# Regression with Squared Loss

$(X, Y)$ in $\mathcal{X} \times \mathbb{R}$

Min. mean square Error (==MMSE==) prediction:

$$\min_{f: \mathcal{X} \to \mathbb{R}} \underbrace{\mathbb{E}\left[(Y - f(X))^2\right]}_{:= L(f)}$$

- optimal predictor if $P_{XY}$ known:

$$f^*(x) = \mathbb{E}[Y | X = x]$$

$$L(f) = L(f^*) + \underbrace{\mathbb{E}\left[(f(x) - f^*(x))^2\right]}_{\geqslant 0}$$

- Learning setting: $P_{XY}$ unknown

$(X_1, Y_1), \ldots, (X_n, Y_n)$ iid

fix $\mathcal{F}$ (class of fcns $f: \mathcal{X} \to \mathbb{R}$)

$\hat{f}_N$ based on data, in $\mathcal{F}$

<u>Goal:</u> $L(\hat{f}_n) - \inf_{f \in \mathcal{F}} L(f)$ small w.h.p.

<u>Set-up:</u> $\mathcal{X}$ arbitrary

$$\mathcal{P} = \left\{ P_{XY} : \mathbb{P}\left[|Y| \leq M\right] = 1 \right\}$$

$\mathcal{F}$ : subset of RKHS $\mathcal{H}_k$

choice of $K$ is a degree of freedom

e.g. $\mathcal{X} = \mathbb{R}^d$, $K(x, x') = 1 + \langle x, x' \rangle$

$\mathcal{H}_K$ consists of $f(x) = \langle w, x \rangle + b$

---

① **Norm - Constrained Predictors**

$$\mathcal{F} = \mathcal{F}_\lambda = \{ f \in \mathcal{H}_K : \|f\|_K \le \lambda \}$$

$0 < \lambda < \infty$ : fixed parameter

ERM: $\displaystyle \hat{f}_n = \operatorname*{argmin}_{f \in \mathcal{F}_\lambda} \underbrace{\frac{1}{n} \sum_{i=1}^{n} (Y_i - f(X_i))^2}_{L_n(f)}$

$$L(\hat{f}_n) - L^*(\mathcal{F}_\lambda)$$

$$= L(\hat{f}_n) - \inf_{f \in \mathcal{F}_\lambda} L(f)$$

$$= L(\hat{f}_n) - L(f^*) \qquad\qquad f^* = \operatorname*{argmin}_{f \in \mathcal{F}_\lambda} L(f)$$

$$\le 2 \sup_{f \in \mathcal{F}_\lambda} | L_n(f) - L(f) |$$

---

Notation:

$$\ell(y, u) := (y - u)^2$$

$$\ell \bullet f(x, y) := \ell(y, f(x))$$

$$= (y - f(x))^2$$

$$\ell \bullet \mathcal{F}_\lambda := \{ \underbrace{\ell \bullet f}_{\text{losses}} : \underbrace{f \in \mathcal{F}_\lambda}_{\text{predictors}} \}$$

$$\sup_{f \in \mathcal{F}_\lambda} |L_n(f) - L(f)| = \sup_{f \in \mathcal{F}} |P_n(\ell \bullet f) - P(\ell \bullet f)|$$

$$= \Delta_n (\ell \bullet \mathcal{F}_\lambda)$$

<span style="color:green">induced losses</span>

where $P_n(h) := \frac{1}{n} \sum_{i=1}^{\sim} h(z_i)$ $\qquad z_i = (x_i, y_i)$

$\qquad\qquad P(h) := \mathbb{E} [h(z)]$

$$h = \ell \bullet f \qquad \text{for some} \quad f \in \mathcal{F}_\lambda$$

$$g(z^n) = \Delta_n (\ell \bullet \mathcal{F}_\lambda)$$

$$= \sup_{f \in \mathcal{F}_\lambda} \left| \frac{1}{n} \sum_{i=1}^{n} \ell \bullet f(z_i) - \mathbb{E} [\ell \bullet f(z)] \right|$$

<u>Claim:</u> $g(\cdot)$ has bdd diffs (McDiarmid)

$z_1, \dots, z_i, \dots, z_n$ $\qquad\qquad z_i = (x_i, y_i) \in \mathcal{X} \times (-M, M]$

