# CS182: Introduction to Machine Learning – Learning Theory (Finite Case & Infinite Case)

Yujiao Shi

SIST, ShanghaiTech

Spring, 2025

# Finite Case

# Statistical Learning Theory Model

1. Data points are generated i.i.d. from some *unknown* distribution

$$\boldsymbol{x}^{(n)} \sim p^*(\boldsymbol{x})$$

2. Labels are generated from some *unknown* function

$$y^{(n)} = c^*\left(\boldsymbol{x}^{(n)}\right)$$

3. The learning algorithm chooses the hypothesis (or classifier) with lowest *training* error rate from a specified hypothesis set, $\mathcal{H}$

4. Goal: return a hypothesis (or classifier) with low *true* error rate

# Types of Error

- True error rate
  - Actual quantity of interest in machine learning
  - How well your hypothesis will perform on average across all possible data points

- Test error rate
  - Used to evaluate hypothesis performance
  - Good estimate of your hypothesis's true error

- Validation error rate
  - Used to set hypothesis hyperparameters
  - Slightly "optimistic" estimate of your hypothesis's true error

- Training error rate
  - Used to set model parameters
  - Very "optimistic" estimate of your hypothesis's true error

# Types of Risk (a.k.a. Error)

- Expected risk of a hypothesis $h$ (a.k.a. true error)

$$R(h) = P_{\boldsymbol{x} \sim p^*}\big(c^*(\boldsymbol{x}) \neq h(\boldsymbol{x})\big)$$

- Empirical risk of a hypothesis $h$ (a.k.a. training error)

$$\hat{R}(h) = P_{\boldsymbol{x} \sim \mathcal{D}}\big(c^*(\boldsymbol{x}) \neq h(\boldsymbol{x})\big)$$

$$= \frac{1}{N}\sum_{n=1}^{N} \mathbb{1}\left(c^*(\boldsymbol{x}^{(n)}) \neq h(\boldsymbol{x}^{(n)})\right)$$

$$= \frac{1}{N}\sum_{n=1}^{N} \mathbb{1}\left(y^{(n)} \neq h(\boldsymbol{x}^{(n)})\right)$$

where $\mathcal{D} = \left\{\left(\boldsymbol{x}^{(n)}, y^{(n)}\right)\right\}_{n=1}^{N}$ is the training data set and $\boldsymbol{x} \sim \mathcal{D}$ denotes a point sampled uniformly at random from $\mathcal{D}$

## Three Hypotheses of Interest

1. The *true function, $c^*$*

2. The *expected risk minimizer,*

$$h^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \, R(h)$$

3. The *empirical risk minimizer,*

$$\hat{h} = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \, \hat{R}(h)$$

Poll Question 1: Which of the following are *always* true?

A. $c^* = h^*$
B. $c^* = \hat{h}$
C. $h^* = \hat{h}$
D. $c^* = h^* = \hat{h}$
E. None of the above
**F. TOXIC**

- The *true function, $c^*$*

- The *expected risk minimizer,*
$$h^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} R(h)$$

- The *empirical risk minimizer,*

$$\hat{h} = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \hat{R}(h)$$

# Key Question

- Given a hypothesis with zero/low training error, what can we say about its true error?

# PAC Learning

- PAC = **P**robably **A**pproximately **C**orrect

- PAC Criterion:

$$P\left(\left|R(h) - \hat{R}(h)\right| \epsilon \geq\right) \mathcal{H} \ni h \, \forall \, \delta - 1 \leq$$

for some $\epsilon$ (difference between expected and empirical risk) and $\delta$ (probability of "failure")

- We want the PAC criterion to be satisfied for $\mathcal{H}$ with small values of $\epsilon$ and $\delta$

# Sample Complexity

- The sample complexity of an algorithm/hypothesis set, $\mathcal{H}$, is the number of labelled training data points needed to satisfy the PAC criterion for some $\delta$ and $\epsilon$

