# The Learning Problem

**Input:** $x \in X \longrightarrow$ features

**Output:** $y \in Y$

**Target function:** $f : X \longrightarrow Y$

**Data:** $(x_1, y_1), \ldots, (x_n, y_n),$ $\underbrace{(x_{n+1}, y_{n+1}), \ldots, (x_m, y_m)}$

$\underbrace{\phantom{(x_1, y_1), \ldots, (x_n, y_n)}}_{\text{Training Data}} \qquad \qquad \text{Testing Data}$

**Loss fn:** $L(x, y, f(x))$, $\quad L : X \times Y \times \mathbb{R} \longrightarrow \mathbb{R}^{\geq 0}$

**Risk:**
(True)
$$R_{L,P}(f) = \int_{X \times Y} L(x, y, f(x)) \; \underbrace{dP(x,y)}_{\substack{\text{Generating pdf of input} \\ \text{\& output} \\ = p(y,x) \, dy \, dx}}$$

$$= \mathbb{E}\left[ L(X, Y, f(X)) \right]$$

$$= \text{expected loss}$$

Here, $\mathbb{P}(x, y) = \underbrace{\mathbb{P}(Y|X)}_{\substack{\text{probabilistic} \\ \text{model of i/p \& o/p}}} \underbrace{\mathbb{P}(X)}_{\substack{\text{probabilistic model} \\ \text{of data generation}}}$

**Bayes risk:**
$$R^{*}_{L,P}(f) = \inf_{f : X \to \mathbb{R}} \int_{X \times Y} L(x, y, f(x)) \, dP(x,y)$$

$$= \inf_{f} \left[ R_{L,P}(f) \right]$$

**Goal** Infer $f_D$ using dataset $D$ whose risk $R_{L,P}(f_D)$
is closest to $R^{*}_{L,P}(f)$

## Consistent Learning

If $f_D$ is the inferred model from data $D$, The learning is said to be universally consistent if

$$R_{L,P}(f_D) \xrightarrow{n} R_{L,P}^* \quad \text{as } n \to \infty \quad \& \ \forall \ P(X,Y)$$

$\longrightarrow$ Stone's theorem (1977)

## No Free Lunch

$\forall$ consistent learning $\land \ \forall$ convergence rate $a_n$,

$\exists \ P(X,Y)$ s.t. convergence rate of this learning method is slower than $a_n$.

$\longrightarrow$ Stochastic bound

Prove $\exists \ M \ \forall \ \epsilon > 0$ s.t.
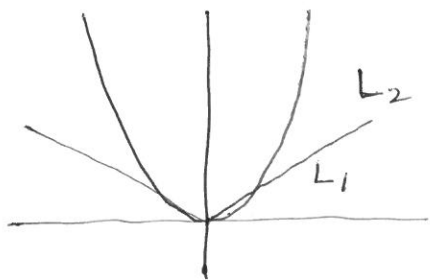
$$P(|X_n| > M) < \epsilon \quad \forall n$$

## Regression

$x \in X \subset \mathbb{R}^d$

$y \in Y = [a, b] \subseteq \mathbb{R}$

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



### Hubber's robust loss

$L(x, y, f(x))$
$$= \begin{cases} \frac{1}{2}(y - f(x))^2 & \text{if } |y - f(x)| < 1 \\ |y - f(x)| - \frac{1}{2} & \text{otherwise} \end{cases}$$

$R_{L,P}(f) = \mathbb{E}\left[(Y - f(x))^2\right]$

## Classification

$x \in X \subset \mathbb{R}^d$

$y \in Y = \{0, 1\} \subset \mathbb{R}$

$L(x, y, f(x)) = \begin{cases} 1 & y \neq f(x) \\ 0 & y = f(x) \end{cases}$

$= 1 - \mathbb{1}(y = f(x))$

$R_{L,P}(f)$

$= 1 - \mathbb{E}\left[\mathbb{1}(y = f(x))\right]$

$= 1 - \mathbb{P}(y = f(x))$

# Empirical risk

$$\widehat{R}_n(f) = \frac{1}{n} \sum_{i=1}^{n} L(Y_i, f(X_i))$$

for simplicity, we write
$$L(x, y, f(x))$$
$$= L(y, f(x))$$

$$= R_n^{train}(f)$$

Goal: $f^* = \underset{f}{\text{argmin}} \; \frac{1}{m-n} \sum_{i=n+1}^{m} L(Y_i, f(X_i))$

$$= \underset{f}{\text{argmin}} \; R_{m-n}^{test}(f)$$

By LLN, $R_n^{train} \xrightarrow{\;n \to \infty\;} R(f)$

$$R_{m-n}^{test} \xrightarrow{\;m-n \to \infty\;} R(f)$$

Thus, we want to minimise $R_n^{train}$ & $R_{m-n}^{test}$

But just minimising $R_n^{train}$ to zero gives extreme overfitting.

