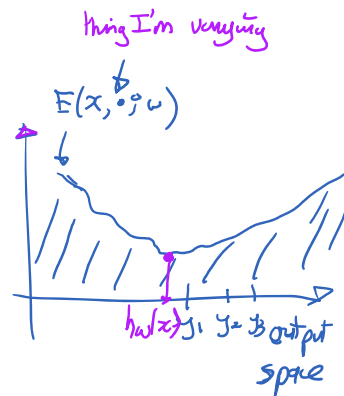


today: energy based methods & surrogate losses
multiclass

energy based methods: [LeCun et al. 2006]

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

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



Ingredients:

modeling {

- 1) what is $E(x, y; w)$? e.g. $S(x, y; w) = \langle w, \phi(x, y) \rangle$
or $E(x, y; w)$ output of a NN with x & y as input
- 2) how do you compute $\underset{y \in \mathcal{Y}(x)}{\operatorname{argmin}} E(x, y; w)$? (computationally) → "inference/decoding"

learning {

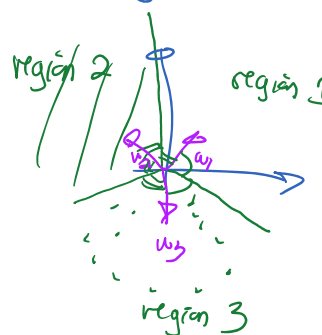
- 3) how to evaluate $E(x, y; w)$ on a training set? → surrogate loss $\mathcal{J}(w)$
in general: $\mathcal{J}(x^{(i)}, y^{(i)}, E(\cdot, \cdot; w))$ "loss functional"
- 4) how to minimize $\mathcal{J}(w)$ to learn w ? → optimization tricks (comp.)

flat multiclass case:

"flat" setting $h_w(x) = \underset{y}{\operatorname{argmax}} \langle w_y, \phi(x) \rangle$ $\in \mathbb{R}^{1 \times d}$

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

visually: $\|w_y\| = 1$



contrast this flat case with structured case:

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

e.g. OCR node feature map $\langle w, \varphi^{(node)}(x,y) \rangle = \sum_p \langle w, \varphi^{(node)}(x_p, y_p) \rangle$
 piece of sequence $\sum_p \mathbb{1}\{y_p = y_p'\} \langle w_{y_p}, \varphi(x) \rangle$
 → here "sharing" of parameters between different pieces of the joint labels → "structure"

aside: in structured prediction, usually absorbs "bias" in parameters
 standard binary classification $\text{sgn}(\langle w, x \rangle + b)$ → $\begin{pmatrix} \tilde{\varphi}(x) \\ 1 \end{pmatrix}$ $\langle \tilde{w}, \tilde{\varphi}(x) \rangle = \langle w, \varphi(x) \rangle + b$
 $\tilde{w} = \begin{pmatrix} w \\ b \end{pmatrix}$

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

Surrogate losses:

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

I) perception loss [Collins & al. 2002 EMNLP] score of ground truth

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

not needed

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

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

by using $\tilde{y} = y$

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

2) averaged perception alg: doesn't converge in general

• run constant step-size stochastic subgradient method on $\hat{\mathcal{L}}(w)$

• output $\hat{w}_T = \frac{1}{T} \sum_{t=0}^{T-1} w_t$ (Polyak avg.)

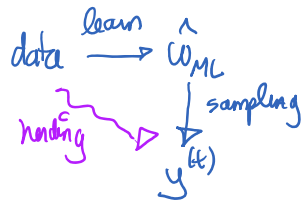
→ will likely converge to $w^* = 0$ when data is not separable

comments: 1) Collins's paper → he gives error bound and generalization error guarantees for perception

2) (aside) connection with the "harding" alg. by Welling & al.

