

Ref: [2014] [S. Shalev, et al] Understanding Machine Learning.
Chap 2 A Gentle Start.

2.1 A Formal Model.

Domain set: X , also called instance space.
pts in X : instances.

Label set: Y

Training set: $S = ((x_1, y_1), \dots, (x_m, y_m))$, a finite seq. of pairs in $X \times Y$.

Each (x_i, y_i) is called a training example.

S is also called a training set (although not really a set since we might have repetitions).

The Learner's Input: (Domain / Label / Training: Learner's Inputs).

A function $h: X \rightarrow Y$, called a predictor, a hypothesis, or a classifier.

Notation: $A(S)$, where A is a learning algorithm.

A "simple" data generating model: (to be modified later)

\mathcal{D} : a probability distribution on X .

\exists "correct" labelling function $f: X \rightarrow Y$.
(to be relaxed later).

x_i are drawn from X according to \mathcal{D} .

$y_i = f(x_i)$.

Measure of success:

$L_{\mathcal{D}, f}(h) := \mathbb{P}_{x \sim \mathcal{D}} [h(x) \neq f(x)] = \mathcal{D}(\{x \mid f(x) \neq h(x)\})$,

called the error of the classifier h .

Remark: The generalization error, the risk, or the true error of h are synonyms of $L_{\mathcal{D}, f}(h)$.

L refers to "loss" of the learner.

Note: The learner is blind to \mathcal{D} .

2.2 Empirical Risk Minimization.

Notation: $h_S: X \rightarrow Y$: the output.

Def
 $L_S(h) := \frac{|\{i \in [m] \mid h(x_i) \neq y_i\}|}{m}$, called the

training error, empirical error, or empirical risk.

Empirical Risk Minimization (ERM) is the ML learning paradigm searching for h that minimizes $L_S(h)$.

Remark:

ERM may fail if one is not careful; overfitting may arise.

2.3 Empirical Risk Minimization w/ Inductive Bias.

A common sol'n to overfitting is to apply ERM to a restricted search space.

Def

\mathcal{H} := a set of certain functions $X \rightarrow Y$, called a hypothesis class.

$ERM_{\mathcal{H}}$ is the algorithm applying ERM only on \mathcal{H} .

Thus $ERM_{\mathcal{H}}(S) \in \argmin_{h \in \mathcal{H}} L_S(h)$.

Remark:

- ① The restriction of applying ERM on certain \mathcal{H} is called an inductive bias.
- ② The choice of which \mathcal{H} to look at is based on some prior knowledge about the problem.
- ③ Fundamental question: Over which \mathcal{H} will $ERM_{\mathcal{H}}$ not result in overfitting?

2.3.1 Finite Hypothesis class.

Goal: Prove that $ERM_{\mathcal{H}}$ will not overfit if $|\mathcal{H}| < \infty$ and $|S|$ is sufficiently large.

Notation: $h_S := ERM_{\mathcal{H}}(S)$

Def

The assumption that " $\exists h^* \in \mathcal{H}$ s.t. $L_{\mathcal{D}, f}(h^*) = 0$ " is called the realizability assumption.

Remark: This assumption is to be relaxed later.

The i.i.d. assumption: S comes from sampling X iid w.r.t. \mathcal{D} . This assumption guarantees S to be representative of \mathcal{D} in probability sense.

Notation: $S \sim \mathcal{D}^m$, where $m = |S|$.

Notation:

- ① δ : probability of "getting a nonrepresentative sample" $(1-\delta)$ is called the confidence parameter.
- ② ϵ : accuracy parameter, measuring the quality of prediction.

The event " $L_{\mathcal{D},f}(h) > \epsilon$ " is interpreted as a failure of the learner.

To prove the main thm, we begin w/ a simple lemma.

Lemma 2.2 (Union Bound)

$$\mathcal{D}(A \cup B) \leq \mathcal{D}(A) + \mathcal{D}(B).$$

Cor 2.3

\mathcal{H} : a finite hypothesis class.

$\delta \in (0,1), \epsilon > 0$.

m : an integer w/ $m \geq \frac{\log(|\mathcal{H}|/\delta)}{\epsilon}$.

$f: \mathcal{X} \rightarrow \mathcal{Y}$: a labeling function.

Assume S is a training set from sampling \mathcal{X} iid w.r.t. a distribution \mathcal{D} and labelling by f of size m , and assume also the realizability assumption.