- Four cases

  - Realizable vs. Agnostic

    - Realizable $\rightarrow c^* \in \mathcal{H}$

    - Agnostic $\rightarrow c^*$ might or might not be in $\mathcal{H}$

  - Finite vs. Infinite

    - Finite $\rightarrow |\mathcal{H}| < \infty$

    - Infinite $\rightarrow |\mathcal{H}| = \infty$

## Theorem 1: Finite, Realizable Case

- For a finite hypothesis set $\mathcal{H}$ s.t. $c^* \in \mathcal{H}$ and arbitrary distribution $p^*$, if the number of labelled training data points satisfies

$$M \geq \frac{1}{\epsilon} \left( \ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right) \right)$$

then with probability at least $1 - \delta$, all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$

## Proof of Theorem 1: Finite, Realizable Case

1. Assume there are $K$ "bad" hypotheses in $\mathcal{H}$, i.e.,

   $h_1, h_2, \ldots, h_K$ that all have $R(h_k) > \epsilon$

2. Pick one bad hypothesis, $h_k$

   A. Probability that $h_k$ correctly classifies the first training data point $< 1 - \epsilon$

   B. Probability that $h_k$ correctly classifies all $M$ training data points $< (1 - \epsilon)^M$

3. Probability that at least one bad hypothesis correctly classifies all $M$ training data points $=$
   $P(h_1$ correctly classifies all $M$ training data points $\cup\ h_2$ correctly classifies all $M$ training data points $\cup$

   $\vdots$

   $\cup\ h_K$ correctly classifies all $M$ training data points$)$

# Proof of Theorem 1: Finite, Realizable Case

$P(h_1$ correctly classifies all $M$ training data points $\cup$

$h_2$ correctly classifies all $M$ training data points $\cup$

$\vdots$

$\cup\ h_K$ correctly classifies all $M$ training data points)

$$\leq \sum_{k=1}^{K} P(h_k \text{ correctly classifies all } M \text{ training data points})$$

by the union bound: $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

$$\leq P(A) + P(B)$$

# Proof of Theorem 1: Finite, Realizable Case

$$\sum_{k=1}^{K} P(h_k \text{ correctly classifies all } M \text{ training data points})$$

$$< k(1-\epsilon)^M \leq |\mathcal{H}|(1-\epsilon)^M$$

because $k \leq |\mathcal{H}|$

3. Probability that at least one bad hypothesis correctly classifies all $M$ training data points $\leq |\mathcal{H}|(1-\epsilon)^M$

4. Using the fact that $1 - x \leq \exp(-x) \ \forall \ x$,
   $$|\mathcal{H}|(1-\epsilon)^M \leq |\mathcal{H}| \exp(-\epsilon)^M = |\mathcal{H}| \exp(-M\epsilon)$$

5. Probability that at least one bad hypothesis correctly classifies all $M$ training data points $\leq |\mathcal{H}| \exp(-M\epsilon)$, which we want to be *low*, i.e., $|\mathcal{H}| \exp(-M\epsilon) \leq \delta$

# Proof of Theorem 1: Finite, Realizable Case

$$|\mathcal{H}| \exp(-M\epsilon) \leq \delta \rightarrow \exp(-M\epsilon) \leq \frac{\delta}{|\mathcal{H}|}$$

$$\rightarrow -M\epsilon \leq \ln\left(\frac{\delta}{|\mathcal{H}|}\right)$$

$$\rightarrow M \geq \frac{1}{\epsilon}\left(-\ln\left(\frac{\delta}{|\mathcal{H}|}\right)\right)$$

$$\rightarrow M \geq \frac{1}{\epsilon}\left(\ln\left(\frac{|\mathcal{H}|}{\delta}\right)\right)$$

$$\rightarrow M \geq \frac{1}{\epsilon}\left(\ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right)\right)$$

# Proof of Theorem 1: Finite, Realizable Case

6. Given $M \geq \frac{1}{\epsilon}\left(\ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right)\right)$ labelled training data points, the probability that ∃ a bad hypothesis $h_k \in \mathcal{H}$ with $R(h_k) > \epsilon$ and $\hat{R}(h_k) = 0$ is $\leq \delta$

