Supervised learning: Obs° $(x_i, y_i)$ $x_i \in \mathbb{R}^p$

$\mathbb{R}$     $\{1 \dots k\}$

Regres°    Classif°
(loss uses ordering) (loss invariant permut°)

Predic° func°    $y_i "\approx" f(x_i, \beta)$ ↳ parameter

Linear models: $f(x, \beta) = \langle x, \beta \rangle$

bias    $\langle x, \beta \rangle + c = \langle [x, 1], [\beta, c] \rangle$

non lin via lifting. $f(x, \beta) = \langle \varphi(x), \beta \rangle$
$(\leadsto \underset{p \to +\infty}{RKHS})$    $x \in \mathbb{R}^{p'}$ $\beta \in \mathbb{R}^{p'}$ $p' \gg p$
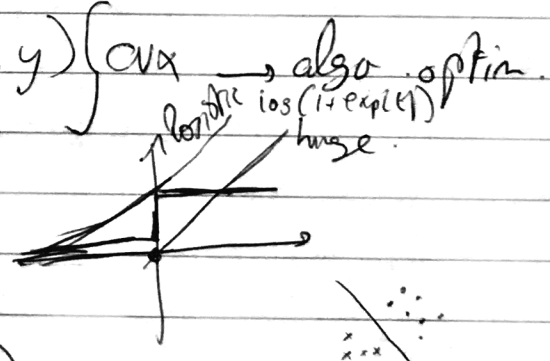
(P) ERM: $\underset{\beta \in \mathbb{R}^p}{\min} \frac{1}{n} \sum_{i=1}^{m} \ell(\langle x_i, \beta \rangle, y_i) + \lambda R(\beta)$

Regression: $\ell(y, y') = |y - y'|^2 / 2$
$\{-1, +1\}$

$\ell(\cdot, y)$ cvx $\longrightarrow$ algo optim.
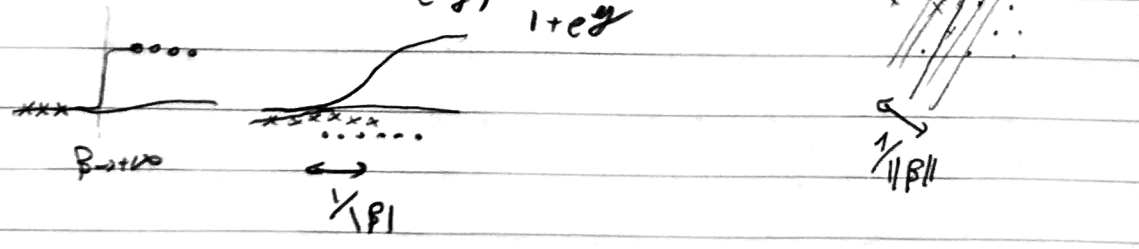$\mathbb{R}$    logistic $\log(1 + \exp(\cdot))$ hinge

Classif: $\ell(\bar{y}, y) = h(-\bar{y}y)$
2 classes

Hard Decis°: $sign(\langle x, \beta \rangle)$
Soft Decis°: $\theta(\langle x, \beta \rangle) = P(y \in \text{class } 1)$
$\theta(y) = \frac{e^y}{1 + e^y}$

$\langle x, \beta \rangle = 0$

$\beta \to +\infty$    $1/|\beta|$    $1/\|\beta\|$

$$\hat{\beta}_\lambda \triangleq \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \frac{1}{m} \sum_i (y_i - \langle \beta, x_i \rangle)^2 + \lambda \|\beta\|^2 = \frac{1}{n} \|y - X\beta\|^2 + \lambda \|\beta\|^2$$

$$= \underbrace{\left( \frac{X^\top X}{n} + \lambda \operatorname{Id}_p \right)^{-1}}_{\hat{C}} \underbrace{\left( \frac{X^\top y}{n} \right)}_{\hat{\mu}} \qquad X = \boxed{\overset{x_i}{\phantom{x}}} \in \mathbb{R}^{n \times p}$$

Hyp. $\quad x_i \overset{iid}{\sim} X \quad E(\|x\|^4)$

$\quad y_i \overset{iid}{\sim} Y \quad E(Y^4) \quad < +\infty \quad \Big\} \rightsquigarrow \text{TCL} \Rightarrow \begin{array}{c} \hat{C} \xrightarrow{n} \\ \hat{\mu} \xrightarrow[P, E^2]{} \end{array} \begin{array}{l} C \triangleq E(xx^\top) \in \mathbb{R}^{p \times p} \\ \mu \triangleq E(xy) \in \mathbb{R}^p \end{array}$