Then, for every ERM hypothesis h_S ,

$$L_{(\mathcal{D},f)}(h_S) \leq \epsilon \text{ w/ probability at least } 1-\delta.$$

<Pf>

Let $A = \{S = (x_1, \dots, x_m) \mid L_{\mathcal{D},f}(h_S) > \epsilon\}$.

Our goal is to figure out condition on m s.t.

$$\mathcal{D}^m(A) \leq \delta.$$

Let $\mathcal{H}_B := \{h \in \mathcal{H} \mid L_{(\mathcal{D},f)}(h) > \epsilon\}$ and

$M := \{S = (x_1, \dots, x_m) \mid \exists h \in \mathcal{H}_B \text{ s.t. } L_S(h) = 0\}$.

(B refers to "Bad" and M refers to "misleading")

For $S \in A$, $L_{\mathcal{D},f}(h_S) > \epsilon \Rightarrow h_S \in \mathcal{H}_B$.

By realizability assumption, $L_{\mathcal{D},f}(h_S) = 0, \forall S$.

Thus $S \in M$ and $A \subseteq M$. $(*)_1$

Note that M can be expressed as

$$M = \bigcup_{h \in \mathcal{H}_B} \{S = (x_1, \dots, x_m) \mid L_S(h) = 0\}. \quad (*)_2$$

By iid assumption,

$$\begin{aligned} \mathcal{D}^m(\{S = (x_1, \dots, x_m) \mid L_S(h) = 0\}) \\ = \prod_{i=1}^m \mathcal{D}(\{x \mid h(x) = f(x)\}) = \prod_{i=1}^m (1 - L_{\mathcal{D},f}(h)) \quad (*)_3 \end{aligned}$$

In particular, for $h \in \mathcal{H}_B$, $(*)_3 < (1-\epsilon)^m$. $(*)_4$

$$\mathcal{D}^m(A) \leq \mathcal{D}^m(M) \leq \sum_{h \in \mathcal{H}_B} \mathcal{D}^m(\{S = (x_1, \dots, x_m) \mid L_S(h) = 0\}) \quad (*)_1$$

$$\leq \sum_{h \in \mathcal{H}_B} (1-\epsilon)^m = |\mathcal{H}_B| \cdot (1-\epsilon)^m \leq |\mathcal{H}| \cdot (1-\epsilon)^m.$$

Note that $1-\epsilon \leq e^{-\epsilon}, \forall \epsilon > 0$.

$$\text{Thus } \mathcal{D}^m(A) \leq |\mathcal{H}| \cdot e^{-m\epsilon}.$$

To make $\mathcal{D}^m(A) \leq \delta$, we may make

$$|\mathcal{H}| \cdot e^{-m\epsilon} \leq \delta. \text{ i.e. } |\mathcal{H}|/\delta \leq e^{m\epsilon}$$

$$\text{i.e. } m \geq \frac{\log(|\mathcal{H}|/\delta)}{\epsilon}, \text{ as assumed. } \quad (*)_5$$

Chap 3 A Formal Learning Model.

3.1 PAC Learning.

PAC means Probably Approximately Correct.

Def

A hypothesis class \mathcal{H} is PAC learnable if

$\exists m_{\mathcal{H}}: (0,1)^2 \rightarrow \mathbb{N}$ and a learning algorithm s.t.

$\forall \epsilon, \delta \in (0,1)$ and every dist. \mathcal{D} over \mathcal{X} , and

\forall labelling function $f: \mathcal{X} \rightarrow \{0,1\}$,

if the realizability assumption holds, then when

running the algorithm on $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ iid

examples generated by \mathcal{D} and labelled by f ,

the returned hypothesis h satisfies

$$L_{(\mathcal{D},f)}(h) \leq \epsilon \text{ w/ prob. at least } 1-\delta. \text{ (over all possible examples)}$$

Remark:

① The accuracy parameter $\epsilon \leftrightarrow$ "approximately correct".

② The confidence parameter $\delta \leftrightarrow$ "probably".

Def

\mathcal{H} : hypothesis class that is PAC learnable.

The sample complexity of learning \mathcal{H} is the minimal

$m_{\mathcal{H}}$ satisfying conditions in the def of PAC learnable.

Cor 3.2

Every finite hypothesis class \mathcal{H} is PAC learnable

$$\text{w/ sample complexity } m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{\log(|\mathcal{H}|/\delta)}{\epsilon} \right\rceil.$$

3.2 A More General Learning Model.

3.2.1 Releasing the Realizability Assumption

- Agnostic PAC Learning.

Motivation:

Possible that we have $(x, y), (x, y') \in S$ w/ $y \neq y'$.

Sol'n:

Replacing target labelling function w/ a "data-labels generating distribution".

