

SESSION 12: STATISTICAL MACHINE LEARNING (II)



Gang Li

Deakin University, Australia

2021-08-07

PAC Learning	3
The Statistical Learning Framework	4
PAC Learnability	5
The General PAC Learning Model	10
Agnostic PAC Learnability	14
PAC versus Agnostic PAC Learning	15
Agnostic Learning Finite Hypothesis Classes	16
Representative Sample	17
Uniform Convergence	19
Agnostic Learning Finite Hypothesis Classes	20
Quiz	22

Table of Content

PAC Learning

- The Statistical Learning Framework
- PAC Learnability
- The General PAC Learning Model
- Agnostic PAC Learnability
- PAC versus Agnostic PAC Learning

Agnostic Learning Finite Hypothesis Classes

- Representative Sample
- Uniform Convergence
- Agnostic Learning Finite Hypothesis Classes

Quiz

PAC Learning

The Statistical Learning Framework

	<p>The learner’s task is to:</p> <p>Input: training data $S = \{(x_1, y_1), \cdots, (x_m, y_m)\} \in (\mathcal{X} \times \mathcal{Y})^m$</p> <p>Output: prediction rule $h : \mathcal{X} \rightarrow \mathcal{Y}$</p> <p>🌀 Measure The error of a prediction rule $h : \mathcal{X} \rightarrow \mathcal{Y}$ can be defined as:</p> <p>Generalization risk $L_{(\mathcal{D}, f)}(h) \stackrel{def}{=} P_{x \sim \mathcal{D}}[h(x) \neq f(x)] \stackrel{def}{=} \mathcal{D}(\{x : h(x) \neq f(x)\})$</p> <p>Empirical risk $L_S(h) \stackrel{def}{=} \frac{ \{i \in [m] : h(x_i) \neq y_i\} }{m}$</p>
ERM	<p>ERM comes up with a predictor h that minimizes $L_S(h)$</p> <p>$ERM_{\mathcal{Y}^{\mathcal{X}}}(S) \in \operatorname{argmin}_{h \in \mathcal{Y}^{\mathcal{X}}} L_S(h)$</p> <p>ERM with Inductive Bias comes up with any $h \in \mathcal{H}$ that minimizes $L_S(h)$</p> <p>$ERM_{\mathcal{H}}(S) \in \operatorname{argmin}_{h \in \mathcal{H}} L_S(h)$</p>

Can only be *Approximately* correct

👉 For any training data S with m i.i.d. examples, we should not hope find an h s.t. $L_{(\mathcal{D},f)}(h) = 0$

Proof.

- For every $\epsilon \in (0, 1)$ take $\mathcal{X} = \{x_1, x_2\}$ and $\mathcal{D}(\{x_1\}) = 1 - \epsilon$, $\mathcal{D}(\{x_2\}) = \epsilon$
- The probability not to see x_2 at all among m i.i.d. examples in S is $(1 - \epsilon)^m \approx e^{-\epsilon m}$
- So if $\epsilon \ll \frac{1}{m}$ we are likely not to see x_2 at all, but then we can not know its label.

□

Relaxation.

- We would be happy with $L_{(\mathcal{D},f)}(h) < \epsilon$, where ϵ is the user-specified **accuracy parameter**.

□

Can only be *Probably* correct

👉 For any training data S with m i.i.d. examples, no algorithm can guarantee $L_{(\mathcal{D},f)}(h) \leq \epsilon$

Proof.

- Recall that the input to the learner is a set of randomly generated examples, there is always a (very small) chance to see the same example again and again.

□

Relaxation.

- We would allow the algorithm to fail with probability δ , where $\delta \in (0, 1)$ is the user-specified **confidence parameter**
- Here, the probability is over the random choice of examples

□

Probably Approximately Correct (PAC) Learnability

A hypothesis class \mathcal{H} is **PAC learnable** if there exists a function $m_{\mathcal{H}} : (0,1)^2 \rightarrow \mathcal{N}$ and a learning algorithm with the following property:

- For every $\epsilon, \delta \in (0,1)$, for every distribution \mathcal{D} over \mathcal{X} , and for every labelling function $f : \mathcal{X} \rightarrow \{0,1\}$, if the **realizable assumption** holds with respect to \mathcal{H} , \mathcal{D} and f , then when we run the algorithm on $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ i.i.d. examples generated by \mathcal{D} and labelled by f , the algorithm returns a hypothesis h such that, with probability of at least $(1 - \delta)$, $L_{(\mathcal{D},f)}(h) \leq \epsilon$.

Key Points.

