PAC Learnability

Spring 2025

Outline

Empirical Risk Minimization

- The learner's input:
 - Domain set (Instances Space): An arbitrary set X.
 - ▶ Domain point (Instance) : $x \in \mathcal{X}$.
 - Label set: $\mathcal{Y} = \{0, 1\}$ or $\mathcal{Y} = \{+1, -1\}$.
 - ▶ Training set: $S = \{(x_i, y_i)\}_{i=1}^m$, where every $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$.
- The learner's output: $h : \mathcal{X} \to \mathcal{Y}$.
- A simple data-generation model: we assume that each pair in the training set S is generated by
 - first sampling a point x_i according to a fixed but unknown distribution D on X,
 - and then labeling it by the "correct" labeling function f, that is, $y_i = f(x_i)$.

Generalization error: a measure of success.

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

• Training error:
$$L_S(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}(h(x_i) \neq y_i)$$

- Hypothesis class H: A set of functions mapping from X to Y.
- The ERM_H Learner: for a given class H, and a training set S, the ERM_H learner uses the ERM rule to choose a predictor h ∈ H, with the lowest possible error over S. Formally,

$$\operatorname{ERM}_{\mathcal{H}}(S) \in \operatorname{argmin}_{h \in \mathcal{H}} L_{S}(h).$$

We also use h_S to denote a result of applying ERM_H to *S*, that is,

$$h_{\mathcal{S}} \in \underset{h \in \mathcal{H}}{\operatorname{argmin}} L_{\mathcal{S}}(h).$$

Definition (The Realizability Assumption) There exists $h^* \in \mathcal{H}$ s.t. $L_{(\mathcal{D},f)}(h^*) = 0$.

This assumption implies that with probability 1, we have

$$\blacktriangleright L_{\mathcal{S}}(h^{\star}) = 0.$$

• $L_S(h_S) = 0$ for every ERM hypothesis h_S .

Theorem

Let \mathcal{H} be a finite hypothesis class. Let $\delta \in (0, 1)$ and $\epsilon > 0$ and let *m* be an integer that satisfies

$$m \geq rac{\log(|\mathcal{H}|/\delta)}{\epsilon}.$$

Then, for any labeling function, f, and for any distribution D, for which the realizability assumption holds, with probability at least $1 - \delta$ over the choice of an i.i.d. sample S of size m, we have that for every ERM hypothesis, h_S , it holds that

$$L_{(\mathcal{D},f)}(h_{\mathcal{S}}) \leq \epsilon.$$

Notes: for a sufficiently large *m*, the ERM_H rule over a finite hypothesis class will be Probably (with confidence 1 – δ) Approximately (up to an error of ε) Correct.

Proof. Let \mathcal{H}_B be the set of "bad" hypotheses, that is,

$$\mathcal{H}_{\mathcal{B}} = \{h \in \mathcal{H} : L_{(\mathcal{D},f)}(h) > \epsilon\}.$$

Let $S|_x = \{x_1, \dots, x_m\}$ be the instances of the training set. Then we upper bound the probability

$$\mathcal{D}^m(\{S|_{x}: L_{(\mathcal{D},f)}(h_S) > \epsilon\}).$$

In addition, let $M = \{S|_x : \exists h \in \mathcal{H}_B, L_S(h) = 0\}$. Note that

$$\{S|_{x}: L_{(\mathcal{D},f)}(h_{S}) > \epsilon\} \subseteq M = \bigcup_{h \in \mathcal{H}_{B}} \{S|_{x}: L_{S}(h) = 0\}.$$

Hence

$$\mathcal{D}^{m}(\{S|_{x}: L_{(\mathcal{D},f)}(h_{S}) > \epsilon\}) \leq \mathcal{D}^{m}(M) = \mathcal{D}^{m}(\bigcup_{h \in \mathcal{H}_{B}} \{S|_{x}: L_{S}(h) = 0\})$$
$$\leq \sum_{h \in \mathcal{H}_{B}} \mathcal{D}^{m}(\{S|_{x}: L_{S}(h) = 0\})$$

Since the instances are sampled i.i.d., we get that

$$\mathcal{D}^{m}(\{S|_{x}: L_{S}(h) = 0\}) = \prod_{i=1}^{m} \mathcal{D}(\{x_{i}: h(x_{i}) = f(x_{i})\}).$$