From now on, $\mathcal{D} :=$ a prob. dist. over $\mathcal{X} \times \mathcal{Y}$.

i.e. the joint distribution on $\mathcal{X} \times \mathcal{Y}$.

Rmk:

$$\mathcal{D}(x, y) = \underbrace{\mathcal{D}((x, y) | x)}_{\text{(conditional)}} \cdot \underbrace{\mathcal{D}_x}_{\text{(marginal)}}$$

Def (redefine)

$$L_{\mathcal{D}}(h) := \mathbb{P}_{(x, y) \sim \mathcal{D}} [h(x) \neq y] = \mathcal{D}(\{(x, y) | h(x) \neq y\}).$$

Rmk:

$$L_S(h) := \frac{|\{i \in [m] | h(x_i) \neq y_i\}|}{m}, \text{ as before.}$$

Def

\mathcal{D} : distribution on $\mathcal{X} \times \{0, 1\}$.

Define

$$f_{\mathcal{D}}(x) = \begin{cases} 1, & \text{if } P[y=1|x] \geq 1/2 \\ 0, & \text{o.w.} \end{cases}$$

$f_{\mathcal{D}}$ is called the Bayes optimal predictor

Prop

For every classifier g , $L_{\mathcal{D}}(f_{\mathcal{D}}) \leq L_{\mathcal{D}}(g)$.

Thus, $f_{\mathcal{D}}$ is optimal. case 1 $P[y=1|x] \geq 1/2$.
case 2 $P[y=1|x] < 1/2$.

Rmk:

$\therefore \mathcal{D}$ is unknown \therefore We can not compute $f_{\mathcal{D}}$.
Compare $P[f_{\mathcal{D}}(x) \neq y]$ and $P[g(x) \neq y]$.

Def 3.3

A hypothesis class \mathcal{H} is agnostic PAC learnable

if $\exists m_{\mathcal{H}}: (0, 1)^2 \rightarrow \mathbb{N}$ and a learning algorithm

s.t. $\forall \epsilon, \delta \in (0, 1)$, \mathcal{D} : dist. over $\mathcal{X} \times \mathcal{Y}$, when

running the algorithm on $m, m_{\mathcal{H}}(\epsilon, \delta)$ iid examples

generated by \mathcal{D} , the returned hypothesis h

satisfies

$$L_{\mathcal{D}}(h) \leq \min_{h' \in \mathcal{H}} L_{\mathcal{D}}(h') + \epsilon, \text{ w/ prob. at least } 1 - \delta.$$

[End of W1]

(over all m training examples).

3.2.2 The Scope of Learning Problems Modeled

For different learning tasks, we need different model, ex: multiclass classification, regression.

Rmk:

- For multiclass classification, the measurement of quality of a hypothesis is as before.
- For regression, we may use the expected square difference.

$$L_{\mathcal{D}}(h) = \mathbb{E}_{(x, y) \sim \mathcal{D}} (h(x) - y)^2, \text{ instead.}$$

Motivation:

We need to generalize the formalism of measure of success to accommodate a wide range of learning tasks.

Def

A function $\ell: \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}_+$ is called a loss function

where $\mathbb{R}_+ := \{\text{nonnegative real numbers}\}$.

Given a distribution \mathcal{D} on \mathcal{Z} , the risk function

$$\text{is } L_{\mathcal{D}}(h) := \mathbb{E}_{z \sim \mathcal{D}} [\ell(h, z)].$$

For $z_1, \dots, z_m \in \mathcal{Z}$ drawn iid wrt \mathcal{D} from \mathcal{Z} , the empirical risk is, where $S = (z_1, \dots, z_m)$,

$$L_S(h) := \frac{1}{m} \sum_{i=1}^m \ell(h, z_i).$$

Example:

① (0-1 loss)

$$\mathcal{Z} := \mathcal{X} \times \mathcal{Y}.$$

$$\ell_{0-1}(h, (x, y)) := \begin{cases} 0, & \text{if } h(x) = y \\ 1, & \text{if } h(x) \neq y. \end{cases} \quad (= \chi_{\{h(x) \neq y\}})$$

The $L_{\mathcal{D}}(h)$ defined here (in this subsection) is the same as before.

② (square loss).

$$\mathcal{Z} := \mathcal{X} \times \mathcal{Y}$$

$$\ell_{sq}(h, (x, y)) := (h(x) - y)^2.$$

Then $L_{\mathcal{D}}(h)$ is the expected square difference.

Rmk:

As demonstrated, the newly defined $L_{\mathcal{D}}(h)$ is a generalization.