Refresher: cv of random variables. ⚡

**Consistency** : $\quad \hat{\beta}_\lambda \xrightarrow[P]{n \to +\infty} \beta_\lambda \triangleq (C + \lambda \operatorname{Id}_p)^{-1} \mu \quad \}$ random

$\qquad\qquad\qquad \overset{\lambda, n \to ?}{\underset{\text{Consistency}}{\nearrow}} \overset{??}{\underset{\searrow}{\downarrow}} \lambda \to 0 \qquad \}$ deterministic

$\qquad\qquad\qquad\qquad\qquad \beta_0 = C^\dagger \mu$

Refresher, pseudo-inverse : $\quad C^\dagger \mu = \underset{\beta}{\operatorname{argmin}} \|\beta\| \text{ st } C\beta = \mu$

· if $\ker C = \{0\}$, $C^\dagger = C^{-1}$. $\quad \underset{\text{eigen}}{\text{SVD}} \; C = U \operatorname{diag}(\sigma_i) U^* \overset{\sigma_i > 0}{\rightsquigarrow} C^\dagger = U \operatorname{diag}(\frac{1}{\sigma_i}) U^*$

**Example** : well specified model : $\quad y_i = \langle \beta^*, x_i \rangle + \varepsilon_i, \quad E(\varepsilon_i) = 0$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \overset{y \perp \text{iid.}}{\underset{\text{indep}}{}}$

$\mu \equiv \underset{xy}{E}(yx) = E_x(\langle \beta^*, x \rangle x) + \underbrace{E(y\varepsilon)}_{= E(y)E(x) = 0}$

$\qquad = E(xx^\top) \beta^* = C\beta^*$

$\rightsquigarrow \beta_0 = C^\dagger C \beta^* = \operatorname{Proj}_{\operatorname{Im}(C)} (\beta^*)$.

$\rightsquigarrow$ P$\text{dm}$ of $\beta^* \notin \operatorname{Im}(C)$ the part of $\beta^*$ in $\ker(C)$ is lost.
$\qquad\qquad \underset{\ker(C)^\perp}{\underbrace{\phantom{xxxxx}}}$

$\rightsquigarrow$ to recover the part of $\beta^*$ inside $\ker(C)$ needs non-linear methods (eg. $\ell^1$, aka LASSO)

defined on same proba. space!!

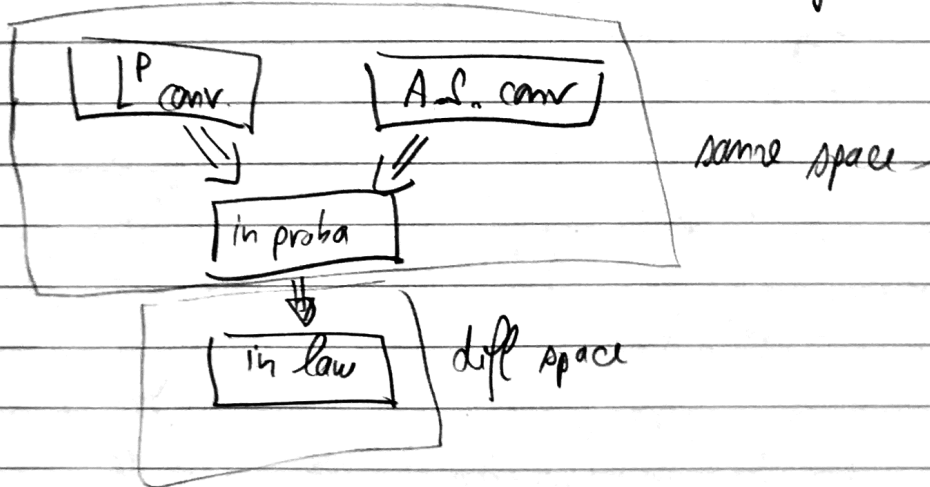Convergence in proba: $\beta_n \xrightarrow{\mathbb{P}} \beta_0 \iff \forall \varepsilon \quad \mathbb{P}(\|\beta_n - \beta_0\| \geq \varepsilon) \xrightarrow{n \to \infty} 0$

Quantified: $\|\beta_n - \beta_0\| = O_{\mathbb{P}}\left(\frac{1}{n^\kappa}\right)$ : "with proba $1 - \eta$, $\|\beta_n - \beta_0\| \leq C(\eta) \times \frac{1}{n^\kappa}$"

$\quad\quad$ very often $\sim \log(1/\eta)$

$L^p$-Convergence in expect$^{\circ}$: $\mathbb{E}(\|\beta_n - \beta_0\|^p) \xrightarrow{n \to \infty} 0$

Rmq: Markov ineq: $\mathbb{P}(\|\beta_n - \beta_0\| \geq \varepsilon) \leq \frac{1}{\varepsilon^p}\mathbb{E}(\|\beta_n - \beta_0\|^p)$

$\quad L^p$ convergence stronger $(\Rightarrow)$ convergence in proba

Other convergence: Almost sure, in law (aka optimal transport).