$\Updownarrow$

Given $M \geq \frac{1}{\epsilon}\left(\ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right)\right)$ labelled training data points, the probability that all hypotheses $h_k \in \mathcal{H}$ with $R(h_k) > \epsilon$ have $\hat{R}(h_k) > 0$ is $\geq 1 - \delta$

# Aside: Proof by Contrapositive

- The contrapositive of a statement $A \Rightarrow B$ is $\neg B \Rightarrow \neg A$

- A statement and its contrapositive are logically equivalent, i.e., $A \Rightarrow B$ means that $\neg B \Rightarrow \neg A$

- Example: "it's raining $\Rightarrow$ Henry brings am umbrella"

    is the same as saying

    "Henry didn't bring an umbrella $\Rightarrow$ it's not raining "

# Proof of Theorem 1: Finite, Realizable Case

7. Given $M \geq \frac{1}{\epsilon}\left(\ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right)\right)$ labelled training data points, the probability that all hypotheses $h_k \in \mathcal{H}$ with $R(h_k) > \epsilon$ have $\hat{R}(h_k) > 0$ is $\geq 1 - \delta$

$\Updownarrow$

Given $M \geq \frac{1}{\epsilon}\left(\ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right)\right)$ labelled training data points, the probability that all hypotheses $h_k \in \mathcal{H}$ with $\hat{R}(h_k) = 0$ have $R(h_k) \leq \epsilon$ is $\geq 1 - \delta$

(proof by contrapositive)

## Theorem 1: Finite, Realizable Case

- For a finite hypothesis set $\mathcal{H}$ s.t. $c^* \in \mathcal{H}$ and arbitrary distribution $p^*$, if the number of labelled training data points satisfies

$$M \geq \frac{1}{\epsilon}\left(\ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right)\right)$$

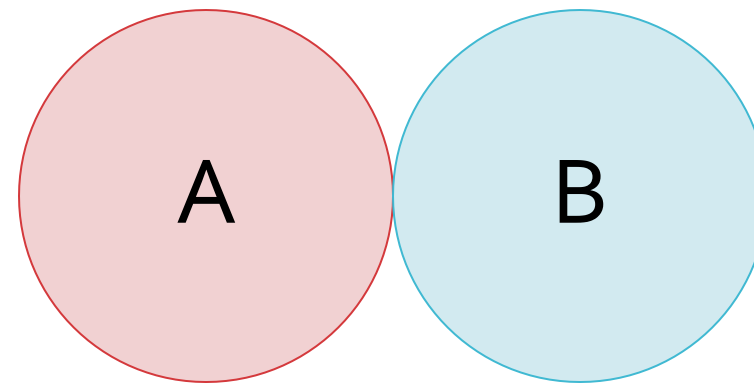then with probability at least $1 - \delta$, all $h \in \mathcal{H}$ with

$\hat{R}(h) = 0$ have $R(h) \leq \epsilon$

Making the bound tight and solving for $\epsilon$ gives...

## Statistical Learning Theory Corollary

- For a finite hypothesis set $\mathcal{H}$ s.t. $c^* \in \mathcal{H}$ and arbitrary distribution $p^*$, given a training data set $S$ s.t. $|S| = M$, all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have

$$R(h) \leq \frac{1}{M}\left(\ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right)\right)$$
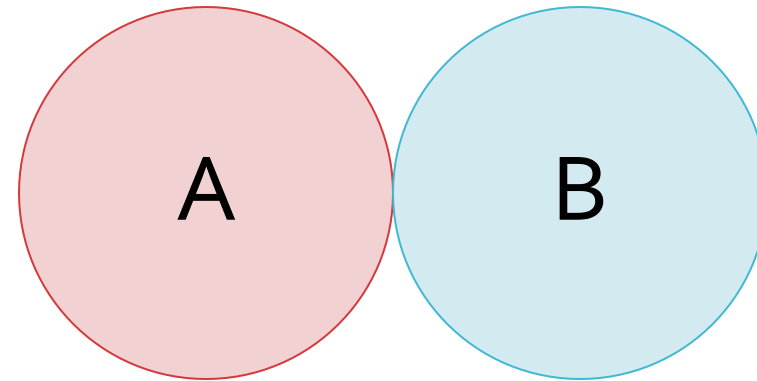
with probability at least $1 - \delta$.

