

DDA5001 Machine Learning

Training versus Testing (Part II)

Xiao Li

School of Data Science
The Chinese University of Hong Kong, Shenzhen



Recap: In-sample Error versus Out-of-sample Error

- **In-sample Error:** Given a set of **training samples** $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$,

$$\text{Er}_{\text{in}} = \frac{1}{n} \sum_{i=1}^n e(f_{\boldsymbol{\theta}}(\mathbf{x}_i), g(\mathbf{x}_i))$$

- **Out-of-sample Error:** Suppose data \mathbf{x} follows a certain distribution \mathcal{D} in an i.i.d. manner,

$$\text{Er}_{\text{out}} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e(f_{\boldsymbol{\theta}}(\mathbf{x}), g(\mathbf{x}))]$$

Remarks:

- The In-sample error Er_{in} is also known as the **training error**.
- The out-of-sample error Er_{out} is **more general than the test error**.
Fortunately, we can use the test error to approximate Er_{out} very well when the test dataset is large enough.

Recap: Concept of Training versus Testing / Generalization

Recall that learning is all about to infer g outside of the seen training dataset, i.e.,

Make the out-of-sample error small

- ▶ But Er_{out} is **not computable** at all.

Here is a simple decomposition:

$$Er_{out} = \underbrace{Er_{out} - Er_{in}}_{\text{generalization error}} + \underbrace{Er_{in}}_{\text{training error}}$$

- ▶ We need to simultaneously make generalization error and training error small, in order to make Er_{out} small.

The goal of generalization is to

explore how out-of-sample error is related to in-sample error, i.e., bound the generalization error

- ▶ The reason to explore this relationship is that Er_{in} is **computable, checkable, and even amenable**.

Recap: The Starting Point

Given training samples $\mathcal{S} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$.

In expectation for any $f \in \mathcal{H}$: (we omit θ in f_θ for simplicity)

$$\mathbb{E}_{\mathcal{S} \sim i.i.d. \mathcal{D}} [\text{Er}_{\text{in}}(f)] = \text{Er}_{\text{out}}(f).$$

- ▶ Law of large numbers: When $n \rightarrow \infty$, we have in-sample error estimates the out-of-sample error accurately.
- ▶ However, this is true only when $n \rightarrow \infty$.

\rightsquigarrow We will derive a **non-asymptotic** (finite n) result for $f \approx g$, which is analogous to the one derived for $\nu \approx \mu$.

Finite Hypothesis Space Generalization

VC Dimension

Additional Assumption: Finite Hypothesis Space

Assumption: Finite hypothesis space

The cardinality $|\mathcal{H}| < +\infty$

where $|\cdot|$ denotes the cardinality (number of elements) of a set. Namely, $|\mathcal{H}|$ measures the number of all possible $f_{\theta} \in \mathcal{H}$.

- ▶ It means that the number of possible f_{θ} in \mathcal{H} is finite. Is it a practical assumption?
- ▶ We will omit θ in f_{θ} in the sequel to ease notation. Remind yourself f is almost always parameterized by some parameter θ .

Generalization for Fixed f : A Lemma

Lemma: High probability bounds for fixed f

Fix any model $f : \mathcal{X} \mapsto \{-1, 1\}$ ($f \in \mathcal{H}$ is fixed). The following inequalities hold for any $t > 0$:

$$\Pr \left[\text{Er}_{\text{in}}(f) - \text{Er}_{\text{out}}(f) \geq t \right] \leq e^{-2nt^2},$$

and

$$\Pr \left[\text{Er}_{\text{in}}(f) - \text{Er}_{\text{out}}(f) \leq -t \right] \leq e^{-2nt^2}.$$

Thus, we have the two side tail probability bound

$$\Pr \left[|\text{Er}_{\text{in}}(f) - \text{Er}_{\text{out}}(f)| \geq t \right] \leq 2e^{-2nt^2}.$$

- ▶ **Non-asymptotic** bounds valid for any n .
- ▶ Equivalently, $\Pr \left[|\text{Er}_{\text{in}}(f) - \text{Er}_{\text{out}}(f)| \leq t \right] \geq 1 - 2e^{-2nt^2}$, which is a high probability bound.

