# Mathematical techniques in data science

PAC Learning

## Supervised learning: standard setting

- Given: a sequence of label data $(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)$ sampled (independently and identically) from an unknown distribution $P_{X,Y}$

- a learning algorithm seeks a function $h : \mathcal{X} \to \mathcal{Y}$, where $\mathcal{X}$ is the input space and $\mathcal{Y}$ is the output space

- The function $h$ is an element of some space of possible functions $\mathcal{H}$, usually called the *hypothesis space*
- In order to measure how well a function fits the training data, a *loss function*

$$L : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}^{\geq 0}$$

  is defined

## Risk and empirical risk

- With a pre-defined loss function, the "optimal hypothesis" is the minimizer over $\mathcal{H}$ of the risk function

$$R(h) = E_{(X,Y) \sim P}[L(Y, h(X))]$$

- Since $P$ is unknown, the simplest approach is to approximate the risk function by the empirical risk
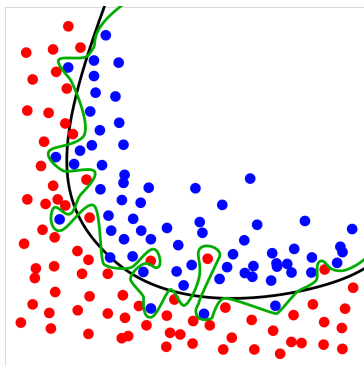
$$R_n(h) = \frac{1}{n} \sum_{i=1}^{n} L(y_i, h(x_i))$$

- The empirical risk minimizer (ERM): minimizer of the empirical risk function (in this lecture, denoted by $\hat{h}_n$)
- Let $h^*$ denotes a minimizer of the risk function

# Failure of ERM

We hope that

$$R(\hat{h}_n) \approx R(h^*),$$

but in general, this might not be true if the hypothesis space $\mathcal{H}$ is too large

## Failure of ERM

- We hope that

$$R(\hat{h}_n) \approx R(h^*),$$

  but in general, this might not be true if the hypothesis space $\mathcal{H}$ is too large

- Question: What does "too large" mean?

- We need to be able to quantify/control the difference between $R(\hat{h}_n)$ and $R(h^*)$

## Modes of estimations

- Analysis

$$\lim_{n \to \infty} x_n = x$$

- Numerical analysis

$$\|x_n - x\| = \mathcal{O}\left(\frac{1}{\sqrt{n}}\right) \quad \text{or} \quad \|x_n - x\| \leq \frac{C}{\sqrt{n}}$$

- PAC (Probably Approximately Correct) learning

$$\|x_n - x\| \leq C(\delta)\frac{1}{\sqrt{n}}$$

with probability at least $1 - \delta$

# PAC learning

## Definition

The probably approximately correct (PAC) learning model typically states as follows: we say that $\hat{h}_n$ is $\epsilon$-accurate with probability $1 - \delta$ if

$$P\left[R(\hat{h}_n) - R(h^*) > \epsilon\right] < \delta.$$

In other words, we have $R(\hat{h}_n) - R(h^*) \leq \epsilon$ with probability at least $(1 - \delta)$.

Probability inequalities

# Markov inequality

### Theorem (Markov inequality)

*For any nonnegative random variable $X$ and $\epsilon > 0$,*

$$P[X \geq \epsilon] \leq \frac{\mathbb{E}[X]}{\epsilon}.$$

# Markov inequality

### Theorem

*For any random variable $X$, $\epsilon > 0$ and $t > 0$*

$$P[X \geq \epsilon] \leq \frac{\mathbb{E}[e^{tX}]}{e^{t\epsilon}}.$$

# Exponential moment of bounded random variables

## Theorem

*If random variable X has mean zero and is bounded in $[a, b]$, then for any $s > 0$,*

$$\mathbb{E}[e^{tX}] \leq \exp\left(\frac{t^2(b-a)^2}{8}\right)$$

### Theorem (Hoeffding's inequality)

Let $X_1, X_2, \ldots, X_n$ be i.i.d copy of a random variable $X \in [a, b]$, and $\epsilon > 0$,

$$P\left[\frac{X_1 + X_2 + \ldots + X_n}{n} - E[X] \geq \epsilon\right] \leq \exp\left(-\frac{n\epsilon^2}{2(b-a)^2}\right).$$

