"Supervised"

# Basic Machine Learning Setup

Given, i.e exist indep. of what we do.
{
- Input Space $\mathcal{X}$
- Output Space $Y$
- Unknown Distribution $\mu$ on $\mathcal{X} \times Y$.
  - We are able to sample from $\mathcal{X} \times Y$ in some way.

Contain our modelling ass. and objects.
{
We provide:
- A class of functions (parametrized by $\Theta$), $f(\cdot, \Theta) : \mathcal{X} \to Y'$.
- A "loss function" $L : Y \times Y' \to \mathbb{R}_{\geq 0}$.

Goal:
- Minimize the Risk:

$$\Theta^* = \arg\min_{\Theta} \int_{\mathcal{X} \times Y} L(y, f(x, \Theta)) \, d\mu.$$

Examples:

$L^2$-regression
· Regression:
- Let $\mathcal{X} = [0, 1]$, $Y = \mathbb{R}$, and $\mu$ given by the pushforward of the $\dots$ under the map $x \to (x, g(x))$ for some unknown function.
- Let $f(\cdot, \Theta)$ be some class of function $C$ (say linear or polynomial)
- Let $L(y, y') = |y - y'|^2$

$$\to \arg\min_{f \in C} \int_0^1 |f(x) - g(x)|^2 \, dx \Big\} = R(f)$$

Hard
· ~~Cross Entropy~~ Classification:
- Let $\mathcal{X} = \mathbb{R}^n$, $Y = \{\pm 1\}$
- Let $l : \mathcal{X} \to Y$ be an (unknown) labelling function, and $\mu$ a dist. on $\mathcal{X}$.
- Let $f(\cdot, \Theta)$ be
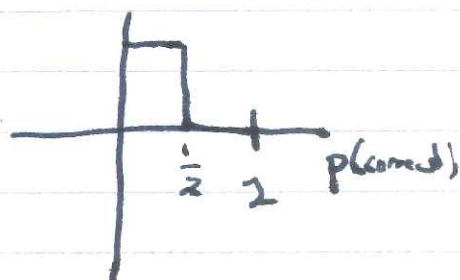- $L(y, y') = \mathbb{1}_{y \neq y'}$

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

- Modified Classification
  :
  
  $\cdot \, Y' = \mathbb{P}\big(Y = \{\pm 1\}\big)$
  
  so $f(\cdot, \Theta): X \rightarrow Y'$
  
  $\cdot L(y, y') = \begin{cases} 1 & \text{if } p_{y'}(y) \leq \frac{1}{2} \\ 0 & \geq \frac{1}{2} \end{cases}$
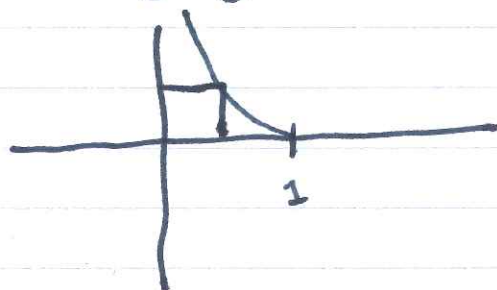


- Cross Entropy for classification
  "
  -"
  "
  
  $\cdot L(y, y') = -\log_2\big(p_{y'}(y)\big)$



Note: Cross-entropy is an upper bound on class. loss.

$\cdot$ We get penalized for our (lack of) additionally confidence.

:

.