Note for every $h \in \mathcal{H}_B$,

$$\mathcal{D}(\{x_i : h(x_i) = f(x_i)\}) = 1 - L_{(\mathcal{D},f)}(h) \le 1 - \epsilon, \text{ and}$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへぐ

$$\mathcal{D}^m(\{S|_x : L_S(h) = 0\}) \le (1 - \epsilon)^m \le e^{-m\epsilon}.$$

Therefore,

$$\mathcal{D}^m(\{S|_x: L_{(\mathcal{D},f)}(h_S) > \epsilon\}) \leq |\mathcal{H}_B|e^{-m\epsilon} \leq |\mathcal{H}|e^{-m\epsilon}.$$

$$|\mathcal{H}|\boldsymbol{e}^{-\boldsymbol{m}\boldsymbol{\epsilon}} \leq \delta,$$

then

Let

$$m \geq rac{\log(|\mathcal{H}|/\delta)}{\epsilon},$$

and

$$1 - \mathcal{D}^m(\{S|_x : L_{(\mathcal{D},f)}(h_S) > \epsilon\}) \ge 1 - \delta. \square$$



Empirical Risk Minimization

Probably Approximately Correct Learning



Definition (PAC Learnability)

A hypothesis class \mathcal{H} is PAC learnable if there exist a function $m_{\mathcal{H}} : (0,1)^2 \to \mathbb{N}$ and a learning algorithm with the following property: For every $\delta, \epsilon \in (0,1)$, for every distribution over \mathcal{X} , and for every labeling function $f : \mathcal{X} \to \{0,1\}$, if the realizable assumption holds with respect to $\mathcal{H}, \mathcal{D}, f$, then when running the algorithm on $m \ge m_{\mathcal{H}}(\epsilon, \delta)$ i.i.d. examples genenated by \mathcal{D} and labeled by f, the algorithms returns a hypothesis h such that, with probability of at least $1 - \delta$ (over the choice of the examples),

 $L_{(\mathcal{D},f)}(h) \leq \epsilon.$

Probably Approximately Correct Learnability

- Approximately Correct: the accuracy parameter e determines how far the output classifier can be from the optimal one.
- Probably: the confidence parameter δ indicates how likely the classifier is to meet that accuracy requirement.

(ロ) (同) (三) (三) (三) (○) (○)

- Sample complexity: How many samples are required to guarantee a probably approximately correct solution.
 - If H is PAC learnable, there are many functions m_H that satisfy the requirements given the definition of PAC learnability.
 - The sample complexity of learning H is defined as minimal function, in the sense that for any ε, δ, m_H(ε, δ) is the minimal integer that satisfies the requirements g of PAC learning with accuracy ε and confidence δ.

(ロ) (同) (三) (三) (三) (○) (○)

Corollary

Every finite hypothesis class is PAC learnable with sample complexity

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

Q: Does the finiteness determine the PAC learnability of a hypothesis class?

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

A: No.



Empirical Risk Minimization

Probably Approximately Correct Learning

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへぐ

Agnostic PAC learnability

To waive the realizability assumption

- ► Recall that the realizability assumption requires that there exists h^{*} ∈ H s.t. L_{D,f}(h^{*}) = 0.
- For practical learning tasks, the realizability assumption may be too strong.
- From PAC learning to Agnostic PAC learning: releasing the realizability assumption.

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

A More Realistic Model for the Data-Generating Distribution

- From the deterministic case of a fixed but unknown distribution over X and a correct labeling function f to the stochastic case.
- Let \mathcal{D} be a probability distribution over $\mathcal{X} \times \mathcal{Y}$.
- Two parts of such a distribution:
 - a marginal distribution \mathcal{D}_x over unlabelled domain points.
 - ► a conditional probability D((x, y)|x) over labels for each point.

