# Basic Machine Learning Setup

"Supervised"

Given:
- Input Space $X$
- Output Space $Y$
- Unknown Distribution "$\mu$" on $X \times Y$.
  - We are able to sample from $X \times Y$ in some way.

Criteria: We specify:
- a class of functions (parametrized by $\Theta$) $f(\cdot, \theta): X \to Y$, parametrized by $Y^?$
- a loss function $L: Y \times Y \to \mathbb{R}_{\geq 0}$.

Goal:
- Minimize the Risk:
$$\Theta^* = \underset{\Theta}{\text{argmin}} \int_{X \times Y} L(y, f(x, \theta)) \, d\mu.$$

# Example:

$\overset{\text{$L^2$-regression}}{\underset{}{\cdot}}$ Regression:
- Let $X = [0,1]$, $Y = \mathbb{R}$, and $\mu$ given by the pushforward of the under the map $x \to (x, g(x))$ for some unknown function.
- Let $f(\cdot, \theta)$ the same class of function $C$ (say lines or polynomial)
- Let $L(y, y') = |y - y'|^2$
$$\to \underset{f \in C}{\text{argmin}} \int |f(x) - g(x)|^2 \, dx \int = R(f)$$
  $\overset{\text{Hard}}{}$

· Classification:
- Let $X = \mathbb{R}^n$, $Y = \{\pm 1\}$
- Let $f: X \to Y$ be an (unknown) underlying function, i.e. $\mu$ a dist. on $X$.
- Let $f(\cdot, \theta)$, $L$
- $L(y, y') = 1_{y \neq y'}$

$$\rightarrow R(f) = \mathbb{P}(f) = \mathbb{P}_\mu(f(x) \neq \ell(x))$$

- Modified Classification

· $Y' = \mathbb{P}(Y)$ $\quad (Y = \{\pm1\})$

so $f(\cdot, \theta): X \to Y'$

$$\ell(y, y') = \begin{cases} 1 & \text{if } P_{y'}(y) \leq \frac{1}{2} \\ 0 & \text{if } P_{y'}(y) > \frac{1}{2} \end{cases}$$



(graph)

- Cross Entropy for classifier

$$\ell(y, y') = -\log_2\left(P_{y'}(y)\right)$$



Note: Cross-entropy is an upper bound on classification loss.
· We get penalized for our confidence.

Empirical Risk

· We don't know access to $\mu$, but we can draw samples $(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)$

$$R_{emp}^n(\theta) = \frac{1}{n}\sum_{i=1}^n L(y_i, f(x_i, \theta))$$

$$G^TG - \frac{\sqrt{2}}{2}(G^TY)^2 - (1-\frac{\sqrt{2}}{2})(G^TY)^2$$

$$-\frac{\sqrt{2}}{2}(Y^TG)^2 - (1-\frac{\sqrt{2}}{2})(Y^TG)^2 + \frac{\sqrt{2}}{2}Y^TGG^TY$$

$$+ \frac{\sqrt{2}}{2}(GG^TY) \quad (K^TG)^2$$

$$+ \frac{\sqrt{2}}{2}(1-\frac{\sqrt{2}}{2})(G^TY)^2 \quad (G^TY)^2$$

Consider instead the regularized empirical loss:

$$R_{emp}^n(\theta) = \frac{1}{n}\sum_{i=1}^n L(y_i, f(x_i, \theta)) \ (+\lambda S(\theta))$$

$$= 2 + \frac{1}{2} + 1 - \sqrt{2} + \frac{1}{2}\sqrt{2} + 1 \stackrel{?}{=} 4 - 2\sqrt{2}$$

and consider $\theta^*_{emp,n} = \underset{\theta}{\text{argmin}} \ R^n_{emp}(\theta)$

Central question:
- Does
$$R^n_{emp}(\theta^*_{emp,n}) \xrightarrow{?} R(\theta^*)$$

- Note that the weak law of large numbers ensures that for a fixed $\theta$,

$$R^n_{emp}(\theta) \xrightarrow{prob} R(\theta)$$

i.e. $\ \mathbb{P}(|R^n_{emp}(\theta) - R(\theta)| \geq \epsilon) \to 0 \quad \forall \epsilon > 0.$

- However, it may be that the convergence isn't uniform. This is the phenomenon of overfitting.

$$V^TX + X^TV$$
$$V' = A^{-1/2}W$$
$$V'^T = A^{-1/2}W'$$

$$W' = A^{-1/2}V' - \frac{1}{2}(V'^TA^{-1/2}X + X^TA^{-1/2}V')$$

$$A^{1/2}W' = V' - \frac{1}{2}X(V'^TA^{-1/2}X + X^TA^{-1/2}V')$$

$$A^{1/2} \qquad X$$

$$\left(I - \frac{1}{2}XX^T\right) = A$$

$$WA^{-1/2} + XA^{-1/2}W = 0$$

$$A^{-1} = (I + XX^T)$$

## Overfitting Example:

- Consider the first example of last class:
$X = [0,1]$, $Y = \mathbb{R}$, $m$ on $X \times Y$ given
by sampling $X$ uniformly and setting
$y = g(x)$ for some unknown function $g$.

- Let $f(\cdot, \theta)$ be the set of
polynomials, i.e. $\Theta = \{ (a_i)_{i=0}^{\infty} \in \mathbb{R}^{\mathbb{N}}$ s.t. $\exists N$, s.t. $a_i = 0, i \geq N \}$
$$f(x, \theta) = \sum_{i=0}^{\infty} a_i x^i := P_\theta(x)$$
$Y' = \mathbb{R}$
$L(y, y') = |y - y'|^2$

- The risk is
$$R(\theta) = \int |P_\theta(x) - g(x)|^2 \, dx.$$

The empirical risk is
$$R^n_{emp}(\theta) = \frac{1}{n} \sum_{i=1}^{n} \left( P_\theta(x_i) - g(x_i) \right)^2.$$

Let $\theta^*_n \in \text{argmin } R_{emp}$ be the
degree $n-1$ interpolate at the points
$x_1, \ldots, x_n$, then
$$R^n_{emp}(\theta^*_n) = 0.$$

Does $R(\theta^*_n) \to 0$? No! Not in general.

Problem: For each fixed $\theta$,
$R^n_{emp}(\theta) \to R(\theta)$, but the minimizers
$\theta^*_n$ of $R^n_{emp}$ is always and that this
convergence is especially slow! Need uniform
convergence.

Note:

- If $g$ is smooth enough: No overfitting
- If the points $x_i$ are sampled differently (see Chebyshev points / CLT ⟶ ⟶) then no overfitting
→ Overfitting depends both on the model and the distribution $p$ on $X \times Y$.