## Empirical Risk

$\cdot$ We don't have access to $\mu$, but we can draw samples $(x_1, y_1), \ldots, (x_n, y_n)$

$$G^T G - \frac{\sqrt{2}}{2}\left(G^T Y\right)^2 - \left(1 - \frac{\sqrt{2}}{2}\right)\left(G^T Y\right)^2$$
$$- \frac{\sqrt{2}}{2}\left(Y^T G\right)^2 - \left(1 - \frac{\sqrt{2}}{2}\right)\left(Y^T G\right)^2 + \frac{1}{2} Y^T G G^T Y$$

Consider instead the (regularized) empirical loss: $\quad + \left(1 - \frac{\sqrt{2}}{2}\right)^2 \left(G G^T Y + \frac{\sqrt{2}}{2}\left(1 - \frac{\sqrt{2}}{2}\right)\right.$

$$R_{emp}^n(\theta) = \frac{1}{n} \sum_{i=1}^{n} L\left(y_i, f(x_i, \theta)\right) \left(+ \lambda S(\theta)\right) \quad \left[\begin{array}{l}\left(Y^T G\right)^2 + \\ \left(G^T Y\right)^2\end{array}\right]$$

$$\cancel{2} \quad 2 + \frac{1}{2} + 1 - \sqrt{2} + \frac{1}{2} \overset{?}{=} \sqrt{2} + 1 \overset{?}{=} 4 - 2\sqrt{2}$$

and consider $\theta_{emp,n}^* = \underset{\theta}{\arg\min}\ R_{emp}^n(\theta)$

Central question:
- Does
$$R_{emp}^n\left(\theta_{emp,n}^*\right) \overset{?}{\longrightarrow} R(\theta^*)$$

- Note that the weak law of large numbers means that for a fixed $\theta$,
$$R_{emp}^n(\theta) \xrightarrow{prob} R(\theta)$$

i.e.
$$\mathbb{P}\left(\left|R_{emp}^n(\theta) - R(\theta)\right| > \varepsilon\right) \to 0 \quad \forall \varepsilon > 0.$$

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

$$V^T X + X^T V \qquad \left(I - \frac{1}{2} X X^+\right) = A$$
$$V = A^{-1/2} W$$
$$V' = A^{1/2} W' \qquad W \cancel{A^{-1/2} X} + \cancel{X A^{-1/2}} W = 0$$

$$W' = A^{-1/2} V' - \frac{1}{2}\left(V'^T A^{1/2} X + X^T A^{1/2} V'\right)$$

$$\cancel{X}$$

$$A^{1/2} W' = V' - \frac{1}{2} X \left(V'^T A^{1/2} X + X^T A^{1/2} V'\right)$$
$$\underset{A^{1/2}}{}$$

$$A^{-1} = \left(I + X X^T\right)$$

## Overfitting Example:

- Consider the first example of last class:

$X = [0, 1]$, $Y = \mathbb{R}$, $\mu$ 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. $\Theta = \{(a_i)_{i=0}^{\infty} \in \mathbb{R}^{\mathbb{N}} \text{ s.t. } \exists N, \text{ s.t. } a_i = 0, i \geq N\}$

$Y' = \mathbb{R}$

$$f(x, \Theta) = \sum_{i=0}^{\infty} a_i x^i := p_\Theta(x)$$

$\cdot L(y, y') = |y - y'|^2$

- The risk is

$$R(\Theta) = \int |p_\Theta(x) - g(x)|^2 \, dx.$$

The empirical risk is

$$\hat{R}_{emp}^n(\Theta) = \frac{1}{n} \sum_{i=1}^{n} (p_\Theta(x_i) - g(x_i))^2.$$

Let $\Theta_n^* \in \arg\min \hat{R}_{emp}^n$ be the

degree $n-1$ interpolation at the points $x_1, ..., x_n$, then

$$\hat{R}_{emp}^n(\Theta_n^*) = 0.$$

Does $R(\Theta_n^*) \to 0$? No! Not ~~suggested~~ for general $g$.

Problem: For each fixed $\Theta$,

$\hat{R}_{emp}^n(\Theta) \longrightarrow R(\Theta)$, but the minimizer

$\Theta_n^*$ of $\hat{R}_{emp}^n$ is always such 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 Chebychev points / Chebyshev measure):
 Also no overfitting
$\rightarrow$ Overfitting depends both on the model and the distribution $\mu$ on $X \times Y$.