Corollary:

$$P\left[\left|\frac{X_1 + X_2 + \ldots + X_n}{n} - E[X]\right| \geq \epsilon\right] \leq 2\exp\left(-\frac{n\epsilon^2}{2(b-a)^2}\right).$$

Generalization bound for finite hypothesis space and bounded loss

## Assumption

- the loss function $L$ is bounded, that is

$$0 \leq L(y, y') \leq c \quad \forall y, y' \in \mathcal{Y}$$

- the hypothesis space is a finite set, that is

$$\mathcal{H} = \{h_1, h_2, \ldots, h_m\}.$$

## Key ideas

- For any $h \in \mathcal{H}$ and $\epsilon > 0$ we have

$$P[|R_n(h) - R(h)| \geq \epsilon] \leq 2 \exp\left(-\frac{n\epsilon^2}{2c^2}\right).$$

- Using a union bound on the "failure probability" associated with each hypothesis, we have

$$P[\exists h \in \mathcal{H} : |R_n(h) - R(h)| \geq \epsilon] \leq 2|\mathcal{H}| \exp\left(-\frac{n\epsilon^2}{2c^2}\right).$$

## Key ideas

- Using a union bound on the "failure probability" associated with each hypothesis, we have

$$P[\forall h \in \mathcal{H} : |R_n(h) - R(h)| < \epsilon]$$

$$\geq 1 - 2|\mathcal{H}| \exp\left(-\frac{n\epsilon^2}{2c^2}\right).$$

- Under this "good event":

$$R(\hat{h}_n) - R(h^*)$$
$$= [R(\hat{h}_n) - R_n(\hat{h}_n)] + [R_n(\hat{h}_n) - R_n(h^*)] + [R_n(h^*) - R(h^*)]$$
$$\leq 2\epsilon$$

- Conclusion: $\hat{h}_n$ is $(2\epsilon)$-accurate with probability $1 - \delta$, where

$$\delta = 2|\mathcal{H}| \exp\left(-\frac{n\epsilon^2}{2c^2}\right)$$

# PAC estimate for ERM

### Theorem

For any $\delta > 0$ and $\epsilon > 0$, if

$$n \geq \frac{8c^2}{\epsilon^2} \log\left(\frac{2|\mathcal{H}|}{\delta}\right)$$

then $\hat{h}_n$ is $\epsilon$-accurate with probability at least $1 - \delta$.

## PAC estimate for ERM

$$n = \frac{8c^2}{\epsilon^2} \log \left( \frac{2|\mathcal{H}|}{\delta} \right)$$

- Fix a level of confidence $\delta$, the accuracy $\epsilon$ of the ERM is

$$\mathcal{O} \left( \frac{1}{\sqrt{n}} \sqrt{\log \left( \frac{1}{\delta} \right) + \log(|\mathcal{H}|)} \right)$$

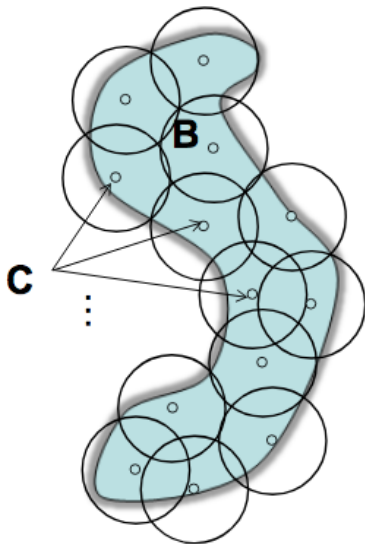- If we want $\epsilon \to 0$ as $n \to \infty$:

$$\log(|\mathcal{H}|) \ll n$$

- The convergence rate will not be better than $\mathcal{O}(n^{-1/2})$

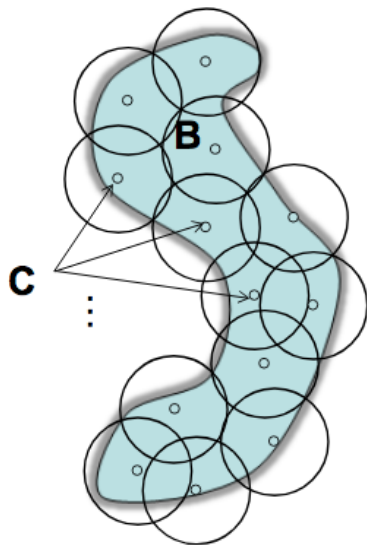Generalization bound using covering number.