Def 3.4

A hypothesis class \mathcal{H} is agnostic learnable w.r.t. a set Z and a loss function $\ell: \mathcal{H} \times Z \rightarrow \mathbb{R}_+$ if $\exists m_{\mathcal{H}}: (0,1)^2 \rightarrow \mathbb{N}$ and a learning algorithm s.t.

$\forall \epsilon, \delta \in (0,1)$, dist. \mathcal{D} over Z , when running the algorithm on $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ iid examples generated by \mathcal{D} , the returned $h \in \mathcal{H}$ satisfies

$$L_{\mathcal{D}}(h) \leq \min_{h' \in \mathcal{H}} L_{\mathcal{D}}(h') + \epsilon \text{ w/ prob. at least } 1-\delta$$

where

$$L_{\mathcal{D}}(h) := \mathbb{E}_{z \sim \mathcal{D}} [\ell(h, z)]$$

(least $1-\delta$ (over the choice of m training examples)).

Rmk:

① Strictly speaking, we need to require, for each fixed $h \in \mathcal{H}$, $\ell(h, z)$ is \mathcal{D} -measurable as a function of z . (Of course, a σ -algebra on Z needs to be preassigned.)

② The authors talk about "representation indep. learning" (or "improper learning") while the learning defined previously is called "proper learning".

Chap 4 Learning via Uniform Convergence.

4.1 Uniform Convergence is SUFF. for Learnability.

Def 4.1

A training set S is called ϵ -representative (w.r.t. domain Z , hypothesis class \mathcal{H} , loss function ℓ and dist. \mathcal{D}) if

$$|L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon, \forall h \in \mathcal{H}.$$

Lemma 4.2

$S: \epsilon/2$ -representative.

$h_S \in \text{ERM}_{\mathcal{H}}(S)$.

Then

$$L_{\mathcal{D}}(h_S) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon.$$

<pf>.

$$L_{\mathcal{D}}(h_S) \leq L_S(h_S) + \epsilon/2 = \min_{h \in \mathcal{H}} L_S(h) + \epsilon/2$$

$$\leq \min_{h \in \mathcal{H}} (L_{\mathcal{D}}(h) + \epsilon/2) + \epsilon/2$$

$$= \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon. \quad \#$$

Rmk:

By 4.2, to guarantee \mathcal{H} to be PAC learnable, it suffices to guarantee, w/ prob. at least $1-\delta$, S is $\epsilon/2$ -representative.

Def 4.3 (Unif. Con.)

A hypothesis class \mathcal{H} is said to have the uniform convergence property (w.r.t. domain Z , loss function ℓ) if

$\exists m_{\mathcal{H}}^{\text{UC}}: (0,1)^2 \rightarrow \mathbb{N}$ s.t. $\forall \epsilon, \delta \in (0,1)$, dist. \mathcal{D} over Z , if S is a sample of $m \geq m_{\mathcal{H}}^{\text{UC}}(\epsilon, \delta)$ examples iid from \mathcal{D} , then S is ϵ -representative w/ prob. $\geq 1-\delta$.

Rmk:

The adj. "uniform" refers to the unif. of $m_{\mathcal{H}}^{\text{UC}}$ over all members of \mathcal{H} and all dist. \mathcal{D} over Z .

Cor 4.4

\mathcal{H} : a hypothesis class w/ UC property.

Then \mathcal{H} is agnostic PAC learnable w/ sample complexity $m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{\text{UC}}(\epsilon/2, \delta)$.

Moreover, $\text{ERM}_{\mathcal{H}}$ is a successful agnostic PAC learner for \mathcal{H} .

4.2 Finite Classes are Agnostic PAC learnable.

We are going to prove the statement as the circle suggests. The following lemma is used:

Lemma 4.5 (Hoeffding's inequality).

$\theta_1, \dots, \theta_m$: iid rv.

$\mathbb{E}[\theta_i] = \mu$, and $P[a \leq \theta_i \leq b] = 1, \forall i$.

Then, $\forall \epsilon > 0$,

$$P\left[\left|\frac{1}{m} \sum_{i=1}^m \theta_i - \mu\right| > \epsilon\right] \leq 2e^{-2m\epsilon^2/(b-a)^2}.$$

Cor 4.6

\mathcal{H} : finite hypothesis class.

\mathcal{Z} : domain.

$\ell: \mathcal{H} \times \mathcal{Z} \rightarrow [a, b]$: a loss function.

