

Lecture 3 - surrogate losses

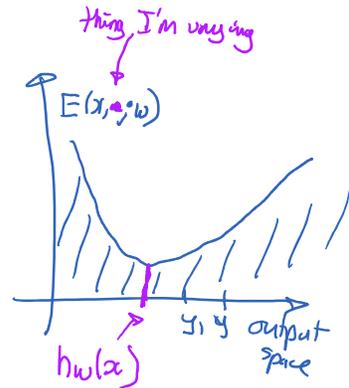
Thursday, January 21, 2021 13:39

today: energy based methods & surrogate losses  
multiclass

energy based methods : [Lecan & al. 2006]

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

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



ingredients:

modeling

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

e.g.  $S(x, y; w) = \langle w, \phi(x, y) \rangle$

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



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

learning

3) how to evaluate "quality"  $E(x, y; w)$  on a training set?  $\rightarrow$  surrogate loss  $\hat{J}(w)$

in general:  $\hat{J}(x^{(i)}, y^{(i)}; E(\cdot, \cdot; w))$  "loss functional"

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

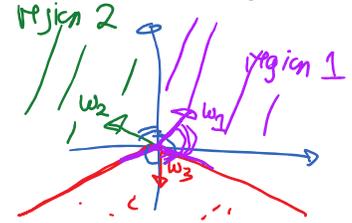
flat multiclass case

"flat" (ie. non structural) setting  $h_w(y) = \underset{y}{\text{argmax}} \langle w_y, \phi(x) \rangle$   $\in \mathbb{R}^d$

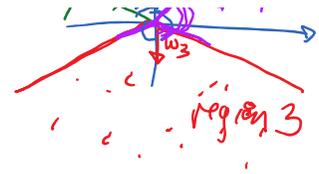
equivalent to  $\phi(x, y) = \begin{pmatrix} 0 \\ \vdots \\ \phi(x) \\ \vdots \\ 0 \end{pmatrix} \in \mathbb{R}^{d \times k}$   $\leftarrow$   $y^{\text{th}}$  position # of classes

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

visually:  $\|w_y\| = 1$



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



contrast this flat case  
with  
structured case:

eg. OCR node feature map  $\langle w, \phi(x, y) \rangle = \sum_p \langle w, \phi^{(node)}(x_p, y_p) \rangle$   
 $\sum_{y_p} \mathbb{I}\{y_p, y_p\} \langle w_{y_p}, \phi(x) \rangle$

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

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

standard binary classification  $\text{sgn}(\langle w, x \rangle + b)$   $\left\{ \begin{matrix} \tilde{\phi}(x) \\ 1 \end{matrix} \right.$

$$\langle \tilde{w}, \tilde{\phi}(x) \rangle = \langle w, \phi(x) \rangle + b$$

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

open question: regularizing or not the bias  
does it matter in struct. pred.?

14h27

surrogate losses

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

I). perception loss [Collins & al. 2002 EMNLP]

$$\ell^{prop.}(x, y; w) = \left[ \max_{\tilde{y} \in \mathcal{Y}(x)} s(x, \tilde{y}, w) - s(x, y; w) \right]_{+}$$

score of ground truth

not needed  
if assume  $y \in \mathcal{Y}(x)$

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

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

by using  $\tilde{y} = y$

observations: 1) degenerate solution to  $\hat{J}(w)$  with  $w=0$  or constant score over  $y$

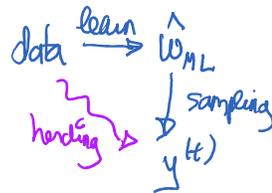
2) averaged perception obj:

- amounts to running constant step-size stochastic subgradient method on  $\hat{J}(w)$
  - output  $\hat{w}_T = \frac{1}{T+1} \sum_{t=0}^T w_t$  (Polyak avg.)
- does not converge in general
-