# Covering numbers

- Assumption: $\mathcal{H}$ is a metric space with distance $d$ defined on it.
- For $\epsilon > 0$, we denote by $\mathcal{N}(\epsilon, \mathcal{H}, d)$ the *covering number* of $(\mathcal{H}, d)$; that is, $\mathcal{N}(\epsilon, \mathcal{H}, d)$ is the minimal number of balls of radius $\epsilon$ needed to cover $\mathcal{H}$.

## Covering numbers

Remark: If $\mathcal{H}$ is a bounded $k-$dimensional manifold/algebraic surface, then we now that

$$\mathcal{N}(\epsilon, \mathcal{H}, d) = \mathcal{O}\left(\epsilon^{-k}\right)$$

# Generalization bound using covering number.

- Assumption: $\mathcal{H}$ is a metric space with distance $d$ defined on it.
- For $\epsilon > 0$, we denote by $\mathcal{N}(\epsilon, \mathcal{H}, d)$ the *covering number* of $(\mathcal{H}, d)$; that is, $\mathcal{N}(\epsilon, \mathcal{H}, d)$ is the minimal number of balls of radius $\epsilon$ needed to cover $\mathcal{H}$.
- Assumption: loss function $L$ satisfies:

$$|L(h(x), y) - L(h'(x), y)| \leq Cd(h, h') \quad \forall, x \in \mathcal{X}; y \in \mathcal{Y}; h, h' \in \mathcal{H}$$

## Key ideas

- If
$$n = \frac{8c^2}{\epsilon^2} \log\left(\frac{2|\mathcal{H}_\epsilon|}{\delta}\right)$$

  then the event
$$|R_n(h) - R(h)| \leq \epsilon, \forall h \in \mathcal{H}_\epsilon$$

  happens with probability at least $1 - \delta$.

- Under this event, consider any $h \in \mathcal{H}$, then there exists $h_0 \in \mathcal{H}_\epsilon$ such that $d(h, h_0) \leq \epsilon$.

# Key ideas

- Since the loss function is Lipschitz

$$|R_n(h) - R_n(h_0)| \leq Cd(h, h_0)$$

and

$$|R(h) - R(h_0)| \leq Cd(h, h_0).$$

- Conclusion:

$$|R_n(h) - R(h)| \leq (2C + 1)\epsilon \quad \forall h \in \mathcal{H}.$$

# Generalization bound using covering number.

### Theorem

*For all $\epsilon > 0$, $\delta > 0$, if*

$$n \geq \frac{c^2}{2\epsilon^2} \log \left( \frac{2\mathcal{N}(\epsilon, \mathcal{H}, d)}{\delta} \right)$$

*then*

$$|R_n(h) - R(h)| \leq (2C + 1)\epsilon \quad \forall h \in \mathcal{H}.$$

*with probability at least $1 - \delta$.*

## Example: Polynomial covering number.

- Assume that

$$\mathcal{N}(\epsilon, \mathcal{H}, d) \le K\epsilon^{-k}$$

  for some $K > 0$ and $k \ge 1$.

- $\hat{h}_n$ is $\epsilon$-accurate with probability at least $1 - \delta$ if

$$n = \frac{c^2(4C + 2)^2}{2\epsilon^2} \left( \log\left(\frac{2K}{\delta}\right) + k \log\left(\frac{4C + 2}{\epsilon}\right) \right)$$

- Homework: Fix $n$ and $\delta$, derive an upper bound for $\epsilon$.

## Remarks

If we want $\epsilon \to 0$ as $n \to \infty$:

$$dimension(\mathcal{H}) \ll n$$

How do we get that?

- Model selection
- Feature selection
- Regularization:
    - Work for the case $dimension(\mathcal{H}) \gg n$
    - Stabilize an estimator $\to$ force it to live in a neighborhood of a lower-dimensional surface
    - Requires a stability bound instead of a uniform generalization bound