## Theorem 2: Finite, Agnostic Case

- For a finite hypothesis set $\mathcal{H}$ and arbitrary distribution $p^*$, if the number of labelled training data points satisfies

$$M \geq \frac{1}{2\epsilon^2}\left(\ln(|\mathcal{H}|) + \ln\left(\frac{2}{\delta}\right)\right)$$

then with probability at least $1 - \delta$, all $h \in \mathcal{H}$ satisfy

$$\left|R(h) - \hat{R}(h)\right| \leq \epsilon$$

- Bound is inversely quadratic in $\epsilon$, e.g., halving $\epsilon$ means we need four times as many labelled training data points

- Again, making the bound tight and solving for $\epsilon$ gives...

# Statistical Learning Theory Corollary

- For a finite hypothesis set $\mathcal{H}$ and arbitrary distribution $p^*$, given a training data set $S$ s.t. $|S| = M$, all $h \in \mathcal{H}$ have

$$R(h) \leq \hat{R}(h) + \sqrt{\frac{1}{2M}\left(\ln(|\mathcal{H}|) + \ln\left(\frac{2}{\delta}\right)\right)}$$

with probability at least $1 - \delta$.

# Infinite Case

What happens when $|\mathcal{H}| = \infty$?

- For a finite hypothesis set $\mathcal{H}$ and arbitrary distribution $p^*$, given a training data set $S$ s.t. $|S| = M$, all $h \in \mathcal{H}$ have

$$R(h) \leq \hat{R}(h) + \sqrt{\frac{1}{2M}\left(\ln(|\mathcal{H}|) + \ln\left(\frac{2}{\delta}\right)\right)}$$

with probability at least $1 - \delta$.

## What happens when $|\mathcal{H}| = \infty$?

- For a finite hypothesis set $\mathcal{H}$ and arbitrary distribution $p^*$, given a training data set $S$ s.t. $|S| = M$, all $h \in \mathcal{H}$ have

$$R(h) \leq \hat{R}(h) + \sqrt{\frac{1}{2M}\left(\ln(|\mathcal{H}|) + \ln\left(\frac{2}{\delta}\right)\right)}$$

with probability at least $1 - \delta$.

$$P\{A \cup B\} \leq P\{A\} + P\{B\}$$

# The Union Bound…

$$P\{A \cup B\} \le P\{A\} + P\{B\}$$

$$P\{A \cup B\} = P\{A\} + P\{B\} - P\{A \cap B\}$$

# The Union Bound is Bad

# Intuition

If two hypotheses $h_1, h_2 \in \mathcal{H}$ are very similar, then the events

- "$h_1$ is consistent with the first $m$ training data points"
- "$h_2$ is consistent with the first $m$ training data points"

will overlap a lot!

# Intuition

If two hypotheses $h_1, h_2 \in \mathcal{H}$ are very similar, then the events

- "$h_1$ is consistent with the first $m$ training data points"
- "$h_2$ is consistent with the first $m$ training data points"

will overlap a lot!

# Labellings

- Given some finite set of data points $S = \left( \boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(M)} \right)$ and some hypothesis $h \in \mathcal{H}$, applying $h$ to each point in $S$ results in a **labelling**

  - $\left( h(\boldsymbol{x}^{(1)}), \ldots, h(\boldsymbol{x}^{(M)}) \right)$ is a vector of $M$ +1's and -1's

- Given $S = \left( \boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(M)} \right)$, each hypothesis in $\mathcal{H}$ induces a labelling but not necessarily a unique labelling