eg ①.
$$f(x) = \begin{cases} Y_i & x = X_i \; \forall \, i=1,\dots,n \\ \text{any real value} & \text{otherwise} \end{cases}$$

$\Rightarrow R_n^{train} = 0$, $R(f) \neq 0$, $R_{m-n}^{test}$ very high

$\hookrightarrow$ Overfitting

Soln: Minimise over a less complex hypotheses set & gradually increase the complexity of the set.

If $\mathcal{H}_t$ is the hypothesis set at an iteration $t$,

$$\underbrace{R(f)}_{\substack{\downarrow \\ \text{risk of} \\ \text{classifier}}} - \underbrace{R^*(f^*)}_{\substack{\downarrow \\ \text{Bayesian} \\ \text{risk}}} = \underbrace{R(f) - \inf_{f \in \mathcal{H}_t} R(f)}_{\text{Term 1}} + \underbrace{\inf_{f \in \mathcal{H}_t} R(f) - R(f^*)}_{\text{Term 2}}$$
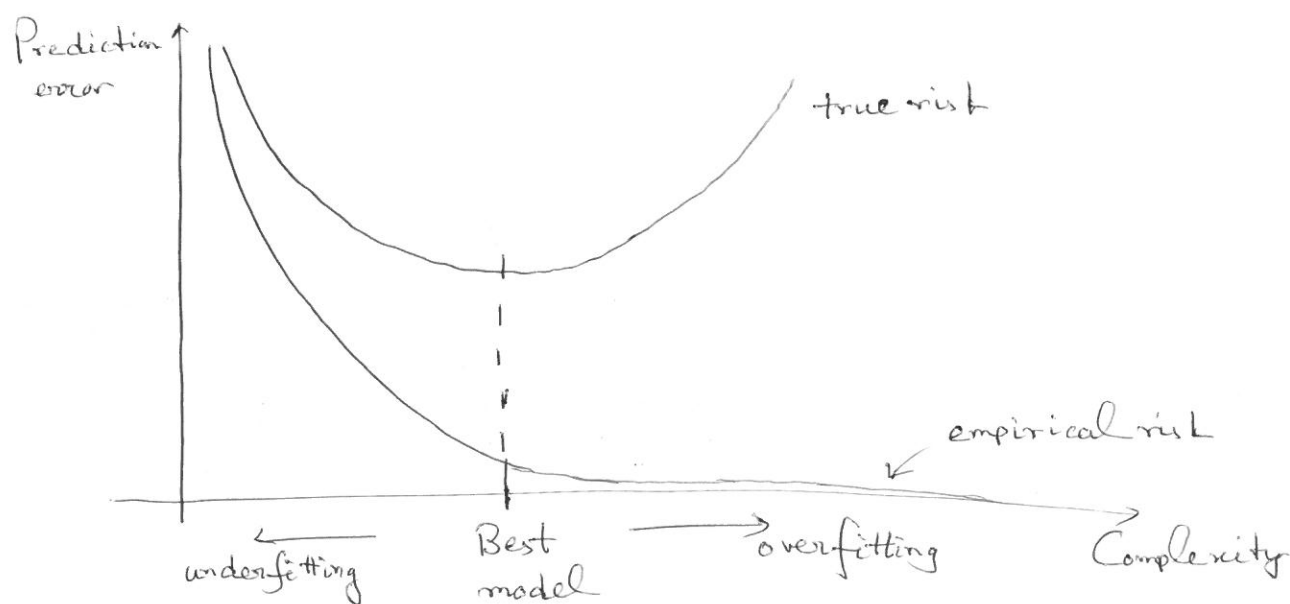
Term 1 = Estimation error

$$= R(f) - \inf_{f \in \mathcal{H}_t} R(f) = R(f) - R^*(f | f \in \mathcal{H}_t)$$

Term 2 = Approximation error

$$= \inf_{f \in \mathcal{H}_t} R(f) - R(f^*) = R^*(f | f \in \mathcal{H}_t) - R(f^*)$$

Term 2 should be non-negative.

Thus, there is this trade-off between complexity of $\mathcal{H}_t$ & minimising the estimation error.



→ Empirical risk alone is not a good performance measure & all loss functions are not eligible for learning.

→ Optimization Problem

$$R_{n,\mathcal{H}}^* = \inf_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} L(Y_i, f(X_i))$$