- It is a distribution free model, i.e. no particular assumption about \mathcal{D}
- Training and test samples are drawn according to the same \mathcal{D} (otherwise transfer learning)
- It deals with the question of learnability for \mathcal{H} not a particular concept, namely the “target labelling function” f .

□

Steps.

- The learner does not know \mathcal{D} and f
- The learner receives the *accuracy* parameter ϵ and the *confidence* parameter δ
- The learner can ask for training data S containing $m_{\mathcal{H}}(\epsilon, \delta)$ examples
 - ◆ the number of examples can depend on ϵ and δ , but not on depend \mathcal{D} and f
- The learner should output a hypothesis h , s.t. with probability of at least $(1 - \delta)$ it holds that $L_{(\mathcal{D},f)}(h) \leq \epsilon$.
 - ◆ the learner should be **P**robably (with probability at least $(1 - \delta)$) **A**pproximately (up to accuracy ϵ) **C**orrect

□

Sample Complexity

The function $m_{\mathcal{H}} : (0,1)^2 \rightarrow \mathcal{N}$ determines the **sample complexity** of learning \mathcal{H} , namely, $m_{\mathcal{H}}(\epsilon, \delta)$ represents how many examples are required to guarantee a PAC solution:

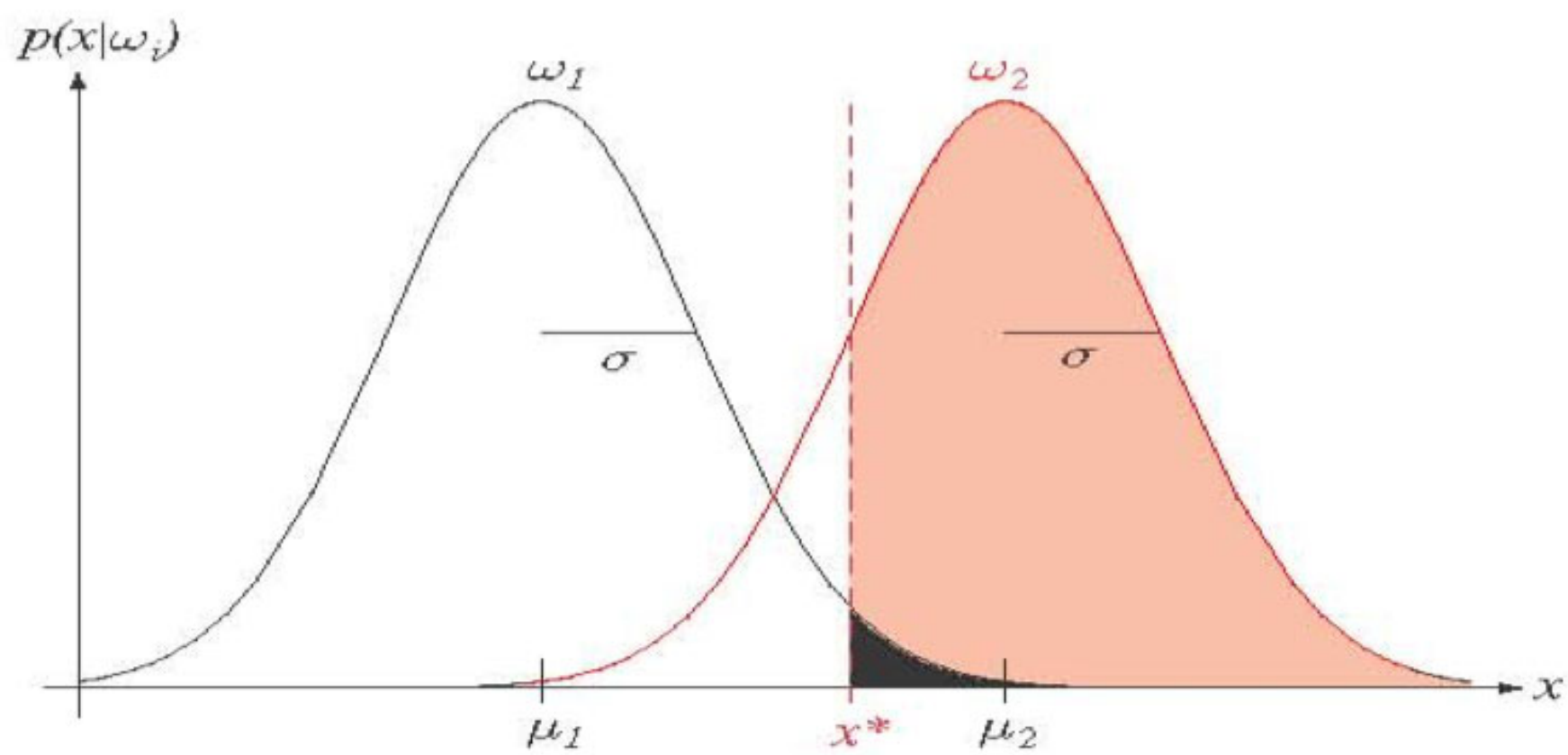
- It is a function of the *accuracy* parameter ϵ and the *confidence* parameter δ
- It also depends on the properties of the hypothesis class \mathcal{H} .
 - ◆ If \mathcal{H} is PAC learnable, there are many functions $m_{\mathcal{H}}$ that satisfy the requirements given in the PAC learnability definition.
 - ◆ We define the sample complexity to be the “minimal function”