  - The set of labellings induced by $\mathcal{H}$ on $S$ is

$$\mathcal{H}(S) = \left\{ \left( h(\boldsymbol{x}^{(1)}), \ldots, h(\boldsymbol{x}^{(M)}) \right) \,\middle|\, h \in \mathcal{H} \right\}$$

# Example: Labellings

$$\mathcal{H} = \{h_1, h_2, h_3\}$$



$h_3$

$h_1$     $h_2$

上海科技大学
ShanghaiTech University

$\mathcal{H} = \{h_1, h_2, h_3\}$

$\left(h_1(\boldsymbol{x}^{(1)}), h_1(\boldsymbol{x}^{(2)}), h_1(\boldsymbol{x}^{(3)}), h_1(\boldsymbol{x}^{(4)})\right)$
$= (-1, +1, -1, +1)$

$\boldsymbol{x}^{(4)}$



$\boldsymbol{x}^{(2)}$

$\boldsymbol{x}^{(1)}$

$\boldsymbol{x}^{(3)}$

$\boldsymbol{x}^{(4)}$

$h_1$

上 海 科 技 大 学
ShanghaiTech University

$\mathcal{H} = \{h_1, h_2, h_3\}$

$\left(h_2\left(x^{(1)}\right), h_2\left(x^{(2)}\right), h_2\left(x^{(3)}\right), h_2\left(x^{(4)}\right)\right)$
$= (-1, +1, -1, +1)$

$x^{(2)}$

$x^{(1)}$

$x^{(3)}$

$x^{(4)}$

$h_2$

# Example: Labellings

$\mathcal{H} = \{h_1, h_2, h_3\}$

$\left( h_3\left(\boldsymbol{x}^{(1)}\right), h_3\left(\boldsymbol{x}^{(2)}\right), h_3\left(\boldsymbol{x}^{(3)}\right), h_3\left(\boldsymbol{x}^{(4)}\right) \right)$
$= (+1, +1, -1, -1)$

$\boldsymbol{x}^{(2)}$

$\boldsymbol{x}^{(1)}$

$h_3$

$\boldsymbol{x}^{(3)}$

$\boldsymbol{x}^{(4)}$

上海科技大学
ShanghaiTech University

$\mathcal{H} = \{h_1, h_2, h_3\}$

$\mathcal{H}(S)$
$= \{(+1, +1, -1, -1), (-1, +1, -1, +1)\}$

$|\mathcal{H}(S)| = 2$

上海科技大学
ShanghaiTech University

$\mathcal{H} = \{h_1, h_2, h_3\}$

$\mathcal{H}(S) = \{(+1, +1, -1, -1)\}$

$|\mathcal{H}(S)| = 1$



$x^{(1)}$

$x^{(2)}$

$h_3$

$x^{(3)}$

$x^{(4)}$

$h_1$ $h_2$

# VC-Dimension

- $\mathcal{H}(S)$ is the set of all labellings induced by $\mathcal{H}$ on $S$
  - If $|S| = M$, then $|\mathcal{H}(S)| \leq 2^M$
  - $\mathcal{H}$ **shatters** $S$ if $|\mathcal{H}(S)| = 2^M$

- The **VC-dimension** of $\mathcal{H}$, $VC(\mathcal{H})$, is the size of the largest set $S$ that can be shattered by $\mathcal{H}$.
  - If $\mathcal{H}$ can shatter arbitrarily large finite sets, then $d_{VC}(\mathcal{H}) = \infty$

