

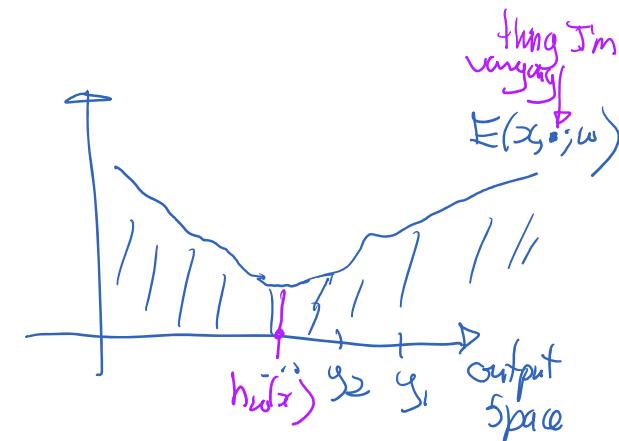
Lecture 3 - scribbles

Tuesday, January 15, 2019 14:12

today: energy based methods & surrogate losses
Multiclass

energy based methods & [Seznec et al. 2006]

$$\begin{aligned} \text{model: } h_w(x) &= \arg \min_{y \in \mathcal{Y}(x)} E(x, y; w) \quad \text{"energy fn."} \\ &= \arg \max_{y \in \mathcal{Y}(x)} S(x, y; w) \quad \text{"score / compatibility fn."} \end{aligned}$$



- ingredients:
- modeling { 1) what is $E(x, y; w)$? e.g. $S(x, y; w) = \langle w, \phi(x, y) \rangle$
 - 2) how do you compute $\arg \min_{y \in \mathcal{Y}(x)} E(x, y; w)$? \rightarrow "inference" / "decoding"
 - learning { 3) how to evaluate $E(x, y; w)$ on a training set? \rightarrow Surrogate Loss $\hat{S}(w)$
"quality" in general: $\hat{S}(x^{(i)}, y^{(i)}, E(\cdot, \cdot))$ "loss functional"
 - 4) how to minimize $\hat{S}(w)$ to learn w ? \rightarrow optimization tricks

flat multiclass case:

$$w \in \mathbb{R}^d$$

flat multiclass case:

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

equivalent to

$$\varphi(x, y) = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ \varphi(x) \\ 0 \end{pmatrix} \in \mathbb{R}^{d+k} \quad \begin{matrix} \# \text{ of classes} \\ \text{g-th position} \end{matrix}$$

OR node vs.
feature map

$$\langle w, \varphi(x, y) \rangle = \sum_p \underbrace{\langle w, \varphi^{(\text{node})}(x, y_p) \rangle}_{\substack{p \in \{y_p = y\} \\ y_p}} + \sum_{y_p \neq y} \langle w_{y_p}, \varphi(x) \rangle$$

→ here "sharing" of parameters between
different pieces of the labels → "structure"

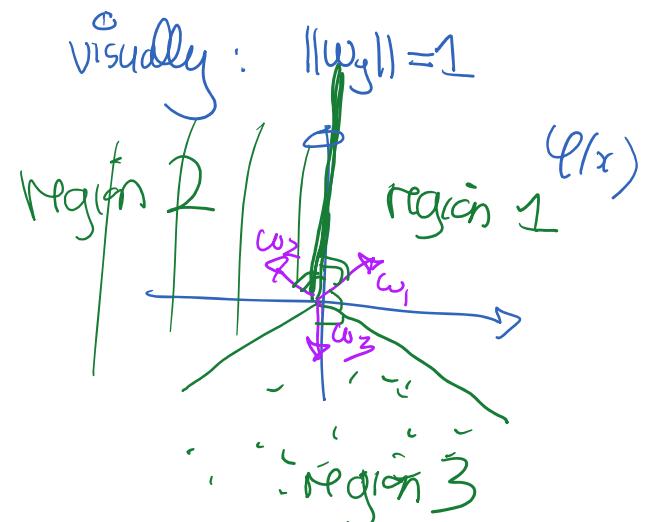
aside: in structured prediction, usually absorb "bias" in parameters

standard binary classification $\operatorname{sgn}\langle w, x \rangle + b$

$$\tilde{\varphi}(x) = \begin{pmatrix} \varphi(x) \\ 1 \end{pmatrix}$$

$$\begin{aligned} \langle \tilde{w}, \tilde{\varphi}(x) \rangle &= \langle w, \varphi(x) \rangle + b \\ \tilde{w} &= \begin{pmatrix} w \\ b \end{pmatrix} \end{aligned}$$

main question: regularizing or not the bias



Open question: Regularizing or not the bias
in structured prediction, does it matter?

surrogate losses:

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

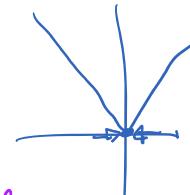
I) perceptron loss [Collins et al. 2002 EMNLP]

$$L(x, y; \omega) = \max_{\tilde{y} \in S(x)} s(x, \tilde{y}; \omega) - \underbrace{s(x, y; \omega)}_{\text{score of ground truth}} \quad (\text{assume } y \in S(x))$$