Generalization Error Revised:

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

・ロト ・ 同 ・ ・ ヨ ・ ・ ヨ ・ うへつ

- ► The Goal: to find some hypothesis, h : X → Y, that (probably approximately) minimizes the generalization error, L_D(h).
- The Bayes Optimal Predictor: Given any distribution D over X × {0,1}, the best label predicting function from X to {0,1} will be

$$f_{\mathcal{D}}(x) = \begin{cases} 1 & \text{if } \mathbb{P}[y=1|x] \ge \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

• It is easy to verify that for every distribution \mathcal{D} ,

$$L_{\mathcal{D}}(f_{\mathcal{D}}) \leq L_{\mathcal{D}}(g)$$

for every classifier $g: \mathcal{X} \to \{0, 1\}.$

- $\blacktriangleright \mathcal{D}$ is a fixed but unknown distribution.
- We cannot utilize the optimal predictor $f_{\mathcal{D}}$.
- Instead, we require that the learning algorithm will find a predictor whose error is not much larger than the best possible error of a predictor in some given benchmark hypothesis class.

(ロ) (同) (三) (三) (三) (○) (○)

Definition (Agnostic PAC Learnability)

A hypothesis class \mathcal{H} is agnostic PAC learnable if there exist a function $m_{\mathcal{H}} : (0,1)^2 \to \mathbb{N}$ and a learning algorithm with the following property: For every $\delta, \epsilon \in (0,1)$, for every distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$, then when running the algorithm on $m \ge m_{\mathcal{H}}(\epsilon, \delta)$ i.i.d. examples genenated by \mathcal{D} , the algorithms returns a hypothesis *h* such that, with probability of at least $1 - \delta$ (over the choice of the *m* training examples),

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

・ロト ・ 同 ・ ・ ヨ ・ ・ ヨ ・ うへつ

- Agnostic PAC learning generalizes the definition of PAC learning.
 - If the realizability assumption holds, agnostic PAC learning provides the same guarantee as PAC learning.
- When the realizability assumption does not hold, no learner can guarantee an arbitrarily small error.
- Under the definition agnostic PAC learning, a learner can still declare success if its error is not much larger than the best error achievable by a predictor from the hypothesis class H.

(ロ) (同) (三) (三) (三) (○) (○)

- Generalized loss functions :
 - Given any set H and some domain Z, let ℓ be any function from H × Z to the set of nonnegative real numbers, ℓ : H × Z → ℝ₊.
 - We call such functions loss functions.
 - For prediction tasks, $Z = \mathcal{X} \times \mathcal{Y}$.

0-1 loss:

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

・ロト ・ 同 ・ ・ ヨ ・ ・ ヨ ・ うへつ

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

► Risk function: the expected loss of a classifier *h* ∈ *H* with respect to A distribution *D* over the domain set *Z*:

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

► Empirical Risk: the expected loss of a classifier h ∈ H over a given a sample S = (z₁, z₂, · · · , z_m) ∈ Z^m:

$$L_{\mathcal{S}}(h) = \frac{1}{m} \sum_{i=1}^{m} \ell(h, z_i).$$

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Definition (Agnostic PAC Learnability for General Loss Functions)

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

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

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



Empirical Risk Minimization

Probably Approximately Correct Learning

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

Agnostic PAC learnability

Uniform Convergence

Definition (ϵ -representative sample)

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

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

Lemma

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

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

Lemma

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

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

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

Proof. For every $h \in \mathcal{H}$,

$$egin{array}{rcl} L_{\mathcal{D}}(h_{\mathcal{S}}) &\leq & L_{\mathcal{S}}(h_{\mathcal{S}})+rac{\epsilon}{2} \ &\leq & L_{\mathcal{S}}(h)+rac{\epsilon}{2} \ &\leq & L_{\mathcal{D}}(h)+rac{\epsilon}{2}+rac{\epsilon}{2} \ &= & L_{\mathcal{D}}(h)+\epsilon. \end{array}$$

(S is ϵ – representative.) (h_S is an ERM predictor.) (S is ϵ – representative.)

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

Definition (Uniform Convergence)

We say that a hypothesis class \mathcal{H} has the *uniform convergence property*(w.r.t. domain Z and loss function I) if there exists a function $m_{\mathcal{H}}^{UC} : (0,1)^2 \to \mathbb{N}$ such that for every $\epsilon, \delta \in (0,1)$ and for every probability distribution \mathcal{D} over Z, if S is a sample of $m \ge m_{\mathcal{H}}^{UC}(\epsilon, \delta)$ examples drawn i.i.d. according to \mathcal{D} , then, with probability at least $1 - \delta$, S is ϵ -representative.