Proof

Recall the Hoeffding's inequality:

Hoeffding's inequality for bounded random variables

Suppose X_i are **independent** random variables with mean μ_i and bounded on $[a_i, b_i]$ for $i = 1, \dots, n$, then for any $t > 0$, we have

$$\Pr \left[\sum_{i=1}^n (X_i - \mu_i) \geq t \right] \leq e^{-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}}$$

Note that $e(f(\mathbf{x}_i), g(\mathbf{x}_i))$ equals to either 0 or 1, we have

$$\begin{aligned} & \Pr \left[\text{Er}_{\text{in}}(f) - \text{Er}_{\text{out}}(f) \geq t \right] \\ &= \Pr \left[\frac{1}{n} \sum_{i=1}^n (e(f(\mathbf{x}_i), g(\mathbf{x}_i))) - \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n e(f(\mathbf{x}_i), g(\mathbf{x}_i)) \right] \geq t \right] \\ &= \Pr \left[\sum_{i=1}^n (e(f(\mathbf{x}_i), g(\mathbf{x}_i))) - \mathbb{E}[e(f(\mathbf{x}_i), g(\mathbf{x}_i))] \geq nt \right] \\ &\leq e^{-2nt^2}. \end{aligned}$$

Generalization for Fixed Model f

Proposition: Generalization for **fixed f**

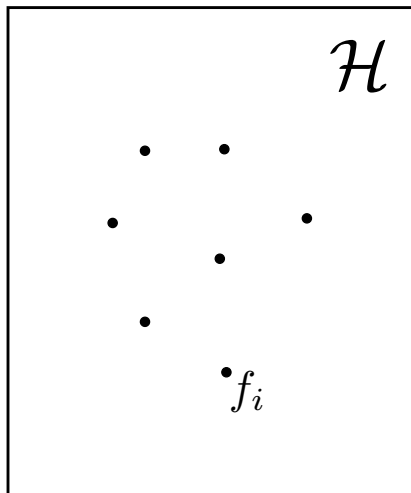
Fix a model $f : \mathcal{X} \mapsto \{-1, 1\}$. For any $\delta > 0$, the following generalization bound holds **with probability at least $1 - \delta$**

$$\text{Er}_{\text{out}}(f) \leq \text{Er}_{\text{in}}(f) + \sqrt{\frac{\log\left(\frac{2}{\delta}\right)}{2n}}$$

Proof: Letting $\delta = 2e^{-2nt^2}$ and solving for t yields the desired result.

Are we done?

What We Need is Uniform Bound for All Possible $f \in \mathcal{H}$



Generalization for Finite Hypothesis Space

Theorem: Generalization for **finite hypothesis space**

Let \mathcal{H} be a **finite** hypothesis space, i.e., $|\mathcal{H}| < \infty$. For any $\delta > 0$, the following generalization bound holds **with probability at least $1 - \delta$**

$$\forall f \in \mathcal{H} \quad \text{Er}_{\text{out}}(f) \leq \text{Er}_{\text{in}}(f) + \sqrt{\frac{\log\left(\frac{2|\mathcal{H}|}{\delta}\right)}{2n}} \quad (1)$$

- ▶ The dependence on δ is only **logarithmically**.
- ▶ The generalization error increase when $|\mathcal{H}|$ grows, but only **logarithmically**.
- ▶ **More samples (larger n) lead to better generalization** (always true in practice).