"3rd way to learn" [see ICML 2012]

.. learn a



15h36

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

Boltzmann dist. in physics

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

suppose $p(y|x;w) \propto \exp(\beta s(x,y;w))$

"inverse temperature" parameter

MCL \rightarrow log-loss

$$g(x,y;w) = \underbrace{-\frac{1}{\beta}}_{\text{rescaling}} \log p(y|x;w) = \underbrace{-\frac{1}{\beta} \log \left(\frac{\exp(\beta s(x,y;w))}{\sum_{\tilde{y}} \exp(\beta s(x,\tilde{y};w))} \right)}_{Z_\beta(x;w) \text{ partition fun.}}$$

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

"log-sum-exp" \rightarrow "softmax" why?

$$\text{let } \hat{y} = \underset{\tilde{y}}{\text{argmax}} s(x,\tilde{y};w)$$

NOTE:

in deep learning book, they call this "softmax"

$$\left(\frac{\exp(\beta s(x,\tilde{y};w))}{\sum_{\tilde{y}} \exp(\beta s(x,\tilde{y};w))} \right)_{y \in \mathcal{Y}}$$

I call this "soft-argmax"

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

$$= s(x,\hat{y};w) + \frac{1}{\beta} \log(\text{stuff})$$

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

thus $\lim_{\beta \rightarrow \infty} \text{log-loss}(\beta) \rightarrow \text{perceptron loss}$

III) Structured hinge loss

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

"loss-augmented decoding"

a) cartoon:

$$\Rightarrow \mathcal{L}^{sum}(x, y; w) = 0$$

b) $f(x, y; w) \geq l(y, h_w(x))$

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 $\hat{y} = \text{argmax}_{\tilde{y} \in \mathcal{Y}(x)} s(\tilde{y}) = h_w(x)$

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

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

binary case

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

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

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

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

structured hinge loss

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

$\mathcal{L}^{sum}(x, y; w) = \max_{\tilde{y} \in \mathcal{Y}(x)} \{ \langle w_+, x \rangle + l(y, \tilde{y}), \langle w_-, x \rangle + l(y, \tilde{y}) \} - \langle w_y, x \rangle$

$w_+ = \tilde{w} + w_-$ $\mathbb{1}\{y \neq +1\}$ $1 - \mathbb{1}\{y \neq +1\}$

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

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

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

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

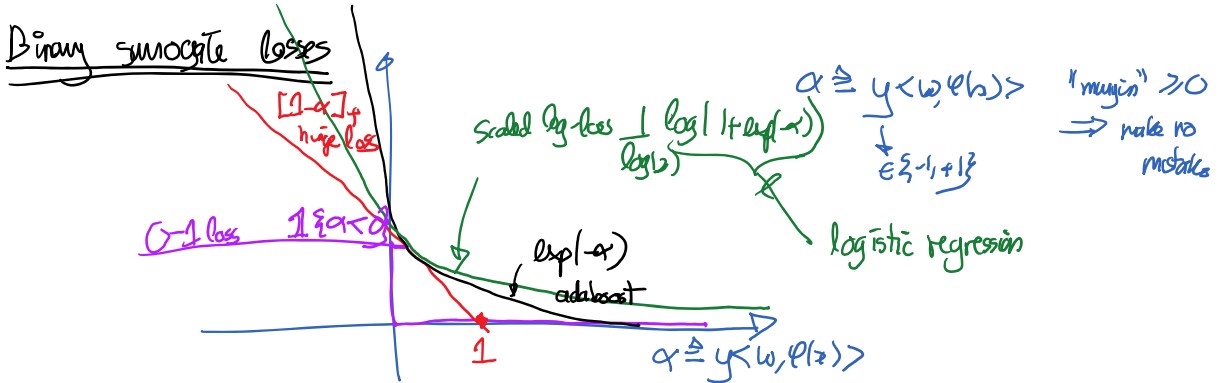
$$\text{case } y = -1 \rightarrow \max\{\langle \tilde{w}, x \rangle + 1, 0\} + 0 = [1 - y \langle \tilde{w}, x \rangle]_+$$

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

structured SVM
 $\mathcal{L}(x, y; w) = [1 - y \langle \tilde{w}, x \rangle]_+$

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

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



[Bartlett & al. 2006] → showed all three methods are consistent