

[bit.ly/IFT6132-W18](https://bit.ly/IFT6132-W18)

## Structured prediction basics & setup

Learning problem:

given a training dataset  $D = (\underset{\substack{\downarrow \\ \text{ex}}}{x^{(i)}}, \underset{\substack{\downarrow \\ \text{es}}}{y^{(i)}})_{i=1}^n$

goal: learn a prediction mapping  $h_{\omega}: X \rightarrow \mathcal{Y}$   
 $\omega$  parameter

that has low generalization error  $L(\omega; P) \triangleq \mathbb{E}_{(x,y) \sim P} [l(y, h_{\omega}(x))]$   
 "test" distribution on  $(x,y)$   
 structured error function

"risk" in ML (Vapnik) / statistical decision loss

[see lecture 4 & 5 of my P6M  $\rightarrow$  review of statistical decision theory]

regularized ERM  
 empirical risk minimization

$$\hat{L}(\omega) = \frac{1}{n} \sum_{i=1}^n l(y^{(i)}, h_{\omega}(x^{(i)})) + R(\omega)$$

↑  
regularizer

not obj. in  $\omega$ ; non-convex  
 messy ... NP hard to minimize in general

$\Rightarrow$  replace it with surrogate loss / contrast fct.  
 ML statistics

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

M-estimator in statistics

surrogate loss e.g. convex in  $w$

↳ examples:   
 • structured hinge loss  $\rightarrow$  structured SVM   
 • log-loss  $\rightarrow$  CRF

## generative vs discriminative learning continuum

generative modeling of  $(x, y)$

$$p_w(x, y)$$

↳ condition  $p_w(y|x) \rightarrow h_w(x) = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} p_w(y|x)$

"more discriminative"

$$p_w(x, y)$$

learn  $\hat{w}$   
by ML

$$p_w(y|x)$$

learn  $\hat{w}$   
by MCL

$$h_w: X \rightarrow \mathcal{Y}$$

learn  $\hat{h}_w$  by using  
surrogate loss minimization  
e.g. structured SVM

related to  
ERM

↳

more assumptions / less robust on classification task

## some important aspects of structured prediction (vs binary classification)

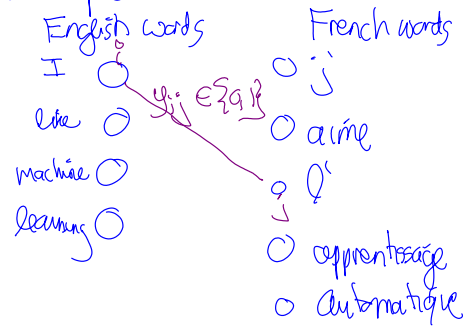
1)  $\mathcal{Y}$  output space is usually exponentially big

2) <sup>structured</sup> error function  $\ell(y, y')$

3) ...

3) sometimes constraints on pieces of  $y$   
 need an "encoding function"  $y = (y_1, \dots, y_p) \in \mathbb{R}^p$

word alignment example:



$$y = (y_{ij})_{(i,j) \in E} \in \{0,1\}^{|E|}$$

↑ possible edges (all pairs)

$$\text{matching constraint} \begin{cases} \sum_j y_{ij} \leq 1 & \forall i \\ \sum_i y_{ij} \leq 1 & \forall j \end{cases}$$

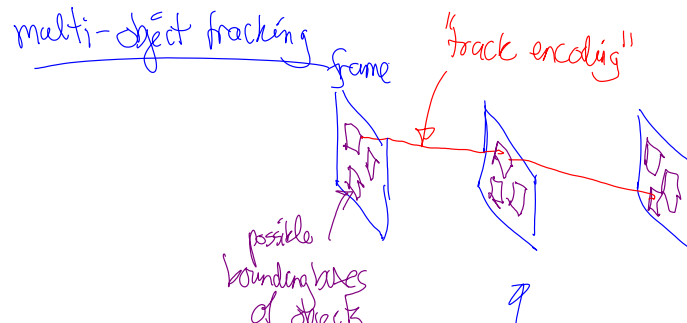
here  $x = (x_1^E, \dots, x_{|E|}^E, x_1^F, \dots, x_{|F|}^F)$

English words

French words

$$\text{here } \gamma(x) = \{y \in \{0,1\}^{|E| \cup |F|} : \sum_j y_{ij} \leq 1, \sum_i y_{ij} \leq 1 \forall j\}$$

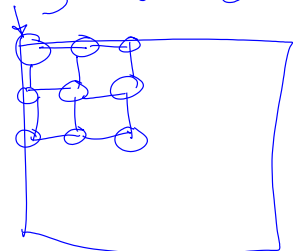
another example:



properties  
of objects

model this with network flow structure

given Image Segmentation:



$x$  = image of RGB values  $L \times L$  pixels

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

background      foreground

standard hwk:

$$h_w(x) \triangleq \underset{y \in \mathcal{Y}(x)}{\operatorname{argmax}} \begin{matrix} \overset{\text{"score"}}{S(x, y; w)} \\ -E(x, y; w) \end{matrix} \quad \begin{matrix} \text{compatibility fct. of } y \text{ with } x \\ \text{energy fct.} \end{matrix}$$

linear model:  $S(x, y; w) = \langle w, \underbrace{\phi(x, y)}_{\text{"joint feature" vector} \in \mathbb{R}^d} \rangle$        $\phi: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$

word alignment example:  $\phi(x, y) = \sum_i u_i \phi(x_i^E, x_i^F)$

$$\phi(x_i^E, x_j^F) \in \mathbb{R}^d$$

features defined on  
English word  $x_i^E$   
and French "  $x_j^F$

- string edit distance  $(x_i^E, x_j^F)$
- $\{ (x_i^E, x_j^F) \text{ are in dictionary} \}$
- distance between  $i$  &  $j$

$$S(x, y; w) = \langle w, \ell(x, y) \rangle = \sum_{i,j} y_{ij} \underbrace{\langle w, \ell(x_i^E, x_j^F) \rangle}_{\text{"score to match } i \text{ to } j"}$$

$$h_w(x) = \arg \max_{y \in \mathcal{Y}(x)} S(x, y; w)$$

$$\begin{aligned} \max_y \quad & \sum_{i,j} y_{ij} s_{ij}(x) \\ \text{s.t.} \quad & y_{ij} \in \{0, 1\} \\ & \sum_j y_{ij} \leq 1 \\ & \sum_i y_{ij} \leq 1 \end{aligned}$$

can be solved  
exactly  
as min cost matching  
problem

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

## Learning $w$ ?

### structured perceptron:

- initialize  $w_0$
- repeat for  $t=0, \dots$ 
  - sample  $x_t$
  - let  $\hat{y}_t = h_{w_t}(x^{(t)}) = \arg \max_{y \in \mathcal{Y}(x^{(t)})} \langle w_t, \ell(x^{(t)}, y) \rangle$
  - $w_{t+1} = w_t + \eta \left( \underbrace{\ell(x^{(t)}, y^{(t)})}_{\substack{\text{step size} \\ \Rightarrow \text{boost ground truth score}}} - \underbrace{\ell(x^{(t)}, \hat{y}_t)}_{\substack{\text{prune prediction score}}} \right)$

decoding oracle

for stability:

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

⊛ structured perceptron can be interpreted as

doing stochastic subgradient optimization on the following non-smooth objective

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

$$f^{\text{percept}}(x, y, w) \triangleq \left[ \max_{\tilde{y} \in \mathcal{Y}} \langle w, \phi(x, \tilde{y}) \rangle - \langle w, \phi(x, y) \rangle \right]_+$$

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  $\geq 0$   
and so  $[\cdot]_+$  is not needed

conditional random field:

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

$$h_w(x) = \arg\max_{y \in \mathcal{Y}} p_w(y|x) = \arg\max_y \langle w, \phi(x, y) \rangle$$

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

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

$$\begin{aligned} f^{\text{CRF}}(x, y; w) &= -\log p_w(y|x) \\ &= \log \left( \sum_{\tilde{y}} \exp(\langle w, \phi(x, \tilde{y}) \rangle) \right) - \langle w, \phi(x, y) \rangle \end{aligned}$$

$$\frac{\underbrace{\quad}_{\text{yes}}}{\log Z_w(x)}$$

Issues: •  $\ell(y, y')$  doesn't appear in it

•  $\sum_{\tilde{y} \in \mathcal{Y}} \exp(w, \phi(x, \tilde{y}))$  can be difficult

#P-complete for  $\mathcal{Y}$  = set of matchings

Structured SVM:

intuition: want  $\langle w, \phi(x^{(i)}, y^{(i)}) \rangle \gg \langle w, \phi(x^{(i)}, \tilde{y}) \rangle + \ell(y^{(i)}, \tilde{y}) \quad \forall \tilde{y} \in \mathcal{Y}_i \triangleq \mathcal{Y}(x^{(i)})$

min  $\|w\|^2$  s.t.  $\uparrow$  "hard margin structured svm"

(binary svm:  $y \in \{-1, +1\}$   $h_w(x) = \text{sgn}(\langle w, \phi(x) \rangle)$   
 $y \langle w, \phi(x) \rangle \geq 1$  hard margin constraints)

soft-margin SVM:

$$P(w) = \frac{1}{n} \sum_{i=1}^n \mathcal{J}^{\text{svm}}(x^{(i)}, y^{(i)}; w)$$

GP with exponential # of constraints

$$\left[ \begin{array}{l} \min_{w, \xi} \\ \xi_i + \langle w, \phi(x^{(i)}, y^{(i)}) \rangle \geq \langle w, \phi(x^{(i)}, \tilde{y}) \rangle + \ell(y^{(i)}, \tilde{y}) \quad \forall \tilde{y} \in \mathcal{Y}_i, \forall i \end{array} \right.$$

$$\frac{\lambda \|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n \xi_i$$

equivalent formulation is  $\min_w \frac{\lambda \|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n \mathcal{J}^{\text{svm}}(x^{(i)}, y^{(i)}; w)$

where  $\mathcal{J}^{\text{svm}}(x, y; w) \triangleq \max_{\tilde{y} \in \mathcal{Y}(x)} \underbrace{[\langle w, \phi(x, \tilde{y}) \rangle + \ell(y, \tilde{y})]}_{\text{"loss-augmented decoding"}} - \langle w, \phi(x, y) \rangle$   
"structured hinge loss"

$\Rightarrow$  hinge not that  $\perp$   $\cdot \ell(i) = \max_{\tilde{y} \in \mathcal{Y}(x^{(i)})} [\langle w, \phi(x^{(i)}, \tilde{y}) \rangle + \ell(y^{(i)}, \tilde{y})] - \langle w, \phi(x^{(i)}, y^{(i)}) \rangle$

it turns out that if  $y^{(i)} \in \mathcal{S}(x^{(i)})$ ; then  $f^{\text{SVM}}(x^{(i)}, y^{(i)}; w) \geq \lambda(y^{(i)}, h w | x^{(i)})$