Trade-off

$$\forall f \in \mathcal{H} \quad \text{Er}_{\text{out}}(f) \leq \text{Er}_{\text{in}}(f) + \sqrt{\frac{\log\left(\frac{2|\mathcal{H}|}{\delta}\right)}{2n}}$$

On the **training** side, we need

more complex hypothesis \mathcal{H} (larger $|\mathcal{H}|$)

On the **generalization** side, we need

less complex hypothesis \mathcal{H} (smaller $|\mathcal{H}|$)

Proof

The proof is done by applying union bound. Let $f_1, \dots, f_{|\mathcal{H}|}$ be the elements of \mathcal{H} . We have

$$\begin{aligned} & \Pr \left[\exists f \in \mathcal{H} \text{ s.t. } |\text{Er}_{\text{in}}(f) - \text{Er}_{\text{out}}(f)| \geq t \right] \\ &= \Pr \left[|\text{Er}_{\text{in}}(f_1) - \text{Er}_{\text{out}}(f_1)| \geq t, \text{ or } \dots \right. \\ & \quad \left. \text{or } |\text{Er}_{\text{in}}(f_{|\mathcal{H}|}) - \text{Er}_{\text{out}}(f_{|\mathcal{H}|})| \geq t \right] \\ & \stackrel{\text{union bound}}{\leq} \sum_{i=1}^{|\mathcal{H}|} \Pr \left[|\text{Er}_{\text{in}}(f_i) - \text{Er}_{\text{out}}(f_i)| \geq t \right] \\ & \leq 2|\mathcal{H}|e^{-2nt^2} \end{aligned}$$

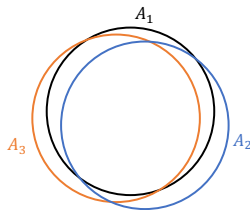
Setting the RHS probability to δ yields the desired result.

Issue with Finite Hypothesis Space

However, $|\mathcal{H}|$ is usually infinite in practice

- ▶ Think about the the simplest (one-dimensional) linear model $f_{\theta}(x) = \theta x$, $\theta \in \mathbb{R}$. We have infinitely many choices for θ .
- ▶ If we let $|\mathcal{H}|$ tends to ∞ , we will have a trivial bound: Some finite thing $< \infty$.

Issue: Union bound can be very coarse.



In this case, $\Pr[A_1 \text{ or } A_2 \text{ or } A_3]$ is much less than $\sum_{i=1}^3 \Pr[A_i]$.

Solution: Find a way to measure the complexity of \mathcal{H} more smartly.

Finite Hypothesis Space Generalization

VC Dimension

Road Map

- ▶ Instead of using $|\mathcal{H}|$ directly to count the complexity of \mathcal{H} , we have to properly account for the overlaps of different hypotheses.
- ▶ In this way, our goal is to replace the number of hypotheses $|\mathcal{H}|$ with an effective number **which is finite** even when $|\mathcal{H}|$ is infinite.

↪ This quantity will be the so-called **VC dimension**, which is of **combinatorial nature**.

- ▶ The VC dimension captures how different the hypotheses in \mathcal{H} are, and hence how much overlap the different hypotheses have.
- ▶ Using this new notion, we will show that we can replace $|\mathcal{H}|$ in the obtained generalization bound with VC dimension.

Before that, let us introduce the related notion called **dichotomy** and **growth function**.

Dichotomy

- ▶ If $f \in \mathcal{H}$ is applied to a **finite** sample set $\{x_1, \dots, x_n\}$, we get n -tuple $\{f(x_1), \dots, f(x_n)\}$ of ± 1 's.
- ▶ Such a n tuple is called a **dichotomy** since it splits $\{x_1, \dots, x_n\}$ into **two** groups: The points for which f is $+1$ and those for which f is -1 .
- ▶ Each $f \in \mathcal{H}$ generates a dichotomy on $\{x_1, \dots, x_n\}$, **but two different f 's may generate the same dichotomy**.

We can now define the dichotomies of the whole hypothesis space \mathcal{H} .