Every **finite hypothesis class** \mathcal{H} is PAC learnable with the sample complexity:

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \lceil \frac{\log(|\mathcal{H}|/\delta)}{\epsilon} \rceil = \lceil \frac{1}{\epsilon} [\log(|\mathcal{H}|) + \log(\frac{1}{\delta})] \rceil$$

Is there a learner?

In many scenarios, there is no perfect learner:



(None)-8119669 (2021-08-07) – 9 / 25

General PAC Learning Model

PAC learning model can be generalized in two aspects:

Relaxing the Realizability Assumption

- We assume that labels are generated by some $f \in \mathcal{H}$, this assumption may be too strong.

Learning beyond Binary Classification

- Many learning tasks involve multiple class classification
- or even prediction of a real valued number.

(None)-8119669 (2021-08-07) – 10 / 25

General PAC Learning — Relaxing the Realizability Assumption

Relaxing the Realizability Assumption:

Intuition.

- Relax the realizability assumption by replacing the “target labelling function” f with a more flexible notion, a data-labels generating distribution.
 - ◆ In PAC model, \mathcal{D} is a distribution over \mathcal{X}
 - ◆ In this aspect, \mathcal{D} is a distribution over $Z = \mathcal{X} \times \mathcal{Y}$
- The *Generalization risk* is then defined as:

$$L_{\mathcal{D}}(h) \stackrel{def}{=} P_{Z \sim \mathcal{D}}[h(x) \neq y] \stackrel{def}{=} \mathcal{D}(\{x : h(x) \neq y\})$$

- The notation of “approximately correct” is now defined as:

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

□

The General PAC Learning — Beyond Binary Classification

Scope of Learning Problems.

Muticlass categorization \mathcal{Y} is a finite set representing $|\mathcal{Y}|$ different classes.

- For example, the degree could be $\mathcal{Y} = \{Bachelor, Honours, Masters, PhD\}$

Regression $\mathcal{Y} = \mathcal{R}$

- For example, one wishes to predict the marks of a student based on the resources access pattern.

□

The General PAC Learning — Loss Functions

- Let $Z = \mathcal{X} \times \mathcal{Y}$
- Given hypothesis $h \in \mathcal{H}$, and an example $(x, y) \in Z$, how good is h on (x, y) ?
- **Loss Function:**
$$l : \mathcal{H} \times Z \rightarrow \mathcal{R}_+$$

0-1 loss $l(h, (x, y)) = \begin{cases} 1 & \text{if } h(x) \neq y \\ 0 & \text{if } h(x) = y \end{cases}$
Squared loss $l(h, (x, y)) = (h(x) - y)^2$
Absolute-value loss $l(h, (x, y)) = |h(x) - y|$
Cost-sensitive loss $l(h, (x, y)) = C_{h(x), y}$, where C is $|\mathcal{Y}| \times |\mathcal{Y}|$ matrix.

Agnostic PAC Learnability

- A hypothesis class \mathcal{H} is **agnostic PAC learnable** with respect to a set Z and a loss function $l : \mathcal{H} \times Z \rightarrow \mathcal{R}_+$, if there exists a function $m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathcal{N}$ and a learning algorithm with the following property:
- For every $\epsilon, \delta \in (0, 1)$, for every distribution \mathcal{D} over Z , when running the algorithm on $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ i.i.d. examples generated by \mathcal{D} , the algorithm returns a hypothesis $h \in \mathcal{H}$ such that, with probability of at least $(1 - \delta)$: $\min_{h^* \in \mathcal{H}} L_{\mathcal{D}}(h^*) + \epsilon$

Probably (at least $(1 - \delta)$ probability) Approximately (up to accuracy ϵ) Correct solve:

- $\min_{h \in \mathcal{H}} L_{\mathcal{D}}(h), \text{ where } L_{\mathcal{D}}(h) \stackrel{def}{=} E_{z \sim \mathcal{D}}[l(h, z)]$
- Learner knows \mathcal{H}, Z and l
 - The learner receives the *accuracy* parameter ϵ and the *confidence* parameter δ
 - The learner can decide on training set size m based on ϵ and δ .
 - The learner does not know \mathcal{D} but can sample $S \sim \mathcal{D}^m$
 - Using S the learner outputs some hypothesis $h \in \mathcal{H}$, with probability of at least $(1 - \delta)$ it holds that $L_{\mathcal{D}}(h) \leq \min_{h^* \in \mathcal{H}} L_{\mathcal{D}}(h^*) + \epsilon$.

PAC versus Agnostic PAC Learning


Table 1: Comparison of PAC and Agnostic PAC

	PAC	Agnostic PAC
Distribution	\mathcal{D} over \mathcal{X}	\mathcal{D} over $\mathcal{X} \times \mathcal{Y}$
Truth	$f \in \mathcal{H}$	not in class or does not exist
Risk	$L_{(\mathcal{D},f)}(h) = \mathcal{D}(\{x : h(x) \neq f(x)\})$	$L_{\mathcal{D}}(h) = \mathcal{D}(\{x : h(x) \neq y\})$
Training set	$(x_1, \dots, x_m) \sim \mathcal{D}^m, \forall i, y_i = f(x_i)$	$((x_1, y_1), \dots, (x_m, y_m)) \sim \mathcal{D}^m$
Goal	$L_{(\mathcal{D},f)}(h) \leq \epsilon$	$L_{\mathcal{D}}(h) \leq \min_{h^* \in \mathcal{H}} L_{\mathcal{D}}(h^*) + \epsilon$

\mathcal{X} : Domain \mathcal{Y} : Range \mathcal{H} : Hypothesis Class
 L : Loss function ϵ : accuracy parameter m : sample size

Agnostic Learning Finite Hypothesis Classes


Representative Sample



A training set S is called **ϵ -representative** w.r.t. domain Z , hypothesis class \mathcal{H} , loss function l and distribution \mathcal{D} , if
$$\forall h \in \mathcal{H}, \quad |L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon$$

- Intuition.*
- The hope is that an h that minimizes the empirical risk with respect to the sample S is a risk minimizer, or has risk close to the minimum, with respect to the true data probability distribution \mathcal{D} .
 - This concept ensures that: uniformly over *all hypotheses* in the hypothesis class \mathcal{H} , the empirical risk will be *close to the true* risk.
-

Representative Sample



Assume that a training set S is $\frac{\epsilon}{2}$ -representative w.r.t. domain Z , hypothesis class \mathcal{H} , loss function l and distribution \mathcal{D} , then, any output of $ERM_{\mathcal{H}}(S)$, namely any $h_S \in \operatorname{argmin}_{h \in \mathcal{H}} L_S(h)$

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


Proof.

- $L_{\mathcal{D}}(h_S) \leq L_S(h_S) + \frac{\epsilon}{2}$
- $L_S(h^*) \leq L_{\mathcal{D}}(h^*) + \frac{\epsilon}{2}$
- Combine them together, we have


$$\begin{aligned} L_{\mathcal{D}}(h_S) &\leq L_S(h_S) + \frac{\epsilon}{2} \\ &\leq L_S(h^*) + \frac{\epsilon}{2} \\ &\leq L_{\mathcal{D}}(h^*) + \frac{\epsilon}{2} + \frac{\epsilon}{2} \\ &= L_{\mathcal{D}}(h^*) + \epsilon \end{aligned}$$

□

Uniform Convergence



A hypothesis class \mathcal{H} has the uniform convergence property if there exists a function $m_{\mathcal{H}}^{UC} : (0, 1)^2 \rightarrow \mathcal{N}$, such that for every $\epsilon, \delta \in (0, 1)$, and every distribution \mathcal{D} , we have: if S is a sample with $m \geq m_{\mathcal{H}}^{UC}(\epsilon, \delta)$ examples drawn i.i.d. according to \mathcal{D} , then with probability of at least $(1 - \delta)$, S is ϵ -representative.



If a class \mathcal{H} has the uniform convergence property with the sample complexity $m_{\mathcal{H}}^{UC}$, then \mathcal{H} is agnostically PAC learnable with the sample complexity

$$m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{UC}\left(\frac{\epsilon}{2}, \delta\right)$$

Furthermore, $ERM_{\mathcal{H}}$ paradigm is a successful agnostic PAC learner for \mathcal{H} .