Then \mathcal{H} has the unif. con. property w/

$$m_{\mathcal{H}}^{uc}(\epsilon, \delta) \leq \left\lceil \frac{\log(2|\mathcal{H}|/\delta)}{2\epsilon^2/(b-a)^2} \right\rceil.$$

Therefore (by Cor 4.4),

\mathcal{H} is agnostic PAC learnable using ERM w/

sample complexity

$$m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{uc}(\epsilon/2, \delta) \leq \left\lceil \frac{2 \log(2|\mathcal{H}|/\delta)}{\epsilon^2/(b-a)^2} \right\rceil.$$

<Pf>

Two steps:

Step 1: Applying the union bound.

Let

$$B = \{S: \exists h \in \mathcal{H} \text{ s.t. } |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\}.$$

$\in (\mathcal{X} \times \mathcal{Y})^m \text{ w/ } \mathcal{D}^m$

Then

$$B = \bigcup_{h \in \mathcal{H}} \{S: |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\}.$$

$$\text{Thus, } \mathcal{D}^m(B) \leq \sum_{h \in \mathcal{H}} \mathcal{D}^m(\{S: |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\}).$$

Step 2: Applying a measure concentrat'n ineq.

(Here, we use Hoeffding's).

Denote $S = (z_1, \dots, z_m)$. Fix $h \in \mathcal{H}$.

$$\theta_i := \ell(h, z_i), \mu := \mathbb{E}_{z \sim \mathcal{D}}[\ell(h, z)] = \mathbb{E}[\ell(h, z_i)], \forall i.$$

$$\text{Note that } L_{\mathcal{D}}(h) = \mu \text{ and } L_S(h) = \frac{1}{m} \sum_{i=1}^m \theta_i.$$

$\therefore z_1, \dots, z_m$ are iid

$\therefore \theta_1, \dots, \theta_m$ are iid.

Thus, by Hoeffding's ineq,

$$\mathcal{D}^m(\{S: |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\}) = \mathbb{P}\left[\left|\frac{1}{m} \sum_{i=1}^m \theta_i - \mu\right| > \epsilon\right] \leq 2e^{-2m\epsilon^2/(b-a)^2}.$$

$$\Rightarrow \mathcal{D}^m(B) \leq 2 \cdot |\mathcal{H}| \cdot e^{-2m\epsilon^2/(b-a)^2}.$$

To make $\mathcal{D}^m(B) < \delta$, it suffices to make

$$2e^{-2m\epsilon^2/(b-a)^2} < \delta. \text{ i.e. } m > \frac{\log(2|\mathcal{H}|/\delta)}{2\epsilon^2/(b-a)^2}. \quad (\#)$$

Remark: ("Discretization Trick").

Though the assumption $|\mathcal{H}| < \infty$ seems unrealistic, due to limitation of finite expression of numbers in computers, it becomes realistic.

Remark:

① Hypothesis classes w/ unif. con. property are also called Glivenko-Cantelli classes.

② (See Chap 6) In binary classification problems, $UC \Leftrightarrow$ learnable.

(Not the case for general learning problems).

Chap 5 The Bias Complexity Tradeoff.

End of W2

Motivation:

Is there a universal learner?

5.1 The No-Free-Lunch Theorem.

Thm 5.1 (No-Free-Lunch)

A : an algorithm for binary classification w.r.t. 0-1 loss and over domain \mathcal{X} .

m : an integer w/ $m < |\mathcal{X}|/2$.

Then \exists dist. \mathcal{D} over $\mathcal{X} \times \{0,1\}$ s.t.

1. $\exists f: \mathcal{X} \rightarrow \{0,1\}$ w/ $L_{\mathcal{D}}(f) = 0$.

2. $L_{\mathcal{D}}(A(S)) \geq 1/8$ w/ prob. $\geq 1/\eta$ over the choice of $S \sim \mathcal{D}^m$.

Remark:

Intuitively, the thm states that

for every learner (the algorithm A), there exists a task on which it fails, while there is another learner that succeeds.

5.1.1 No-Free-Lunch and Prior Knowledge.

Remark:

A preassigned hypothesis class \mathcal{H} represents a prior knowledge. Thus, by taking $\mathcal{H} = \{\text{all functions from } \mathcal{X} \text{ to } \{0,1\}\}$, we are in a state of lack of prior knowledge.

Cor 5.2

\mathcal{X} : an infinite domain set.

$\mathcal{H} := \{\text{all functions } \mathcal{X} \rightarrow \{0,1\}\}$.

Then \mathcal{H} is not PAC learnable.

<Pf>
Suppose \mathcal{H} is PAC learnable.