Dichotomies of \mathcal{H}

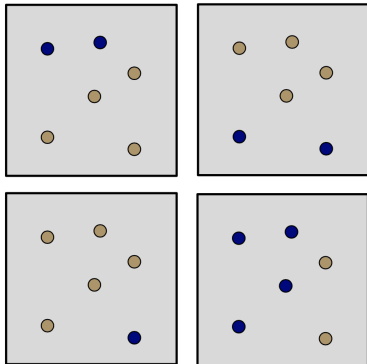
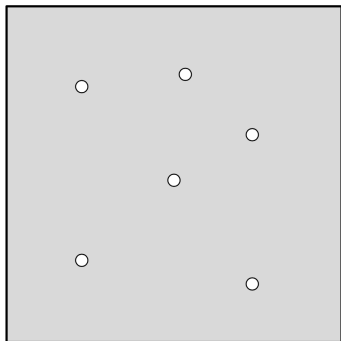
Given $\{x_1, \dots, x_n\}$, the dichotomies generated by \mathcal{H} on these points are defined by

$$\mathcal{H}(x_1, \dots, x_n) = \{(f(x_1), \dots, f(x_n)) : f \in \mathcal{H}\}.$$

- ▶ One can think of the dichotomies $\mathcal{H}(x_1, \dots, x_n)$ as a set of hypothesis just like \mathcal{H} is, except that the hypotheses are only seen through the n data points.
- ▶ Larger $\mathcal{H}(x_1, \dots, x_n)$ means that \mathcal{H} is more diverse / rich.

Example of Dichotomies

We have five points. There are four different $f \in \mathcal{H}$ on the points.



Growth Function

Growth function is a **number**, which is defined based on the number of dichotomies.

Growth function

The growth function for the hypothesis set \mathcal{H} is defined as:

$$\mathcal{G}_{\mathcal{H}}(n) = \max_{\{x_1, \dots, x_n\} \subseteq \mathcal{X}} |\mathcal{H}(x_1, \dots, x_n)|,$$

where $|\cdot|$ denotes the cardinality (number of elements) of a set.

Instead of counting the size of \mathcal{H} by $|\mathcal{H}|$, the idea of growth function is:
Using \mathcal{H} , what is the maximum number of ways we can label a n -points dataset?

Properties of Growth Function

- ▶ $G_{\mathcal{H}}(n)$ counts the **most dichotomies** that can possibly be generated on **any** n points in \mathcal{X} .
- ▶ To compute $\mathcal{G}_{\mathcal{H}}(n)$, we consider all possible choice of n points, and pick the one that gives us the most dichotomies, which is of **combinatorial** nature.
- ▶ Similar to $|\mathcal{H}|$, $\mathcal{G}_{\mathcal{H}}(n)$ is a measure of the **richness** of the hypothesis set \mathcal{H} . The difference is that it is now considered on n points rather than the entire input space \mathcal{X} .
- ▶ Since $\mathcal{H}(\mathbf{x}_1, \dots, \mathbf{x}_n) \subset \{-1, +1\}^n$ (the set of all possible dichotomies on any n points). Clearly, we have

$$\mathcal{G}_{\mathcal{H}}(n) \leq 2^n.$$

- ▶ If \mathcal{H} is capable of generating all possible dichotomies on $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, then $\mathcal{H}(\mathbf{x}_1, \dots, \mathbf{x}_n) = \{-1, +1\}^n$, i.e., $\mathcal{G}_{\mathcal{H}}(n) = 2^n$, and we say that \mathcal{H} can **shatter** the data points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$.

Vapnik-Chervonenkis (VC) Dimension

We now introduce a well known notion — Vapnik-Chervonenkis (VC) dimension.

VC dimension

The VC dimension of a hypothesis space \mathcal{H} , denoted by $d_{VC}(\mathcal{H})$ or simply d_{VC} , is the largest n so that it can be shattered by \mathcal{H} , i.e.,

$$d_{VC}(\mathcal{H}) := \max\{n : \mathcal{G}_{\mathcal{H}}(n) = 2^n\}.$$

If $\mathcal{G}_{\mathcal{H}}(n) = 2^n$ for all n , then $d_{VC}(\mathcal{H}) = \infty$.

- ▶ By definition, VC dimension indicates the representation power of \mathcal{H} .
- ▶ $d_{VC} + 1$ counts the number of data points n that \mathcal{H} starts to not shatter.

Fact:

Bounding Growth Function using VC Dimension

$$\mathcal{G}_{\mathcal{H}}(n) \leq n^{d_{VC}} + 1.$$

↪ Next lecture: VC dimension generalization.