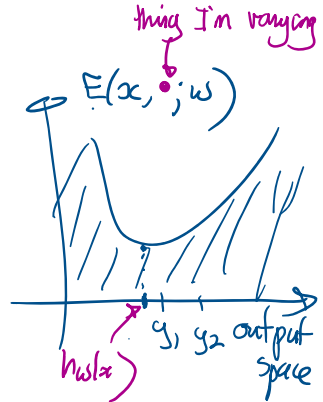


today: • energy based methods & surrogate losses  
• multi class

energy based methods: [LeCun & al. 2006]

model:  $h_w(z) = \underset{\tilde{y} \in \mathcal{Y}(z)}{\operatorname{argmin}} E(x, \tilde{y}; w)$  "energy f.f."  
 $= \underset{\tilde{y} \in \mathcal{Y}(z)}{\operatorname{argmax}} s(x, \tilde{y}; w)$  "score / compatibility"



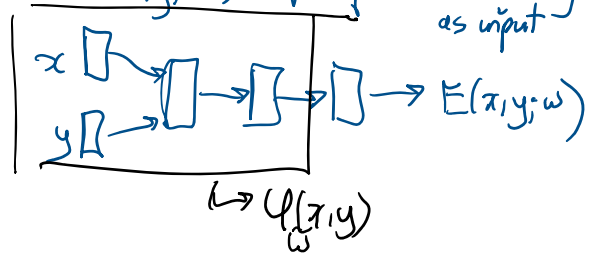
ingredients:

mapping

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

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

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



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

→ "decoding" / "inference"

learning

3) how to evaluate "quality"  $E(x, y; w)$  on a training set?

→ surrogate loss

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

"loss functional"

4) how to minimize  $\tilde{\mathcal{J}}(w)$  to learn  $\hat{w}$ ?

→ optimization tricks

flat multiclass case

"flat" (i.e. non-structured) setting  $h_w(y) = \underset{y}{\operatorname{argmax}} \langle w_y, \phi(x) \rangle$   $w \in \mathbb{R}^d$

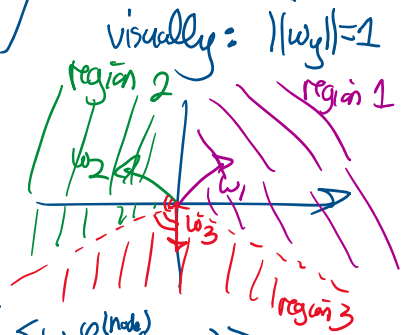
equivalent to  $\phi(x, y) = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \end{pmatrix} \in \mathbb{R}^{d \times k}$

(to simplify)

equivalent to  $\phi(x,y) = \begin{pmatrix} 0 \\ 0 \\ \phi(x) \\ 0 \\ 0 \end{pmatrix} \in \mathbb{R}^{d+k}$   $\xrightarrow{y^{\text{th}} \text{ position}} \begin{pmatrix} w_y \end{pmatrix}$  (to simplify picture)

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

contrast this flat case with structured case:



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

$$\sum_{y_p} \mathbb{1}\{y_p, y\} \langle w_{y_p}, \phi(x) \rangle$$

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

aside: in structured prediction, usually absorbs "bias" in parameters

standard binary classification  $\text{sgn}(\langle w, x \rangle + b)$

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

$$\tilde{w} = \begin{pmatrix} w \\ b \end{pmatrix}, \quad \tilde{\phi}(x) = \begin{pmatrix} \phi(x) \\ 1 \end{pmatrix}$$

open question: regularizing or not the bias  
does it matter in structured prediction?

Summed Losses

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

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

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

score of ground truth

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

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

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

$\triangleq -\psi(\tilde{y})$  by using  $\tilde{y}=y$

$$\hat{y} \triangleq -\nabla_{\tilde{y}} \psi(\tilde{y}) \quad \text{by using } \tilde{y} = y$$

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

2) averaged perceptron alg

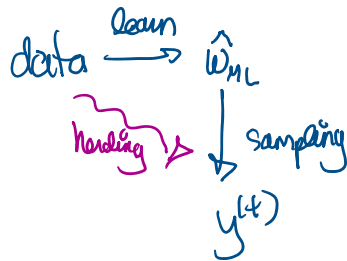
- amounts to running constant step size stochastic subgradient method on  $\mathcal{J}(w)$
- output  $\hat{w}_T = \frac{1}{T+1} \sum_{t=0}^T w_t$  (Polyak alg.) ( $R(w)=0$ )

↳ 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 "herding" alg. Welling & al. [ICML 2012]  
"3rd way to learn"



10n3)

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

suppose  $p(\tilde{y}|x; w) \propto \exp(\beta s(x, \tilde{y}; w))$   
*"inverse temperature parameter"* (pink arrow points to  $\beta$ )

Boltzmann dist.

$$\beta = \frac{1}{k_B T} \text{ temperature}$$

MCL → log-loss

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

*rescaling* (pink arrow points to  $\frac{1}{\beta}$ )

$Z_\beta(x; w)$  "partition fct."