Then \exists algorithm A and $m_{\mathcal{H}} : (0,1)^2 \rightarrow \mathbb{N}$ w/
the "PAC learnable properties".

Take $\epsilon < 1/8$ and $\delta < 1/4$.

Let \mathcal{D} be a distribution on $\mathcal{X} \times \{0,1\}$ corresponding
to A as in the No-Free-Lunch Thm.

PAC learnability implies that

$L_{\mathcal{D}}(A(S)) < \epsilon < 1/8$ w/ prob. $> 1 - \delta > 3/4$ over the
i.e. choice of $S \sim \mathcal{D}^m$.

$$\mathcal{D}^m(\{S \mid L_{\mathcal{D}}(A(S)) < 1/8\}) > 3/4.$$

$$\Leftrightarrow 1 - \mathcal{D}^m(\{S \mid L_{\mathcal{D}}(A(S)) \geq 1/8\}) > 3/4.$$

$$\Leftrightarrow \mathcal{D}^m(\{S \mid L_{\mathcal{D}}(A(S)) \geq 1/8\}) < 1/4$$

$$\Leftrightarrow \sim(\mathcal{D}^m(\{S \mid L_{\mathcal{D}}(A(S)) \geq 1/8\}) \geq 1/4)$$

$$\Leftrightarrow \sim(L_{\mathcal{D}}(A(S)) \geq 1/8 \text{ w/ prob. } \geq 1/4 \text{ over choice of } S \sim \mathcal{D}^m)$$

This contradicts w/ 2. of
No-Free-Lunch Thm. (#)

Rmk:

By Cor 5.2, w/o using any prior knowledge will fail
learning. We may impose prior knowledge via restricting
 \mathcal{H} .

5.2 Error Decomposition.

Q: How should we choose a good hypothesis class
?

① \mathcal{H} should be large enough to include a hypothesis
w/ small $\min_{h \in \mathcal{H}} L_{\mathcal{D}}(h)$.

② \mathcal{H} should be small enough to be learnable.
(as suggested by No-Free-Lunch Thm.)

Def \mathcal{H} : a hypothesis class

Given $h_S \in \text{ERM}_{\mathcal{H}}(S)$.

Write $L_{\mathcal{D}}(h_S) = \epsilon_{\text{app}} + \epsilon_{\text{est}}$, where

$$\epsilon_{\text{app}} := \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) \text{ and } \epsilon_{\text{est}} := L_{\mathcal{D}}(h_S) - \epsilon_{\text{app}}.$$

ϵ_{app} is called the approximation error.

ϵ_{est} is called the estimation error.

Rmk:

① ϵ_{app} is the minimum risk achievable by \mathcal{H} .
i.e. it is exactly the inductive bias.

If \mathcal{H} : realizable, then $\epsilon_{\text{app}} = 0$.

In agnostic case, ϵ_{app} may be large.

② ϵ_{est} results because $L_{\mathcal{D}}(h_S)$ is just an
estimation of $\min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) (= \epsilon_{\text{app}})$.

Observation:

	ϵ_{app}	ϵ_{est}	$L_{\mathcal{D}}(h_S) = \epsilon_{\text{app}} + \epsilon_{\text{est}}$
$\mathcal{H} \uparrow$	\downarrow	? might \uparrow	?
$\mathcal{H} \downarrow$	\uparrow	? might \downarrow	?

Def

on \mathcal{H}
The tradeoff we face when trying to minimize
 $L_{\mathcal{D}}(h_S)$ is called the bias-complexity tradeoff.

Rmk:

Too large \mathcal{H} may result in overfitting while
"small \mathcal{H} may result " underfitting.

Chap 6 The VC Dimension.

Q: Which classes \mathcal{H} are PAC learnable?

6.1 Infinite-Size Class Can Be Learnable.

Def

A function of the form $h_a: \mathbb{R} \rightarrow \{0,1\}$ w/
 $h_a(x) = \mathbb{1}_{[x < a]}$ is called a threshold function.

Lemma 6.1

$\mathcal{H} := \{\text{all threshold functions}\}$.

Then \mathcal{H} is PAC learnable, using ERM, w/
sample complexity $m_{\mathcal{H}}(\epsilon, \delta) \leq \lceil \log(2/\delta)/\epsilon \rceil$.

Rmk:

The "PAC learnable" here is the very first
(non-agnostic) learnability w/ specific target
 f .

<Pf> (Intuition).

Let S be a sample wrt a dist. \mathcal{D} over \mathcal{X} .