- To prove that $VC(\mathcal{H}) = d$, you need to show
  1. $\exists$ some set of $d$ data points that $\mathcal{H}$ can shatter and
  2. $\nexists$ a set of $d+1$ data points that $\mathcal{H}$ can shatter
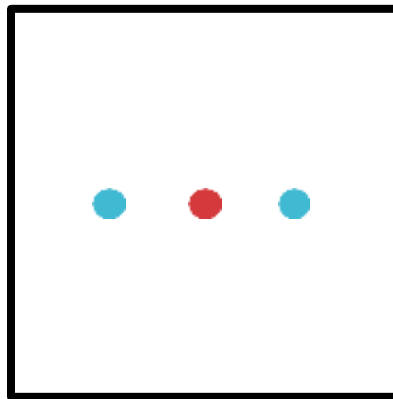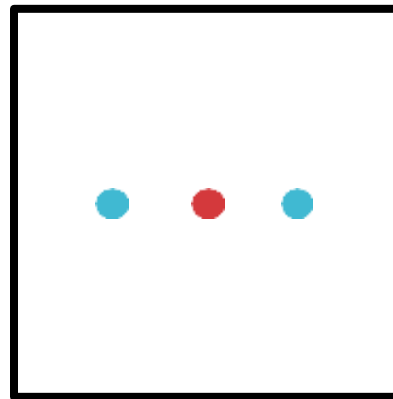
## VC-Dimension: Example

- $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators

- What is $VC(\mathcal{H})$?

  - Can $\mathcal{H}$ shatter some set of 1 point?
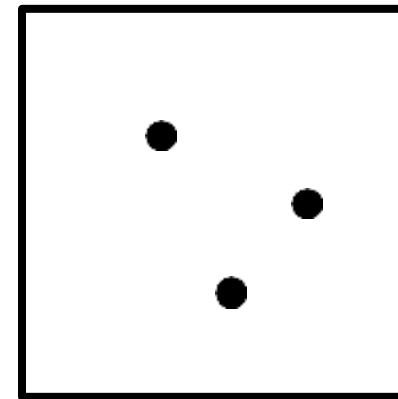
$S$

## VC-Dimension: Example

- $x \in \mathbb{R}^2$ and $\mathcal{H} = $ all 2-dimensional linear separators

- What is $VC(\mathcal{H})$?
  - Can $\mathcal{H}$ shatter some set of 1 point?
  - Can $\mathcal{H}$ shatter some set of 2 points?

$S$

## VC-Dimension: Example

- $x \in \mathbb{R}^2$ and $\mathcal{H} = $ all 2-dimensional linear separators

- What is $VC(\mathcal{H})$?
  - Can $\mathcal{H}$ shatter some set of 1 point?
  - Can $\mathcal{H}$ shatter some set of 2 points?
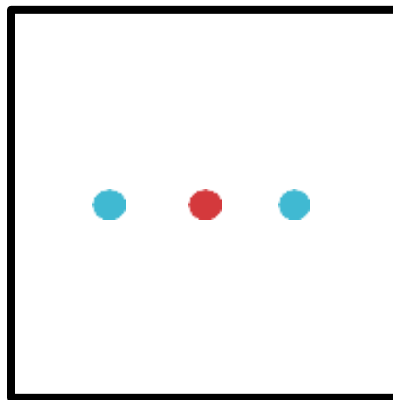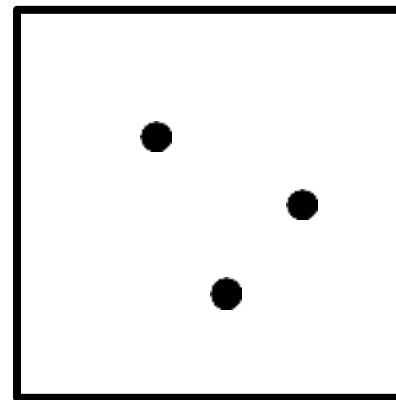  - Can $\mathcal{H}$ shatter some set of 3 points?



$S$

# VC-Dimension: Example

- $x \in \mathbb{R}^2$ and $\mathcal{H} = $ all 2-dimensional linear separators

- What is $VC(\mathcal{H})$?
  - Can $\mathcal{H}$ shatter some set of 1 point?
  - Can $\mathcal{H}$ shatter some set of 2 points?
  - Can $\mathcal{H}$ shatter some set of 3 points?