$\quad \hookrightarrow$ need not be def. on same space!



$L^p$ conv.  $\quad$ A.S. conv  $\quad\quad\quad$ same space
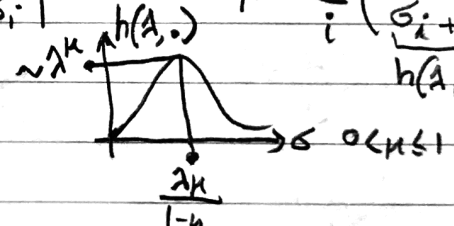
in proba

in law  $\quad$ diff space

$Dp°$: $\hat{\beta}_\lambda \triangleq (\hat{C}+\lambda Id)^{-1}\hat{u}$ ; $\tilde{\beta}_\lambda \triangleq (C+\lambda Id)^{-1}\hat{u}$ ; $\beta_\lambda \triangleq (C+\lambda Id)^{-1}u$

**VARIANCE**

$\hat{\beta}_\lambda - \beta_0 = \tilde{\beta}_\lambda - \beta_\lambda \xrightarrow{②} (C+\lambda Id)^{-1}(\hat{u}-u)$     $\boxed{\delta \triangleq \frac{1}{\sqrt{n}}}$

in proba  $\|\cdot\| \leq \frac{1}{\lambda}$   $\|\cdot\| \leq \delta$  ①

$+\hat{\beta}_\lambda - \tilde{\beta}_\lambda \xrightarrow{①} [(\hat{C}+\lambda Id)^{-1}-(C+\lambda Id)^{-1}]\times\hat{u}$

in proba   $\sim \frac{\delta}{\lambda}$   (why not $\frac{\delta}{\lambda^2}$ ??) ②

**BIAS**

$+\beta_\lambda - \beta_0 \xleftarrow{③} [(C+\lambda Id)^{-1}-C^+]u$

$\underbrace{\phantom{xxxxxxxxxxxxxx}}$

need hyp. on $\beta_0$

$\overset{(approx°)}{\underline{\text{lemma}}}$ : if $\beta_0 \in Im(C^\mu) \overset{0\leq\mu\leq 1}{\phantom{x}} \text{ ie } \exists q, \beta_0 = C^\mu q$, then $\|\beta_\lambda - \beta_0\| \leq c(k)\cdot\|q\|\lambda^\mu$ ③'

$\underline{\text{Proof}}$ : $\|\beta_\lambda - \beta_0\|^2 = \sum_i \left|\frac{1}{\sigma_i+\lambda} - \frac{1}{\sigma_i}\right|^2 \times \sigma_i^{2(\mu+1)} \times q_i^2 = \sum_i \left(\frac{\lambda\sigma_i^\mu}{\sigma_i+\lambda}\right)^2\cdot u_i^2$

$\boxed{u = C\beta_0 = C^{\mu+1}q}$     $\underbrace{}_{h(\lambda,\sigma)}$

$\sim \lambda^\mu$   $h(\lambda,\cdot)$



$\frac{\lambda\mu}{1-\mu}$   $\sigma$   $0<\mu\leq 1$

$\underline{\text{Thm [EST°]}}$ : (with high proba) $\|\hat{\beta}_{\lambda_n}-\beta_0\| = O_P\left(\left(\frac{1}{\sqrt{n}}\right)^{\frac{\mu}{1+\mu}}\right)$  $\lambda_n \triangleq n^{-\frac{1/2}{1+\mu}}$

$\underline{\text{Proof}}$ : ①+②+③ $\Rightarrow \|\hat{\beta}_\lambda - \beta_0\| = O\left(\frac{\delta}{\lambda}+\lambda^\mu\right) \rightsquigarrow$ choose $\lambda = \delta^{\frac{1}{1+\mu}}$