↳ will converge to  $w^* = 0$   
when data is not separable

Comments 1) Collins's paper → he gives error bound and generalization error guarantees for perceptron

2) (aside) connection with the "harding" alg. by Welling & d. [ICML 2012]  
 "3rd way to learn"



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

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

Boltzmann dist.  
 $\beta = \frac{1}{k_B T_{temp}}$

MCL → log-loss

$$J(x,y;w) = -\log p(y|x;w) = -\frac{1}{\beta} \log \left( \frac{\exp(\beta s(x,y;w))}{\sum_{\tilde{y}} \exp(\beta s(x,\tilde{y};w))} \right)$$

partition function

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

"log-sum-exp" → "soft max." why?  
 let  $\hat{y} = \arg \max_{\tilde{y}} s(x,\tilde{y};w)$

$$\frac{1}{\beta} \log \left( \exp(\beta s(x,\hat{y};w)) \left[ \sum_{\tilde{y}} \exp(\beta (s(x,\tilde{y};w) - s(x,\hat{y};w))) \right] \right)$$

$$= \frac{1}{\beta} s(x,\hat{y};w) + \frac{1}{\beta} \log \left( \sum_{\tilde{y}} \exp(\beta (s(x,\tilde{y};w) - s(x,\hat{y};w))) \right) \leq |S|$$

as  $\beta \rightarrow \infty$  (i.e. zero temp. limit)  $\frac{1}{\beta} \log \left( \sum_{\tilde{y}} \exp(\beta s(x,\tilde{y};w)) \right) \rightarrow \max_{\tilde{y}} s(x,\tilde{y};w)$

Note: in deep learning book, they call "soft-max"

$$\left( \frac{\exp(s(y))}{\sum_{\tilde{y}} \exp(s(\tilde{y}))} \right)_{y \in Y}$$

I call this "soft-argmax" thus  $\lim_{\rho \rightarrow \infty} \log\text{-loss}(\rho) \rightarrow \text{perceptron loss}$

### III) structured hinge loss

$$f^{\text{svm}}(x, y; w) = \max_{\tilde{y} \in \mathcal{Y}(x)} [s(x, \tilde{y}; w) + l(y, \tilde{y})] - s(x, y; w)$$

↳ loss-augmented decoding

a)  $\frac{\text{carban.}}{\Rightarrow l(y, \tilde{y}_{\text{next}})} \left\{ \begin{array}{l} s(y) \\ s(\tilde{y}_{\text{next}}) \\ \vdots \end{array} \right. \Rightarrow f^{\text{svm}}(x, y; w) = 0$

b)  $f^{\text{svm}}(x, y; w) \geq l(y, h_w(x))$

why?  $f(x, y; w) = \max_{\tilde{y}} [s(\tilde{y}) + l(\tilde{y})] - s(y)$   
 $\Rightarrow s(\hat{y}) + l(\hat{y}) - s(y)$  let  $\hat{y} = \underset{\tilde{y}}{\text{argmax}} s(\tilde{y}) = h_w(x)$   
 if  $y \in \mathcal{Y}(x) \Rightarrow s(\hat{y}) \geq s(y)$   
 $\Rightarrow l(\hat{y}) = l(y, h_w(x))$  //

binary case: for structured hinge loss  $\rightarrow$  see notes from last year

$$y \in \{-1, +1\} \quad w = \begin{pmatrix} w_+ \\ w_- \end{pmatrix} \quad \varphi(x, +1) = \begin{pmatrix} \varphi(x) \\ 0 \end{pmatrix}$$

$$h_w(x) = \underset{\text{argmax}}{\{ \langle w_+, \varphi(x) \rangle, \langle w_-, \varphi(x) \rangle \}}$$

predict +1 if  $\langle w_+, \varphi(x) \rangle \geq \langle w_-, \varphi(x) \rangle$

$$\Leftrightarrow \langle w_+ - w_-, \varphi(x) \rangle \geq 0$$

(see notes to show)

