

today:

- surrogate loss
- generalization bounds
- consistency

last line: $\text{Spdist}(x, y, w) = \mathbb{E}_{\substack{\epsilon \sim \text{N}(0, I) \\ (x, y) \text{ invariant}}} [\ell(y, h_{w+\epsilon}(x))] \quad h_w(x) = \underset{y \in \mathcal{Y}}{\text{argmax}} s_w(x, y)$

convex surrogates:
so far

$$\begin{aligned} \text{perception}(x, y; w) &= \max_{\tilde{y} \in Y} s(x, \tilde{y}; w) - s(y) \\ &= \max_{\tilde{y} \in Y} [-n] \end{aligned}$$

Structural hinge loss
(SVMstruct)

they are upper
bounds on $Q(y, h_w(x))$

$$\max_{\tilde{y} \in \mathcal{Y}} [s(\tilde{y}) + \rho \min_{\tilde{y}} \{ \ell(y, \tilde{y}) - \eta \tilde{y} \}] - s(y) \quad \text{"margin rescaling"}$$

"margin rescaling"

$$\begin{aligned} \max_{\tilde{y} \in \mathcal{Y}} \ell(y, \tilde{y}) [1 + s(\tilde{y}) - s(y)] & \quad \text{"slack re-reading"} \\ \max_{\tilde{y}} \ell(y, \tilde{y}) [1 - m(\tilde{y})] & \end{aligned}$$

"slack rescaling"

$$\text{log-loss (CRF)} \quad \left[\frac{1}{P} \log \left(\sum_{\tilde{y}} \exp(\tilde{P} s(\tilde{y})) \right) - s(y) \right] \quad [-\log p_w(y|x)]$$

2) (soft-max)

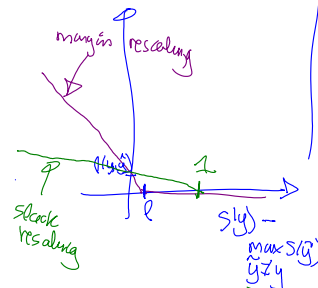
$\beta \rightarrow 10 \Rightarrow$ perceptron loss

$$\frac{1}{\beta} \log \left(\sum_y \exp(-\beta m(y)) \right)$$

"smooched huge"

$$\frac{1}{\beta} \log \left(\sum_{\tilde{y}} \exp(\beta [l(y, \tilde{y}) - m(\tilde{y})]) \right)$$

[e.g. Plešcher *et al.* 2010]



what are ^{theoretical} properties?

- a) generalization error bounds
- b) consistency properties of calibration functions

⊛ why structured score functions?

$$s(x, y) = \sum_{c \in \mathcal{C}} s_c(x, y_c)?$$

motivations similar to graphical models

1) Statistical efficiency: less # of parameters (simpler score functions S_L)
 \Rightarrow easier to learn see Data bias

2) computational \parallel to compute the $\arg\max_{\vec{y}} s(\vec{y})$ [Cortes & al. 2016]

but compare to what happens for Hamming loss:

given true conditional $q_x(y) \triangleq p(y|x)$

expected error of prediction \tilde{y} : $\mathbb{E}_{q_x(y)} [\ell(y, \tilde{y})] \triangleq \ell_{q_x}(\tilde{y})$

for Hamming loss: $\ell_{q_x}(\tilde{y}) = \mathbb{E}_{q_x(y)} \left[\frac{1}{2} \frac{1 - \mathbb{1}_{y=\tilde{y}}}{1 + \mathbb{1}_{y=\tilde{y}}} \right] = \frac{1}{2} (1 - q_x(\tilde{y}))$

\Rightarrow best decision is for each p , $\hat{y}_p = \arg\max_{\tilde{y}_p} p(\tilde{y}_p|x)$

\rightarrow (if no constraints), no need of "consistency" between parts
 - could just train independent models for each part marginal $p(y_p|x)$

"marginal decoding"

but \rightarrow this function might be too complicated

$$\left(\begin{array}{l} \min_x f(x) \\ \text{s.t. } \|x\|^2 \leq 1 \end{array} \right) \rightarrow \left(\begin{array}{l} f(x) + \lambda (\|x\|^2 - 1) \\ x_\lambda^* = \arg\min_x \end{array} \right)$$

claim: $x_\lambda^* \in \arg\min_{x: \|x\|^2 \leq \|x_\lambda^*\|^2} f(x)$

$$\min_x \max_{\lambda} f(x, \lambda) = \min_x f(x) \text{ s.t. constraints} = p^*$$

$$\max_{\lambda} \min_x f(x, \lambda) \leq p^*$$

weak duality

generalization error bounds:

for binary classification, a classical bound is
 with prob $1-\delta$ over D_n

$$\forall w \quad L_{\text{out}}(w) \leq L_n(w) + \frac{1}{\sqrt{n}} \sqrt{d \log \frac{1}{\delta} + \log \frac{2}{\delta}}$$

where d is VC-dimension
 of $\mathcal{H} = \{h_w : w \in W\}$

$\triangleq \max \{m : \# \text{ of different prediction functions on } m \text{ examples} \}$
 $= 2^m$
 for linear functions of p parameters
 $VC \text{-dim} = p+1$

bound above is true for any distribution

\Rightarrow too loose

motivates going to data distribution dependent measure of complexity

example: empirical Rademacher complexity

$$\hat{R}_n(\mathcal{H}) \triangleq \mathbb{E}_{\sigma} \left[\sup_{h \in \mathcal{H}} \left| \frac{2}{n} \sum_{i=1}^n \sigma_i 1\{y_i \neq h(x_i)\} \right| \right]$$