$b_0 := \max \{x : (x, 1) \in S\}. (\max(\emptyset) := -\infty).$

$b_1 := \min \{x : (x, 0) \in S\}. (\min(\emptyset) := \infty).$

Then $\nexists b$ w/ $(b, 0) \in S$ or $(b, 1) \in S$, and

$h_{b_S} \in \text{ERM}_{\mathcal{H}}(S)$, where b_S is arbitrary in

Assume $a^* \in \mathbb{R}$ w/ $L_{\mathcal{D}}(h_{a^*}) = 0$. (b_0, b_1) .

"Then", as $|S| = m \uparrow$, b_S is closer to a^* w/

Therefore, \mathcal{H} is PAC learnable. high prob.

6.2 The VC-Dimension.

Motivation:

As suggested by L 6.1, while $|\mathcal{H}| < \infty$ is sufficient for learnability, it is NOT necessary.

Def 6.2, 6.3

\mathcal{H} : a hypothesis class of fun. from \mathcal{X} to $\{0, 1\}$.

$C = \{c_1, \dots, c_m\} \subseteq \mathcal{X}$, i.e. C is a finite subset of \mathcal{X} .

① The restriction of \mathcal{H} to C is the set

$$\mathcal{H}_C = \{f|_C : C \rightarrow \{0, 1\} \mid f \in \mathcal{H}\}.$$

② \mathcal{H} is said to shatter C if \mathcal{H}_C is the set of all possible functions from C to $\{0, 1\}$ i.e. $|\mathcal{H}_C| = 2^{|C|}$.

Ex 6.2

$\mathcal{H} := \{\text{all threshold functions}\}.$

① $C = \{c_1\} \subseteq \mathbb{R}.$

Then $h_{c_1-1}(c_1) = 0$ and $h_{c_1+1}(c_1) = 1.$

Thus \mathcal{H} shatters C .

② $C = \{c_1, c_2\}$ w/ $c_1 < c_2.$

Then $\nexists h \in \mathcal{H}$ s.t. $h(c_1) = 0$ and $h(c_2) = 1.$

$\Rightarrow \mathcal{H}$ does not shatter C . Ⓢ

Cor 6.4 (of Proof of No-Free-Lunch Thm).

\mathcal{H} : a hypothesis class of fun. from \mathcal{X} to $\{0, 1\}.$

m : a training set size.

$\exists C \subseteq \mathcal{X}$ of size $2m$ s.t. C is shattered by $\mathcal{H}.$

Then, \forall algorithm A , \exists dist. \mathcal{D} on $\mathcal{X} \times \{0, 1\}$ and $\exists h \in \mathcal{H}$ s.t.

(1) $L_{\mathcal{D}}(h) = 0$, but

(2) $L_{\mathcal{D}}(A(S)) \geq 1/8$ w/ prob. $\geq 1/4$ (over the choice of $S \sim \mathcal{D}^m$).

Roughly speaking, Cor 6.4 tells us that

if \mathcal{H} shatters some set of size $2m$, then we cannot learn \mathcal{H} using m examples.

Def 6.5

The VC-dimension of a hypothesis class \mathcal{H} , denoted $\text{VCdim}(\mathcal{H})$ is the maximal size of a set $C \subseteq \mathcal{X}$ that can be shattered by $\mathcal{H}.$

i.e.

$$\text{VCdim}(\mathcal{H}) = \max \{|C| : C \subseteq \mathcal{X}, \mathcal{H} \text{ shatters } C\}.$$

$$\text{VCdim}(\mathcal{H}) = \infty \text{ if } \forall n \in \mathbb{N}, \exists C \text{ w/ } |C| = n$$

Rmk: s.t. \mathcal{H} shatters C .

VC refers to Vapnik-Chervonenkis.

Thm 6.6

\mathcal{H} : a hypothesis class w/ $\text{VCdim}(\mathcal{H}) = \infty.$

Then \mathcal{H} is not PAC learnable.

<Pf>

For training size m , $\exists C \subseteq \mathcal{X}$ w/ size $2m$ s.t.

\mathcal{H} shatters C .

By Cor 6.4, the result follows. Ⓢ

Rmk:

We shall see later in this Chap that the converse is also true.

6.3 Examples.

6.3.1 Threshold fun.

$$\text{VCdim}(\mathcal{H}) = 1.$$

6.3.2 Intervals

$$\text{VCdim}(\mathcal{H}) = 2.$$

6.3.3 Axis-Aligned Rectangles (on \mathbb{R}^2).

$$\text{VCdim}(\mathcal{H}) = 4.$$

6.3.4 Finite Classes.

\mathcal{H} : a finite hypothesis class.