Corollary

If a class \mathcal{H} has the uniform convergence property with a function $m_{\mathcal{H}}^{UC}$ then the class is agnostically PAC learnable with the sample complexity $m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{UC}(\frac{\epsilon}{2}, \delta)$. Furthermore, in that case, the ERM_{\mathcal{H}} paradigm is a successful agnostic PAC learner for \mathcal{H} .

・ロト ・ 同 ・ ・ ヨ ・ ・ ヨ ・ うへつ

Corollary

Let \mathcal{H} be a finite hypothesis class, let Z be a domain, and let $\ell : \mathcal{H} \times Z \rightarrow [0, 1]$ be a loss function. Then \mathcal{H} enjoys the uniform convergence property with sample complexity

$$m_{\mathcal{H}}^{\mathit{UC}}(\epsilon,\delta) \leq \lceil rac{\log(2|\mathcal{H}|/\delta)}{2\epsilon^2}
ceil.$$

Furthermore, the class is agnostically PAC learnable using the ERM algorithm with sample complexity

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

◆□▶ ◆□▶ ▲□▶ ▲□▶ ■ ののの

Theorem (Hoeffding's Inequality)

Let $\theta_1, \dots, \theta_m$ be a sequence of *i.i.d.* random variables and assume that for all *i*, $\mathbb{E}[\theta_i] = \mu$ and $\mathbb{P}[a \le \theta_i \le b] = 1$. Then, for any $\epsilon > 0$,

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

(ロ) (同) (三) (三) (三) (○) (○)

Proof.

Fix some $\epsilon, \delta \in (0, 1)$. We need to find a sample size *m* that guarantees that for any \mathcal{D} , with probability of at least $1 - \delta$ of the choice of $S = (z_1, \dots, z_m)$ sampled i.i.d. from \mathcal{D} we have that for all $h \in \mathcal{H}$, $|L_S(h) - L_{\mathcal{D}}(h)| \le \epsilon$. That is,

$$\mathcal{D}^m(\{\boldsymbol{S}: \forall h \in \mathcal{H}, |\boldsymbol{L}_{\mathcal{S}}(h) - \boldsymbol{L}_{\mathcal{D}}(h)| \leq \epsilon\}) \geq 1 - \delta.$$

Equivalently, we need to show that

$$\mathcal{D}^m(\{\boldsymbol{S}: \exists h \in \mathcal{H}, |L_{\mathcal{S}}(h) - L_{\mathcal{D}}(h)| > \epsilon\}) < \delta.$$

Notice that

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

Applying Hoeffding's inequality, then we obtain that

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

Hence

$$egin{aligned} &\mathcal{D}^m(\{m{S}:\exists h\in\mathcal{H},|L_{m{S}}(h)-L_{\mathcal{D}}(h)|>\epsilon\})\ &\leq &\sum_{h\in\mathcal{H}}\mathcal{D}^m(m{S}:|L_{m{S}}(h)-L_{\mathcal{D}}(h)|>\epsilon)\ &\leq &2|\mathcal{H}|\exp(-2m\epsilon^2). \end{aligned}$$

Finally, if we choose

$$m \geq rac{\log(2|\mathcal{H}|/\delta)}{2\epsilon^2},$$

then

$$\mathcal{D}^{m}(\{\boldsymbol{S}: \exists h \in \mathcal{H}, |L_{\mathcal{S}}(h) - L_{\mathcal{D}}(h)| > \epsilon\}) < \delta. \square$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへぐ

Outline

Empirical Risk Minimization

Probably Approximately Correct Learning

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

Agnostic PAC learnability

Uniform Convergence

The Bias-Complexity Tradeoff

Theorem (No-Free-Lunch)

Let A be any learning algorithm for the task of binary classification with respect to the 0-1 loss over a domain \mathcal{X} . Let m be any number smaller than $|\mathcal{X}|/2$, representing a training set size. Then, there exists a distribution \mathcal{D} over $\mathcal{X} \times \{0,1\}$ such that:

- (1) There exists a function $f : \mathcal{X} \to \{0, 1\}$ with $L_{\mathcal{D}}(f) = 0$.
- (2) With probability of at least $\frac{1}{7}$ over the choice of $S \sim \mathcal{D}^m$ we have that $L_{\mathcal{D}}(\mathcal{A}(S)) \geq \frac{1}{8}$.