"1 correlations with random noise"

$\sigma_i \stackrel{i.i.d.}{\sim} \{-1, 1\}$ uniformly "Rademacher" R.V.

$$\text{then } \forall w \quad L(w) \leq \hat{L}_n(w) + \hat{R}_n(\mathcal{H}) + \frac{1}{\sqrt{n}} 3 \sqrt{\log \frac{2}{\delta}}$$

complexity depends on \mathcal{H}_n (implicitly on \mathcal{P})

"double sample trick" \rightarrow use second sample \mathcal{D}_n for

$$L(w) = \mathbb{E}_{\mathcal{D}_n} [\hat{L}_n(w)]$$

"Symmetrization trick"

\rightarrow bound the sup of differences between $L(w)$ & $\hat{L}_n(w)$

union bound as usual

Structured prediction generalization bounds [Cortes et al., NIPS 2016]

graphical model / factor graph

general loss $\ell(y, y')$ s.t. $\ell(y, y') \neq 0 \Leftrightarrow y \neq y'$; suppose $S(x, y) = \sum_{C \in \mathcal{C}} S(x, y_C)$
 $\mathcal{C} \subseteq \mathcal{C}_x$ a set of cliques

Thm. 7: with prob. $\geq 1 - \delta$

$$\forall w \in \mathcal{W} \quad L(w) \leq \sup_{y \in \mathcal{Y}} \ell(y, y') + 4\sqrt{2} \hat{R}_n^{\mathcal{C}}(Hw) + 3 \max_{y \in \mathcal{Y}} \sqrt{\log \frac{2}{\delta n}}$$

$$\text{where } \hat{R}_n^{\mathcal{C}} = \frac{1}{n} \mathbb{E}_{\sigma} \left[\sup_{w \in \mathcal{W}} \sum_{i=1}^n \sqrt{|\mathcal{C}|} \sum_{C \in \mathcal{C}} \sum_{y_C \in \mathcal{Y}_C} \sigma_{i,C,y} S_C(x_i, y_C; w) \right]$$

only uses $\{x^{(i)}\}_{i=1}^n$ "empirical factor graph complexity"

Thm. 2: if $S_C(x_i, y_C) = \langle w, \phi_C(x_i, y_C) \rangle$

and consider $W_\Delta \triangleq \{w : \|w\|_2 \leq \Delta\}$; let $R = \max_{i,c,y} \|\phi_c(x_i; y_c)\|_2$

$$\text{then } \hat{R}_{\text{on}}(H_{W_\Delta}) \leq \frac{R \Delta}{\sqrt{n}} |\mathcal{C}| \sqrt{\max_c |\mathcal{Y}_c|}$$

(so small cliques?)

$$\hat{S}_{\text{reg}}(w) + \underbrace{\left(\frac{R |\mathcal{C}|}{\sqrt{n}} \sqrt{\max_c |\mathcal{Y}_c|} \right)}_{\Delta} \underbrace{\|w\|_2}_{\Delta} \quad (\text{model selection})$$

SVM struct can be interpreted as minimizing upper bound on generalization error

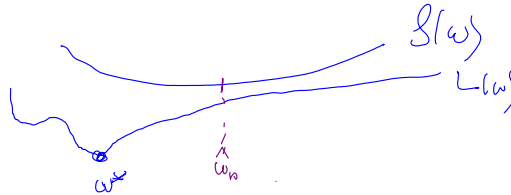
properties:

- minimize upper bound
- hope that minimize $L(w)$

- can evaluate bound b/gt guarantees

careful: minimizing an upperbound is not same as minimizing $L(w)$

→ next consistency



pointers

- sidenote: relationship between **constrained** and **regularized / penalized** formulations:
 - see section 4.7.3 (Pareto) and 4.7.4 (scalarization) of [Boyd's book](#) for the formal relationships
 - (in French): see exercise 2 from this homework from my ENS class: http://www.di.ens.fr/~slacoste/teaching/apprentissage-fall2015/TP/TP_5.pdf
- VC dimension / Rademacher complexity for binary case:
 - see slides from presentation by John Shawe-Taylor at MLSS 2009: http://mlg.eng.cam.ac.uk/mlss09/mlss_slides/ShawTaylor_1.pdf
 - VC dimension definition: slide 38
 - generalization error bound for binary classification with VC dimension: slide 46
 - Rademacher complexity: slide 85
 - generalization error bound with Rademacher complexity: slide 87
- structured prediction generalization bound:
 - Structured Prediction Theory Based on Factor Graph Complexity
 - Corinna Cortes, Vitaly Kuznetsov, Mehryar Mohri, Scott Yang
 - [NIPS 2016](#)