$z_1, \dots, z_i', \dots, z_n$ $\qquad\qquad z_i' = (x_i', y_i') \in \mathcal{X} \times (-M, M]$

$$g(z_1, \dots, z_i, \dots, z_n) - g(z_1, \dots, z_i', \dots, z_n)$$

$$\leq \frac{1}{n} \sup_{f \in \mathcal{F}_\lambda} \left| (y_i - f(x_i))^2 - (y_i' - f(x_i'))^2 \right|$$

$$\leq \frac{1}{n} \sup_{x \in \mathcal{X}} \sup_{|y| \leq M} \sup_{f \in \mathcal{F}_\lambda} (y - f(x))^2$$

$$\leq \frac{2}{n} \left( M^2 + \sup_{f \in \mathcal{F}_\lambda} \sup_{x \in \mathcal{X}} |f(x)|^2 \right) \qquad \begin{array}{l} \color{green}(a+b)^2 \\ \color{green}\leq 2a^2 + 2b^2 \end{array}$$

$$K : \mathcal{X} \times \mathcal{X} \to \mathbb{R} \qquad \text{Mercer kernel}$$

$$C_K := \sup_{x \in \mathcal{X}} \sqrt{K(x,x)} < \infty$$

$$f \in \mathcal{H}_K : \qquad \|f\|_\infty := \sup_{x \in \mathcal{X}} |f(x)|$$

**Lemma** $\qquad \|f\|_\infty \leq C_K \|f\|_K$

**Proof** $\qquad f \in \mathcal{H}_K$

$$
\begin{aligned}
|f(x)| &= |\langle f, K_x \rangle_K| & \text{(repr. prop.)} \\
&\leq \|f\|_K \|K_x\|_K & \text{(Cauchy-Schwarz)} \\
&= \|f\|_K \sqrt{K(x,x)} \\
&\leq C_K \|f\|_K.
\end{aligned}
$$

$$\sup_{f \in \mathcal{F}_\lambda} \sup_{x \in \mathcal{X}} |f(x)|^2 = \sup_{f : \|f\|_K \leq \lambda} \|f\|_\infty^2$$

$$\leq C_K^2 \lambda^2$$

$\Rightarrow g(z^n) = \Delta_n(\ell \bullet \mathcal{F}_\lambda)$ has bdd diff.

with $\qquad c_1 = \cdots = c_n = \frac{2}{n}(M^2 + C_K^2 \lambda^2)$

McDiarmid: $\qquad \forall t > 0$

$$\mathbb{P}\left\{ \Delta_n(\ell \bullet \mathcal{F}_\lambda) \geq \mathbb{E}\,\Delta_n(\ell \bullet \mathcal{F}_\lambda) + t \right\} \leq \exp\left( -\frac{n\,t^2}{2(M^2 + C_K^2 \lambda^2)^2} \right)$$

- $t = \sqrt{\dfrac{2}{n}(M^2 + C_K^2 \lambda^2)^2 \log\left(\frac{1}{\delta}\right)}$

$= (M^2 + C_K^2 \lambda^2)\sqrt{\dfrac{2 \log(1/\delta)}{n}}$

$$\left| \quad \Delta_n(\ell \bullet \mathcal{F}_\lambda) \leq \mathbb{E}\Delta_n(\ell \bullet \mathcal{F}_\lambda) + (M^2 + C_K^2 \lambda^2)\sqrt{\dfrac{2\log(1/\delta)}{n}} \right.$$

$\left| \quad \text{w.p.} \geq 1 - \delta \right.$

- Symmetrization:

$\mathbb{E}\Delta_n(\ell \bullet \mathcal{F}_\lambda) \leq 2\mathbb{E}R_n(\ell \bullet \mathcal{F}_\lambda(z^n))$

$\ell \bullet \mathcal{F}_\lambda(z^n) = \left\{ (\ell \bullet f(z_1), \ldots, \ell \bullet f(z_n)) : f \in \mathcal{F}_\lambda \right\}$

$= \left\{ ((y_1 - f(x_1))^2, \ldots, (y_n - f(x_n))^2) : f \in \mathcal{F}_\lambda \right\}$

Claim: for $f \in \mathcal{F}_\lambda$,

$\qquad |y - f(x)| \leq M + C_K \lambda$

for all $x \in \mathcal{X}$, $|y| \leq M$.