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

"log-sum-exp" → "soft max" why?

"log-sum-exp" → "soft max" why?

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

$$\frac{1}{\beta} \log \left[ \exp(\beta s(\hat{y})) \left[ \sum_y \exp(\beta (s(y) - s(\hat{y}))) \right] \right] \leq 0$$

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

note:  
in deep learning book they call soft max

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

as  $\beta \rightarrow \infty$  (i.e. zero temp. limit)  $\frac{1}{\beta} \log \left( \sum_y \exp(\beta s(y)) \right) \xrightarrow{\beta \rightarrow \infty} \max_y s(y)$

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

### III) structured hinge loss

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

loss-augmented decoding

a) cartoon:  $\begin{matrix} \text{---} & s(y) \\ \text{---} & s(\tilde{y}_{\text{max}}) \\ \text{---} & \\ \text{---} & \end{matrix} \Rightarrow g^{\text{sum}}(x, y; w) = 0$   
 want:  $\geq l(y, \tilde{y}_{\text{max}})$

$$b) \boxed{g^{\text{sum}}(x, y; w) \geq l(y, h_w(x))}$$

why?  $g(x, y; w) = \max_{\tilde{y}} [s(\tilde{y}) + l(y, \tilde{y})] - s(y)$  let  $\hat{y} = \underset{y}{\text{argmax}} s(y) = h_w(x)$

$$\geq s(\hat{y}) + l(y, \hat{y}) - s(y)$$

if  $y \in \mathcal{Y}(x) \Rightarrow s(\hat{y}) \geq s(y)$

$$\geq l(\hat{y}) = l(y, h_w(x)) //$$

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

↳ hybrid between CRF loss & SUM loss

binary case: for structured hinge loss

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

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

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

$$\Leftrightarrow \langle \underbrace{w_+ - w_-}_{\tilde{w}}, \varphi(x) \rangle \geq 0$$

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

Let's show  $\sum_{\text{SVM}}(x, y; w) = [1 - y \langle \tilde{w}, x \rangle]_+$  where  $\tilde{w} = w_+ - w_-$

(recopy boring derivation here) "

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

(recopied from 2020 notes -- derivations for binary SVM:)

structured hinge loss

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

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

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

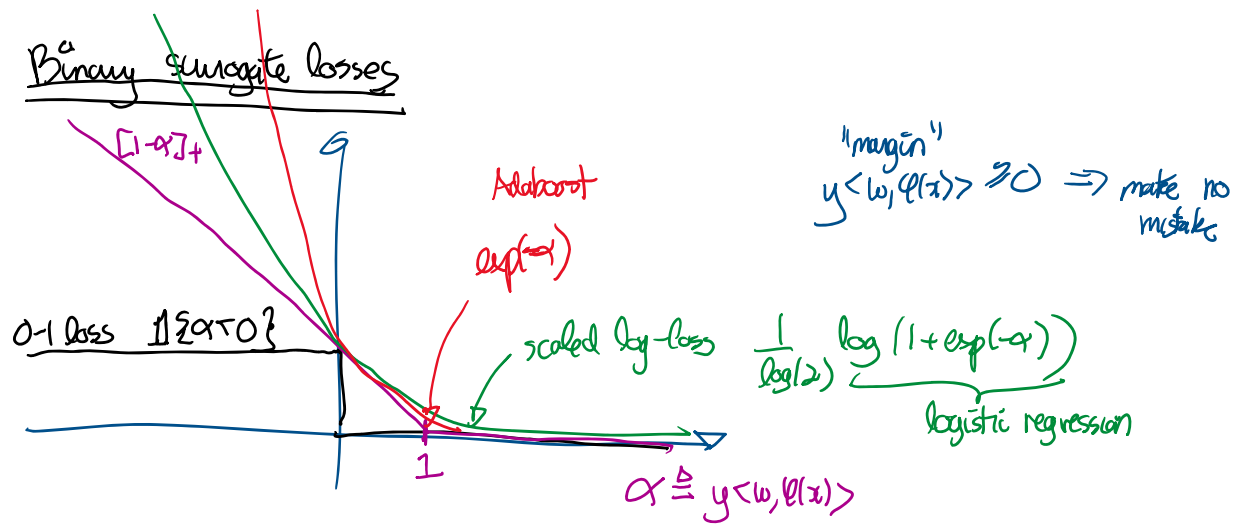
$$= \max \left\{ \langle \tilde{w}, x \rangle + 1, 1 - 1 \right\} + \langle w_-, x \rangle - \langle w_y, x \rangle$$

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

# Binary hinge losses



[Zweigert & El. 2006]  $\rightarrow$  showed all these methods are consistent