・ロト ・ 同 ・ ・ ヨ ・ ・ ヨ ・ うへつ

- How does the No-Free-Lunch result relate to the need for prior knowledge?
- Let us consider an ERM predictor over the hypothesis class H of all functions f from X to {0,1}.
- This class represents lack of prior knowledge: Every possible function from X to {0, 1} is considered a good candidate.

Corollary

Let \mathcal{X} be an infinite domain set and let \mathcal{H} be the set of all functions from \mathcal{X} to $\{0,1\}$. Then, \mathcal{H} is not PAC learnable.

・ロト ・ 同 ・ ・ ヨ ・ ・ ヨ ・ うへつ

Proof. Assume, by way of contradiction, that the class is learnable. Choose some $\epsilon < 1/8$ and $\delta < 1/7$. By the definition of PAC learnability, there must be some learning algorithm *A* and an integer $m = m(\epsilon, \delta)$, such that for any data-generating distribution over $\mathcal{X} \times \{0, 1\}$, if for some function $f : \mathcal{X} \to \{0, 1\}$, $L_{\mathcal{D}}(f) = 0$, then with probability greater than $1 - \delta$ when *A* is applied to samples *S* of size *m*, generated i.i.d. by \mathcal{D} , $L_{\mathcal{D}}(A(S)) \le \epsilon$. However, applying the No-Free-Lunch theorem, since

 $|\mathcal{X}| > 2m$, for the algorithm *A*, there exists a distribution \mathcal{D} such that with probability greater than $1/7 \ge \delta$, $L_{\mathcal{D}}(A(S)) > 1/8 > \epsilon$, which leads to the desired contradiction. \Box

Error Decomposition:

• Let h_S be an ERM_H hypothesis, then

$$L_{\mathcal{D}}(h_{\mathcal{S}}) = [L_{\mathcal{D}}(h_{\mathcal{S}}) - \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h)] + \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h).$$

$$\mathcal{L}_{\mathcal{D}}(h_{\mathcal{S}}) - \epsilon_{\text{Bayes}} = [\mathcal{L}_{\mathcal{D}}(h_{\mathcal{S}}) - \min_{h \in \mathcal{H}} \mathcal{L}_{\mathcal{D}}(h)] + [\min_{h \in \mathcal{H}} \mathcal{L}_{\mathcal{D}}(h) - \epsilon_{\text{Bayes}}].$$

 Approximation Error: measures how much inductive bias we have.

$$\epsilon_{\mathrm{app}} = \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h).$$

$$\epsilon_{\mathrm{app}} = \min_{h \in \mathcal{H}} \mathcal{L}_{\mathcal{D}}(h) - \epsilon_{\mathrm{Bayes}}.$$

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

 Enlarging the hypothesis class can decrease the approximation error. Error Decomposition:

▶ Let h_S be an ERM_H hypothesis, then

$$L_{\mathcal{D}}(h_{\mathcal{S}}) = [L_{\mathcal{D}}(h_{\mathcal{S}}) - \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h)] + \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h).$$

Estimation Error: the difference between the minimum risk achievable by a predictor in the hypothesis class and the error achieved by the ERM predictor.

$$\epsilon_{\rm est} = L_{\mathcal{D}}(h_{\mathcal{S}}) - \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h).$$

- The quality of estimation error depends on the training set size and the size, or complexity, of the hypothesis set.
- ► For a finite hypothesis case, eest increases (logarithmically) with |*H*| and decrease with *m*.

The bias-complexity tradeoff

- Choosing H to be a very rich class decreases the approximation error but at the same time might increase the estimation error, as a rich H might lead to overfitting.
- Choosing H to be a very small set reduce the estimation error but might increase the approximation error or, in other words, might lead to underfitting.
- Why not choose the class containing only the Bayes optimal classifier?
- Learning theory studies how rich we can make H while still maintaining reasonable estimation error.

In many cases, empirical research focuses on designing good hypothesis classes for a certain domain.