Proof: $\qquad |y - f(x)| \leq |y| + |f(x)|$

$\qquad\qquad\qquad \leq M + \|f\|_\infty$

$\qquad\qquad\qquad \leq M + C_K \lambda.$ 📖

- Apply contraction principle w/ $\varphi(t) = t^2$

on $[-(M + C_K \lambda), \underbrace{M + C_K \lambda}_{:= A}]$

$$\left( (Y_1 - f(X_1))^2, \ldots, (Y_n - f(X_n))^2 \right)$$

$$= \left( \varphi(Y_1 - f(X_1)), \ldots, \varphi(Y_n - f(X_n)) \right)$$

<span style="color:green">$\in [-(M + C_k \lambda), M + C_k \lambda]$</span>

$\longrightarrow$ Lip-const. on $[-A, A]$ $\overset{f, \varphi}{\downarrow}$

$$\Rightarrow R_n(\ell \bullet F_\lambda(z^n)) \leq 2 \cdot 2A \cdot \mathbb{E}_\varepsilon \left[ \frac{1}{n} \sup_{f \in F_\lambda} \left| \sum_{i=1}^{n} (Y_i - f(X_i)) \varepsilon_i \right| \right]$$

$$\frac{1}{n} \mathbb{E}_\varepsilon \left[ \sup_{f \in F_\lambda} \left| \sum_{i=1}^{n} \varepsilon_i (Y_i - f(X_i)) \right| \right]$$

$$\leq \frac{1}{n} \mathbb{E}_\varepsilon \left[ \left| \sum_{i=1}^{n} \varepsilon_i Y_i \right| \right] + \frac{1}{n} \mathbb{E}_\varepsilon \left[ \sup_{f \in F_\lambda} \left| \sum_{i=1}^{n} \varepsilon_i f(X_i) \right| \right]$$

$$=: \frac{1}{n} \left( \boxed{I} \qquad\qquad + \qquad \boxed{II} \right)$$

where:

$\boxed{I} \quad \mathbb{E}_\varepsilon \left[ \left| \sum_{i=1}^{n} \varepsilon_i Y_i \right| \right]$ $\qquad\qquad |Y_i| \leq M$

$$= \mathbb{E}_\varepsilon \sqrt{\left( \sum_{i=1}^{n} \varepsilon_i Y_i \right)^2}$$

$$\leq \sqrt{\mathbb{E}_\varepsilon \left( \sum_{i=1}^{n} \varepsilon_i Y_i \right)^2}$$

$$= \sqrt{\sum_{i=1}^{n} Y_i^2}$$

$$\leq M \sqrt{n}$$

$$\text{II} \quad \mathbb{E}_{\varepsilon}\left[\sup_{f \in \mathcal{F}_{\lambda}} \left| \sum_{i=1}^{n} \varepsilon_i f(X_i) \right| \right]$$

$$= \mathbb{E}_{\varepsilon}\left[\sup_{f \in \mathcal{F}_{\lambda}} \left| \sum_{i=1}^{n} \varepsilon_i \langle f, K_{X_i} \rangle_{K} \right| \right]$$

$$= \mathbb{E}_{\varepsilon}\left[\sup_{f \in \mathcal{F}_{\lambda}} \left| \langle f, \sum_{i=1}^{n} \varepsilon_i K_{X_i} \rangle_{K} \right| \right]$$

$$\leq \mathbb{E}_{\varepsilon} \left\| \sum_{i=1}^{n} \varepsilon_i K_{X_i} \right\|_{K} \cdot \sup_{f \in \mathcal{F}_{\lambda}} \| f \|_{K}$$

$$\leq \lambda \cdot \mathbb{E}_{\varepsilon} \left\| \sum_{i=1}^{n} \varepsilon_i K_{X_i} \right\|_{K}$$

$$= \lambda \cdot \mathbb{E}_{\varepsilon} \sqrt{\langle \sum_{i=1}^{n} \varepsilon_i K_{X_i}, \sum_{i=1}^{n} \varepsilon_i K_{X_i} \rangle_{K}}$$

$$\leq \lambda \, c_K \sqrt{n}$$

$$\Rightarrow \quad R_n(\ell \cdot \mathcal{F}_{\lambda}) \leq 4(M + c_K \lambda) \cdot \frac{1}{n}\left( \text{I} + \text{II} \right)$$

$$\leq \frac{4(M + c_K \lambda)^2}{\sqrt{n}}$$

$$\mathbb{E} \Delta_n(\ell \cdot \mathcal{F}_{\lambda}) \leq 2 \, \mathbb{E} R_n(\ell \cdot \mathcal{F}_{\lambda})$$

$$\leq \frac{8(M + c_K \lambda)^2}{\sqrt{n}}$$

$$\therefore \quad \text{w.p.} \geq 1 - \delta,$$

$$L(\hat{f}_n) \leq L^*(\mathcal{F}_\lambda) + \frac{16 (M + c_K \lambda)^2}{\sqrt{n}}$$

$$+ (M^2 + c_K^2 \lambda^2) \sqrt{\frac{8 \log(1/\delta)}{n}}.$$

Comments:
- ERM over a ball $\mathcal{F}_\lambda$ in $\mathcal{H}_K$
- radius $\lambda$: hard constraint on hypothesis space complexity