■ $m_{\mathcal{H}}^{UC}$ measures the minimal sample complexity of obtaining the uniform convergence.

Agnostic Learning Finite Hypothesis Classes

Assume \mathcal{H} is **finite** and the range of the loss function is $[0, 1]$, then \mathcal{H} is **agnostic PAC learnable** using the $ERM_{\mathcal{H}}$ algorithm with sample complexity:

👍

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \lceil \frac{2 \log(2|\mathcal{H}|/\delta)}{\epsilon^2} \rceil = \lceil \frac{2}{\epsilon^2} [\log(2|\mathcal{H}|) + \log(\frac{1}{\delta})] \rceil$$

Proof. It suffices to show that \mathcal{H} has the **uniform convergence property** with

$$m_{\mathcal{H}}^{UC}(\epsilon, \delta) \leq \lceil \frac{\log(2|\mathcal{H}|/\delta)}{2\epsilon^2} \rceil$$

1. To show uniform convergence, we need: $\mathcal{D}^m(\{S : \exists h \in \mathcal{H}, |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\}) < \delta$

2. From the union bound, we have:

$$\begin{aligned} &\mathcal{D}^m(\{S : \exists h \in \mathcal{H}, |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\}) \\ &= \mathcal{D}^m(\bigcup_{h \in \mathcal{H}} \{S : |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\}) \\ &\leq \sum_{h \in \mathcal{H}} \mathcal{D}^m(\{S : |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\}) \end{aligned}$$

□

Proof.

3. $L_{\mathcal{D}}(h) = E_{z \sim \mathcal{D}}[l(h, z)]$ and $L_S(h) = \frac{1}{m} \sum_{i=1}^m l(h, z_i)$, let $\theta_i = l(h, z_i)$

4. For all i , $E[\theta_i] = L_{\mathcal{D}}(h)$

5. From *Hoeffding's inequality*:

$$\mathcal{D}^m(\{S : |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\}) \leq 2e^{-2m\epsilon^2}$$

6. We have:

$$\mathcal{D}^m(\{S : \exists h \in \mathcal{H}, |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\}) \leq \sum_{h \in \mathcal{H}} \mathcal{D}^m(\{S : |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\}) \leq 2|\mathcal{H}|e^{-2m\epsilon^2}$$

7. So if $m \geq \frac{\log(2|\mathcal{H}|/\delta)}{2\epsilon^2}$, we have the right hand side is at most δ as required.

□

(None)-8119669 (2021-08-07) – 20 / 25

The Discretization Trick

■ Suppose \mathcal{H} is parametrized by d numbers.

■ Suppose we are happy with a representation of each number using b bits

■ Then $|\mathcal{H}| \leq 2^{db}$, and so

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \lceil \frac{2db + 2 \log(2/\delta)}{\epsilon^2} \rceil$$

■ While not very elegant, it is a great tool for upper bounding *sample complexity*.

(None)-8119669 (2021-08-07) – 21 / 25

10

PAC Learning

?

Theoretical analysis:
1. If the range of the loss function is $[a, b]$, then the sample complexity satisfies:
$$m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{UC}(\epsilon/2, \delta) \leq \lceil \frac{2 \log 2 |\mathcal{H}| \delta (b - a)^2}{\epsilon^2} \rceil.$$

2. Given any probability distribution \mathcal{D} over $\mathcal{X} \times \{0, 1\}$, the *Bayes Optimal Predictor* is defined as: $f_{\mathcal{D}}(x) = \begin{cases} 1 & \text{if } P[y = 1|x] \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$ Show that for every probability distribution \mathcal{D} , the *Bayes Optimal Predictor* $f_{\mathcal{D}}$ is optimal, in the sense that for every classifier g from \mathcal{X} to $\{0, 1\}$, we have $L_{\mathcal{D}}(f_{\mathcal{D}}) \leq L_{\mathcal{D}}(g)$.

(None)-8119669 (2021-08-07) – 23 / 25




Questions?

(None)-8119669 (2021-08-07) – 24 / 25

Contact Information

Associate Professor **GANG LI**
School of Information Technology
Deakin University
Geelong, Victoria 3216, Australia



-  GANGLI@TULIP.ORG.AU
-  [OPEN RESOURCES OF TULIP-LAB](#)
-  [TEAM FOR UNIVERSAL LEARNING AND INTELLIGENT PROCESSING](#)