$\underline{\text{Rates}}$ : $\overset{\mu=1}{\underset{\mu=1/2}{\nearrow}} \begin{array}{l}1/2 \text{ (optimal)} \\ 1/3 \text{ (usual one)}\end{array}$     balance   $\underline{\text{Rem}}$: to have faster (eg linear) rates, needs NL mth (e+)

---

$\underline{\text{Prediction}}$ : $\underset{x}{\mathbb{E}}(\langle x,z\rangle^2) = \langle Cz,z\rangle$     $\mu\leq 1/2$

③ $\rightsquigarrow \langle C(\beta_\lambda-\beta_0),\beta_\lambda-\beta_0\rangle = \sum_i \left(\frac{\lambda\sigma_i^{\mu+1/2}}{\sigma_i+\lambda}\right)^2 \mu_i^2 = O\left(\|u\|^2\cdot\lambda^{2\mu+1}\right)$ Ⓐ

② $\rightsquigarrow \langle C(\tilde{\beta}_\lambda-\beta_\lambda),\tilde{\beta}_\lambda-\beta_\lambda\rangle = \sum_i \sigma_i \frac{1}{(\sigma_i+\lambda)^2}\times(\hat{u}_i-u_i)^2 = O\left(\frac{\|\hat{u}-u\|^2}{\lambda}\right) = O\left(\frac{1}{n\lambda}\right)$ Ⓑ



$\frac{\sigma}{(\sigma+\lambda)^2}$
$\frac{1}{\lambda}$
$\sigma=\lambda$

① $\rightsquigarrow \langle C(\hat{\beta}_\lambda-\tilde{\beta}_\lambda),(\hat{\beta}_\lambda-\tilde{\beta}_\lambda)\rangle = ?? = O\left(1/n\lambda\right)$ Ⓒ

$\underline{\text{Thm}}$  $\mathbb{E}(\langle x,\hat{\beta}_\lambda-\hat{\beta}_0\rangle^2)^{1/2} = O\left((1/\sqrt{n})^{\frac{\mu+1/2}{\mu+1}}\right)$  $\lambda_n \triangleq n^{-\frac{1/2}{\mu+1}}$  $\underline{\text{Rate}} \overset{\mu=1}{\underset{\mu=1/2}{\nearrow}} \begin{array}{l}3/4 \\ 2/3\end{array}$

$\underline{\text{Proof}}$ Ⓐ Ⓑ Ⓒ $\rightarrow \langle C(\hat{\beta}_\lambda-\beta_0),\hat{\beta}_\lambda-\beta_0\rangle^{1/2} = O\left(\frac{\delta}{\sqrt{\lambda}}+\lambda^{\mu+1/2}\right) \rightarrow \lambda = \delta^{\frac{1}{\mu+1}}$

Discussion on the source cond$^{\underline{c}}$ : $\beta_0 = C^\mu q = U \, diag(\sigma_i^\mu) U^T q$.

- In finite dimension, $\beta_0 \in Im(C) \Rightarrow \forall \mu, \, \beta_0 \in Im(C^\mu)$ !

  indeed, $\beta_0 = C^\mu q$ for $q = (C^+)^\mu \beta_0$ !

  But can be very bad bound $\|q\| \sim \frac{1}{\lambda_{min}^\mu} \ldots$

  $\lambda_{min}$ can goes to 0 very fast with dim$^\circ$ $p$.

  so the goal is to have $\|q\|$ small

- In $\infty$ dim$^\circ$ (eg RKHS) , $\beta_0 \in Im(C^\mu)$ not always true, typically $\lambda_{min} = 0$, $\sigma_i \underset{i \to \infty}{\longrightarrow} 0$ !

- Rule of thumb: if C trans$^\circ$ invariant (convol$^\circ$), $U = fft$, and $\sigma_i$ small for high freq, so that large $\mu$ corresponds to smoother $\beta_0$ ($\sim$ SOBOLEV space).

## COMPARISON w/k IP

$$y = X\beta_0 + \varepsilon \qquad \delta = \|\varepsilon\|$$

$\underline{\text{EST}^\circ}: \|\hat{\beta}_\lambda - \beta_0\| = O\left(\frac{\delta}{\sqrt{\lambda}} + \lambda^\mu\right) = O\left(\delta^{\frac{\mu}{\mu + 1/2}}\right).$

$\underline{\text{PRED}^\circ}: \frac{1}{\sqrt{n}}\|X(\hat{\beta}_\lambda - \beta_0)\| = O\left(\delta + \lambda^\mu\right) = O(\delta).$ (trivial bound)

indep $\lambda$ !!