$$f^{\text{svm}}(x, y; w) = [1 - y \langle \tilde{w}, x \rangle]_+$$

$$\tilde{w} \triangleq w_+ - w_-$$

$$h_w(x) = \text{sign}(\langle \tilde{w}, x \rangle)$$

(recopied from last year:)

structured hinge loss

$$\tilde{w} = w_+ - w_- \quad h_w(x) = \text{sgn}(\langle \tilde{w}, x \rangle)$$

$$f_{\text{SUM}}(x, y; w) = \max \left\{ \langle w_+, x \rangle + \underbrace{\ell(y, +)}_{\mathbb{1}\{y \neq +1\}}, \langle w_-, x \rangle + \underbrace{\ell(y, -)}_{1 - \mathbb{1}\{y \neq +1\}} \right\} - \langle w_y, x \rangle$$

$w_+ = \tilde{w} + w_-$

$$= \max \left\{ \langle \tilde{w}, x \rangle + \langle w_-, x \rangle + \mathbb{1}\{y \neq +1\}, \langle w_-, x \rangle + 1 - \mathbb{1}\{y \neq +1\} \right\} - \langle w_y, x \rangle$$

$$= \max \left\{ \langle \tilde{w}, x \rangle + 1, 1 - 1 \right\} + \langle w_-, x \rangle - \langle w_y, x \rangle$$

(case  $y = +1$ ):  $\max \{ \langle \tilde{w}, x \rangle, 1 \} - \langle \tilde{w}, x \rangle = [1 - y \langle \tilde{w}, x \rangle]_+$

(case  $y = -1$ ):  $\max \{ \langle \tilde{w}, x \rangle + 1, 0 \} + 0 = [1 - y \langle \tilde{w}, x \rangle]_+$

overall:  $[1 - y \langle \tilde{w}, x \rangle]_+$

∴ ∴ ∴

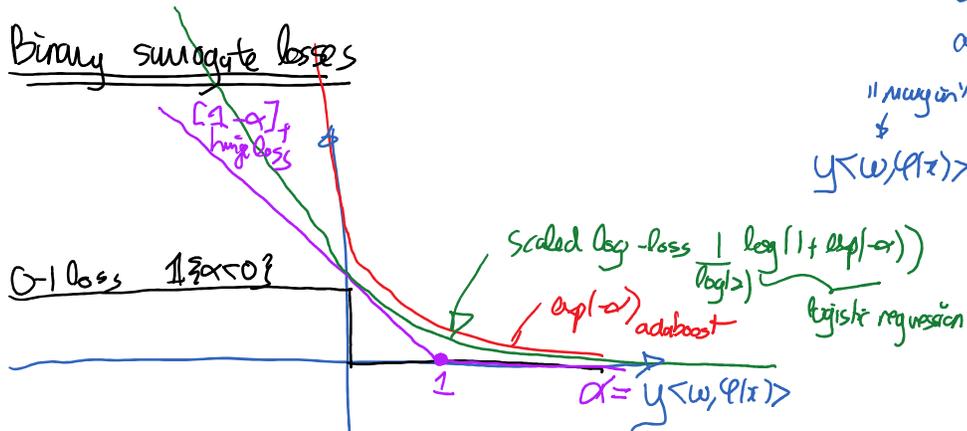
$$f_{\text{SUM}}(x, y; w) = [1 - y \langle \tilde{w}, x \rangle]_+$$

where  $\tilde{w} = w_+ - w_-$

i.e. structured hinge loss reduces to binary SVM hinge loss when using  $\ell(y, y') = \mathbb{1}\{y \neq y'\}$  and  $y = \{-1, +1\}$

"margin"  $\neq$   
 $y \langle w, \phi(x) \rangle \geq 0 \Rightarrow$  make no mistake

### Binary surrogate losses



[Bartlett & al. 2006]  $\rightarrow$  showed all these methods are consistent