phi(x_i, y) \rangle + \ell(y_i, y) \} \quad \text{non-linear constraint}$$

equivalent problem

$\textcircled{II}$

$$\min_w \left[ \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \left[ \max_{y \in \mathcal{Y}_i} \{ \langle w, \phi(x_i, y) \rangle + \ell(y_i, y) \} - \langle w, \phi(x_i, y_i) \rangle \right] \right]$$

structured hinge loss