

IFT6132

TODO: fill survey <http://bit.ly/IFT6132-W23> by tomorrow!

structured prediction basis & setup

learning problem:

gives a training data dataset $D = (x^{(i)}, y^{(i)})_{i=1}^n$
 \downarrow \downarrow
 $\in X$ $\in Y$

goal: learn a prediction mapping $h_w: X \rightarrow Y$
 \downarrow
 parameter

that has a low classification error
 or
generalization error

$$L(w; P) \triangleq \mathbb{E}_{(x,y) \sim P} [\ell(y, h_w(x))]$$

test distribution

structured error function

statistical decision loss
 "risk" in ML
 (Vapnik risk)

[see lecture 5 of my PGM \rightarrow review statistical decision theory]

example of learning approach:

regularized ERM
 "empirical risk minimization"

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

not cts. in w \uparrow
regularizer

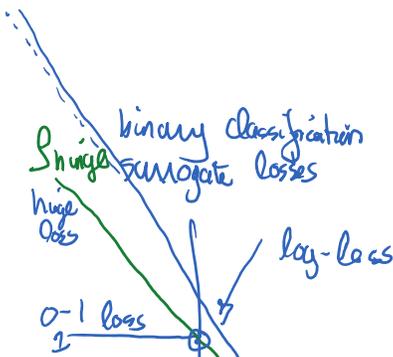
messy, non-convex
 LMP hard to minimize in general]

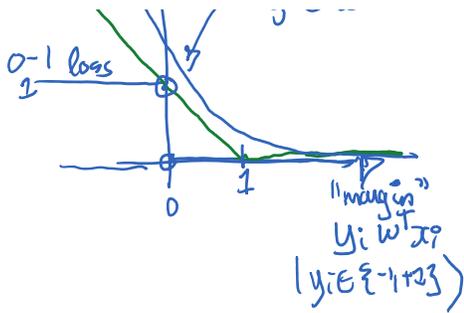
\Rightarrow replace \hat{L} with a surrogate loss / convex fct. \leftarrow statistics

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

\uparrow surrogate loss e.g. convex in w

M-estimator

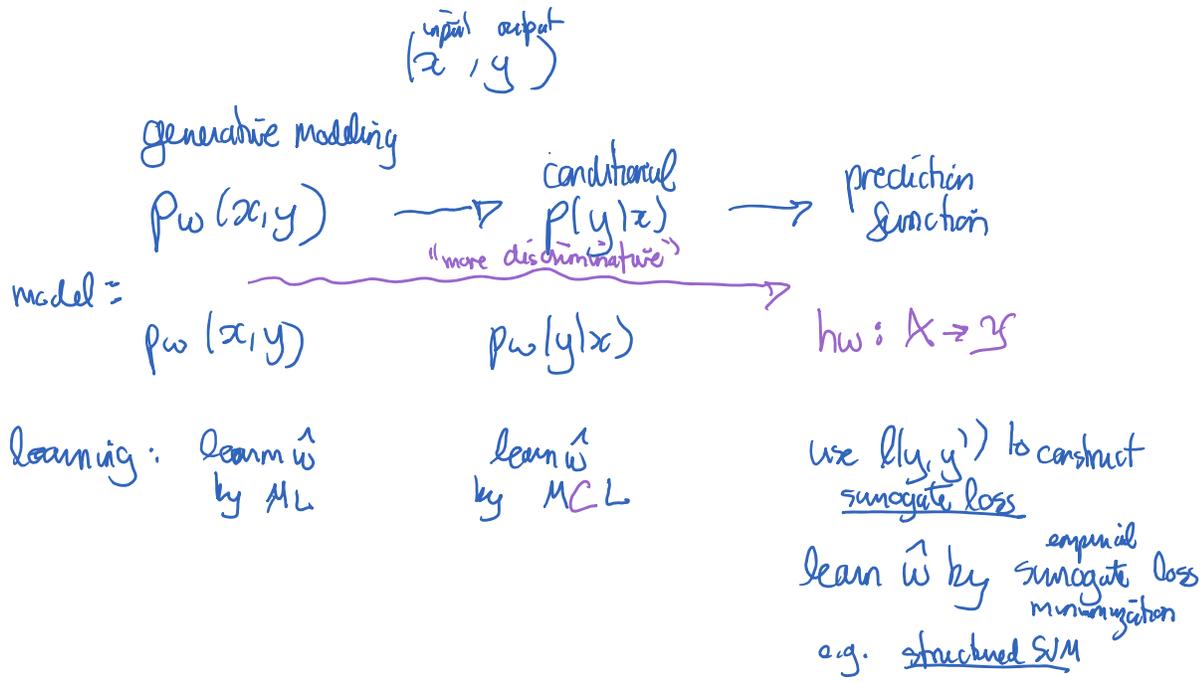




M-estimator
in statistics

- examples:
- structured hinge loss \rightarrow structured SVM
 - log-loss \rightarrow CRF
aka. "cross-entropy"

generative vs. discriminative learning continuum



\leftarrow more assumptions / less robust classification tasks

some important aspects of structured prediction (vs. binary classification)

- 1) Y output is usually exponentially big (in natural size of input)
- 2) structured error function $l(y, y')$ [$l(x, y, y')$]
e.g. Hamming loss $l(y, y') = \sum_p \mathbb{1}_{\{y_p \neq y'_p\}}$
- 3) sometimes constraints on pieces of y

\hookrightarrow e.g. consider word alignment example:

English words French words
 $\dots a_i \quad y_i \in \{0, 1, 2, \dots, j\}$

note: "need encoding fct."

$y = (y_1, \dots, y_p) \in \mathbb{R}^p$

$\dots \quad \dots \quad \dots$