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

$$\max_{\tilde{y}} \underbrace{\langle \omega, \varphi(x, \tilde{y}) \rangle - \langle \varphi(x, y) \rangle}_{-\mathcal{L}_x(\tilde{y})} > \frac{\eta}{q} \circlearrowleft$$

by using $\tilde{y} = y$



Observations: 1) degenerate solution $\omega = 0$ or constant score over y

2) averaged perceptron alg.: → does not converge in general

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

- output $\hat{\omega}_T = \frac{1}{T+1} \sum_{t=0}^T \omega_t$ (Polyak avg.) → will converge to $\omega^* = 0$ when data is not separable

- Comments:
- 1) Collins' paper → he gives error bound and generalization error guarantees for perceptron
 - 2) (aside) connection with the "hacking" alg. by Welling et al.
"3rd way to learn" [see TCM 2012]



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

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

inverse temperature parameter

Boltzmann dist. in physics

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

MCL → log-loss

$$\mathcal{L}(x, y; \omega) = \underbrace{-\frac{1}{\beta} \log p(y|x; \omega)}_{\text{...}} = -\frac{1}{\beta} \log \left(\frac{\exp(\beta s(x, y))}{\sum \exp(\beta s(x, \tilde{y}))} \right) \} Z_\beta(x; \omega) \text{ partition}$$

$s(x, y; \omega)$

$\xrightarrow{\beta}$
rescaling

$$\beta \leftarrow \frac{\sum_{\tilde{y}} \exp(\beta s(\tilde{y}))}{Z_p(x; w)} \quad \left\{ Z_p(x; w) \text{ partition fn.} \right.$$

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

"log-sum-exp" \rightsquigarrow "soft-max" why?
Def $\hat{y} = \arg \max_y s(\tilde{y})$

NOTE:
in deep learning book

$$\left(\frac{\exp(s(\tilde{y}))}{\sum_{\tilde{y}} \exp(s(\tilde{y}))} \right)_{\text{year}}$$

I call this
"soft-argmax"

$\beta \rightarrow \infty$ (ie. zero temp limit)

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

$$= s(\hat{y}) + \frac{1}{\beta} \log \left(\underbrace{\text{stuff}}_{\leq 1/\beta} \right)$$

$$\frac{1}{\beta} \log \left(\sum_{\tilde{y}} \exp(\beta s(\tilde{y})) \right) \xrightarrow{\beta \rightarrow \infty} \max_{\tilde{y}} s(\tilde{y})$$

thus: $\lim_{\beta \rightarrow \infty} \text{log-loss}(\beta) \rightarrow \text{perception loss}?$

III) structured hinge loss

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

$$\mathcal{L}(x, y; \omega) = \max_{\tilde{y} \in \mathcal{Y}(x)} [L_S(x, \tilde{y}; \omega) + \ell(y|\tilde{y})] - s(x, y; \omega)$$

"loss-augmented decoding"

a)

cartoon:

$$\geq \ell(y, y_{\text{next}}) \quad \begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \quad \Rightarrow \quad \mathcal{L}_{\text{SUM}}(x, y; \omega) = 0$$

b) $\mathcal{L}(x, y; \omega) \geq \ell(y, h_w(x))$

why? $\mathcal{L}(x, y; \omega) = \max_{\tilde{y}} [s(\tilde{y}) + \ell(y|\tilde{y})] - s(y)$

$$\geq s(\hat{y}) + \ell(\hat{y}) - s(y) \quad \text{using } \hat{y} = \arg\max_{\tilde{y} \in \mathcal{Y}(x)} s(\tilde{y}) = h_w(x)$$

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

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

binary case:

$$y \in \{-1, +1\}$$

$$\omega = \begin{pmatrix} w_+ \\ w_- \end{pmatrix} \quad \varphi(x_1+1) = \begin{pmatrix} \varphi(x) \\ 0 \end{pmatrix}$$

$$h_w(x) = \arg\max \left\{ \langle w_+, x \rangle, \langle w_-, x \rangle \right\}$$

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

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

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

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

$$l_{\text{SVM}}(x, y; w) = \max \left\{ \underbrace{\langle w_+, x \rangle + l(y, +)}_{\substack{\text{if } y=+1 \\ \text{if } y=-1}}, \underbrace{\langle w_-, x \rangle + l(y, -)}_{\substack{\text{if } y=+1 \\ \text{if } y=-1}} \right\} - \langle w_y, x \rangle$$

$$\max \left\{ \langle \tilde{w}, x \rangle + \langle w_-, x \rangle + \begin{cases} 1 & y=+1 \\ -1 & y=-1 \end{cases}, \langle w_-, x \rangle + 1 - \begin{cases} 1 & y=+1 \\ -1 & y=-1 \end{cases} \right\} - \langle w_y, x \rangle$$

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

$$\stackrel{\text{case}}{=} y=+1 \quad \max \left\{ \langle \tilde{w}, x \rangle, 1 \right\} - \langle \tilde{w}, x \rangle = [1 - \langle \tilde{w}, x \rangle]_+$$

$$\Rightarrow y=-1 \quad \max \left\{ \langle \tilde{w}, x \rangle + 1, 0 \right\} + 0 = \max \{ 0, 1 - y \langle \tilde{w}, x \rangle \}$$

structured

SVM

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

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

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

i.e. structured hinge loss
reduces to binary SVM hinge loss
when using $l(y, y') = \begin{cases} 1 & y \neq y' \\ 0 & y = y' \end{cases}$
and $y = \begin{cases} 1 & y=+1 \\ -1 & y=-1 \end{cases}$

Binary surrogate losses

Different surrogate losses