② <u>Norm - penalized predictors</u>

$$\mathcal{F} = \mathcal{H}_K \quad (\text{entire RKHS})$$

$$\hat{f}_n = \underset{f \in \mathcal{H}_K}{\text{argmin}} \left\{ L_n(f) + \gamma \| f \|_K^2 \right\}$$

$$\quad (\gamma > 0: \text{tunable parameter})$$

$$J_\gamma(f) := L(f) + \gamma \| f \|_K^2$$

$$J_{n,\gamma}(f) := L_n(f) + \gamma \| f \|_K^2$$

$$= \frac{1}{n} \sum_{i=1}^{n} (Y_i - f(x_i))^2 + \gamma \| f \|_K^2$$

$$L(f) \leq J_\gamma(f)$$

$$L_n(f) \leq J_{n,\gamma}(f)$$

**Observation:** suppose $f \in \mathcal{H}_K$ minimizes $J_\gamma(f)$ or $J_{n,\gamma}(f)$. Then $\|f\|_K \leq \frac{M}{\sqrt{\gamma}}$.

**Corollary**
$$\min_{f \in \mathcal{H}_K} J_\gamma(f) = \min_{f \in \mathcal{F}_{M/\sqrt{\gamma}}} J_\gamma(f)$$

$$\min_{f \in \mathcal{H}_K} J_{n,\gamma}(f) = \min_{f \in \mathcal{F}_{M/\sqrt{\gamma}}} \overline{J}_{n,\gamma}(f)$$

**Proof**

Assume $f_\gamma^* = \operatorname*{argmin}_{f \in \mathcal{H}_K} J_\gamma(f)$

$$J_\gamma(f_\gamma^*) \leq J_\gamma(0) = L(0) = \mathbb{E}|Y|^2 \leq M^2$$

$$J_\gamma(f) \geq \gamma \|f\|_K^2 \qquad \forall f$$

$$\Rightarrow \quad \|f_\gamma^*\|_K^2 \leq \frac{M^2}{\gamma}$$

$$J_\gamma(\widehat{f}_n) - J_\gamma(f_\gamma^*) \qquad \qquad \widehat{f}_n, f_\gamma^* \in \mathcal{F}_{M/\sqrt{\gamma}}$$

$$\leq 2 \sup_{f \in \mathcal{F}_{M/\sqrt{\gamma}}} |J_{n,\gamma}(f) - J_\gamma(f)|$$

$$= 2 \sup_{f \in \mathcal{F}_{M/\sqrt{\gamma}}} |L_n(f) + \gamma \|f\|_K^2 - L(f) - \gamma \|f\|_K^2|$$

$$= 2 \sup_{f \in \mathcal{F}_{M/\sqrt{\gamma}}} |L_n(f) - L(f)|$$

— already did this w/ $\lambda > 0$ arbitrary

— take $\lambda = M/\sqrt{\gamma}$

$\therefore$ w.p. $\geq 1 - \delta$,

$\bullet \; J_\gamma(\hat{f_n}) \leq J_\gamma(f_\gamma^*) + \dfrac{16\left(M + c_k \frac{M}{\sqrt{\gamma}}\right)^2}{\sqrt{n}}$

$$+ \left(M^2 + c_k^2 \frac{M^2}{\gamma}\right)\sqrt{\dfrac{8\log(1/\delta)}{n}}$$

$\bullet \; L(\hat{f_n}) \leq J_\gamma(\hat{f_n})$

$$\leq J_\gamma(f_\gamma^*) + \cdots$$

$$J_\gamma(f_\gamma^*) = \min_{f \in \mathcal{H}_K} \left\{ L(f) + \gamma \|f\|_K^2 \right\}$$

$$=: L^*(\mathcal{H}_K) + A(\gamma)$$

where $L^*(\mathcal{H}_K) = \inf_{f \in \mathcal{H}_K} L(f)$, so w.p. $\geq 1 - \delta$:

$L(\hat{f_n}) - \inf_{f \in \mathcal{H}_K} L(f)$

$$\leq A(\gamma) + \dfrac{16 M^2 \left(1 + \frac{c_k}{\sqrt{\gamma}}\right)^2}{\sqrt{n}}$$

$$+ M^2\left(1 + \frac{c_k^2}{\gamma}\right)\sqrt{\dfrac{8\log(1/\delta)}{n}}$$

$\hat{f_n} = \hat{f}_{n,\gamma}$

$= \underset{f \in \mathcal{H}_K}{\text{argmin}} \; J_{n,\gamma}(f)$