For $C \subseteq \mathcal{X}$, $|\{ \text{all fun. } C \rightarrow \{0,1\} \}| = 2^{|C|}$.

Thus, if $2^{|C|} > |\mathcal{H}|$, then \mathcal{H} cannot shatter C .

$$\Rightarrow VCdim(\mathcal{H}) \leq \log_2(|\mathcal{H}|) < \infty.$$

Rmk: ^① For finite \mathcal{H} , $VCdim(\mathcal{H})$ may be significantly smaller than $\log_2(|\mathcal{H}|)$.

② $VCdim(\mathcal{H})$ is not necessarily the number of parameters defining \mathcal{H} . (see 6.3.5).

6.4 The Fundamental Theorem of PAC learning.

Thm 6.7 (The Fundamental Thm of Statistical Learning).

\mathcal{H} : a hypothesis class of fun. from \mathcal{X} to $\{0,1\}$,
w/ 0-1 loss function.

Then T.F.A.E.

- (1) \mathcal{H} has the UC property.
- (2) ERM is a successful agnostic PAC learner for \mathcal{H} .
- (3) \mathcal{H} : agnostic PAC learnable.
- (4) \mathcal{H} : PAC learnable.
- (5) ERM is a successful PAC learner for \mathcal{H} .
- (6) $VCdim(\mathcal{H}) < \infty$.

Thm 6.8 (The F.T. of S.L. - Quantitative Version).

\mathcal{H} : hypothesis class of fun. $\mathcal{X} \rightarrow \{0,1\}$. w/ 0-1 loss.

$$VCdim(\mathcal{H}) = d < \infty.$$

Then \exists abs. const. C_1, C_2 s.t.

(1) \mathcal{H} has UC property w/

$$C_1 \frac{d + \log(1/\delta)}{\epsilon^2} \leq m_{\mathcal{H}}^{UC}(\epsilon, \delta) \leq C_2 \frac{d + \log(1/\delta)}{\epsilon^2}.$$

(2) \mathcal{H} : agnostic PAC learnable w/

$$C_1 \frac{d + \log(1/\delta)}{\epsilon^2} \leq m_{\mathcal{H}}(\epsilon, \delta) \leq C_2 \frac{d + \log(1/\delta)}{\epsilon^2}.$$

(3) \mathcal{H} : PAC learnable w/

$$C_1 \frac{d + \log(1/\delta)}{\epsilon} \leq m_{\mathcal{H}}(\epsilon, \delta) \leq C_2 \frac{d + \log(1/\delta)}{\epsilon}.$$

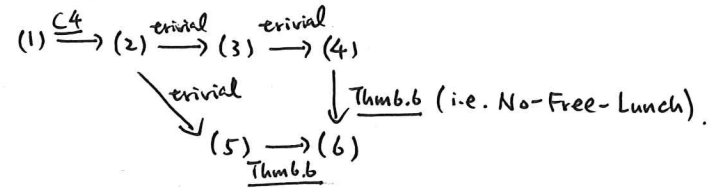
<Pf> in Chap 28.

Rmk:

Similar F.T. holds for "some" other learning problems but "not all".

6.5 Proof of Thm 6.7.

So far, we have



We are going to prove (6) \rightarrow (1).

6.5.1 Sauer's Lemma and Growth Function.

Def 6.9 \mathcal{H} : a hypothesis class.

The growth function of \mathcal{H} , denoted $\mathcal{Z}_{\mathcal{H}}: \mathbb{N} \rightarrow \mathbb{N}$, is

$$\mathcal{Z}_{\mathcal{H}}(m) := \max_{\substack{C \subseteq \mathcal{X} \\ |C|=m}} |\mathcal{H}_C|.$$

Rmk:

If $VCdim(\mathcal{H}) = d$, then $\forall m \leq d$, $\mathcal{Z}_{\mathcal{H}}(m) = 2^m$, exponential in m .

However, when $m > d$, we have $\mathcal{Z}_{\mathcal{H}}(m)$ polynomial in m :

Lemma 6.10 (Sauer - Shelah - Perles).

\mathcal{H} : hypothesis class w/ $VCdim(\mathcal{H}) \leq d < \infty$.

Then, $\forall m$,

$$\mathcal{Z}_{\mathcal{H}}(m) \leq \sum_{i=0}^d \binom{m}{i}.$$

In particular, if $m > d+1$, then

$$\mathcal{Z}_{\mathcal{H}}(m) \leq (e^m/d)^d.$$

End of W3