$S$

# VC-Dimension: Example

- $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators

- What is $VC(\mathcal{H})$?
  - Can $\mathcal{H}$ shatter some set of 1 point?
  - Can $\mathcal{H}$ shatter some set of 2 points?
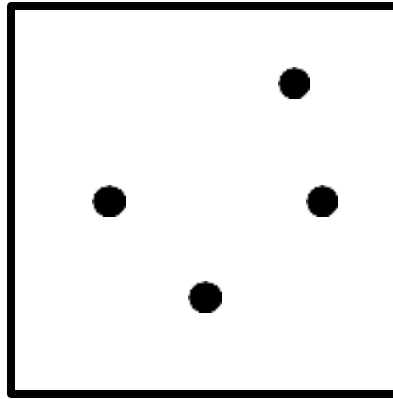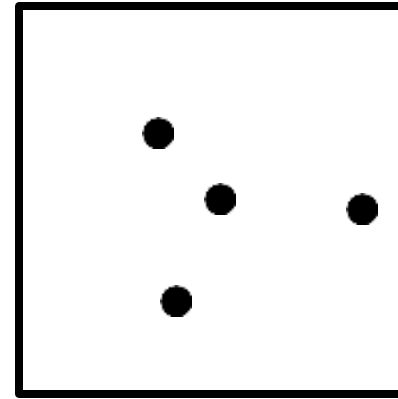  - Can $\mathcal{H}$ shatter some set of 3 points?

$S$

## VC-Dimension: Example

- $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators

- What is $VC(\mathcal{H})$?
  - Can $\mathcal{H}$ shatter some set of 1 point?
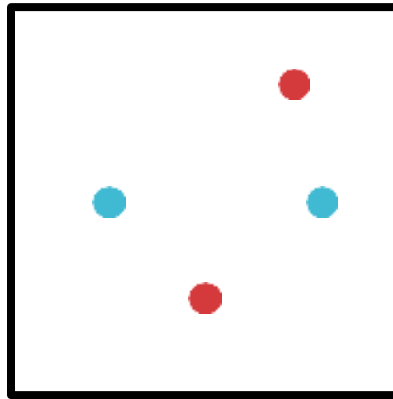  - Can $\mathcal{H}$ shatter some set of 2 points?
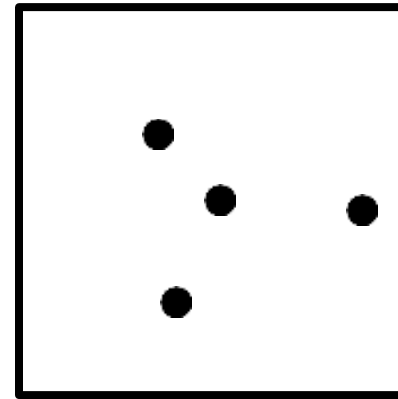  - Can $\mathcal{H}$ shatter **some** set of 3 points?

$S$

# VC-Dimension: Example

- $x \in \mathbb{R}^2$ and $\mathcal{H} = $ all 2-dimensional linear separators

- What is $VC(\mathcal{H})$?

  - Can $\mathcal{H}$ shatter some set of 1 point?
  - Can $\mathcal{H}$ shatter some set of 2 points?
  - Can $\mathcal{H}$ shatter some set of 3 points?

$S_1$

$S_2$

- $x \in \mathbb{R}^2$ and $\mathcal{H} = $ all 2-dimensional linear separators

- What is $VC(\mathcal{H})$?

  - Can $\mathcal{H}$ shatter some set of 1 point?
  - Can $\mathcal{H}$ shatter some set of 2 points?
  - Can $\mathcal{H}$ shatter some set of 3 points?

## VC-Dimension: Example



$$|\mathcal{H}(S_1)| = 6 \qquad\qquad |\mathcal{H}(S_2)| = 8$$

# VC-Dimension: Example

- $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators

- What is $VC(\mathcal{H})$?
  - Can $\mathcal{H}$ shatter some set of 1 point?
  - Can $\mathcal{H}$ shatter some set of 2 points?
  - Can $\mathcal{H}$ shatter some set of 3 points?
  - Can $\mathcal{H}$ shatter some set of 4 points?
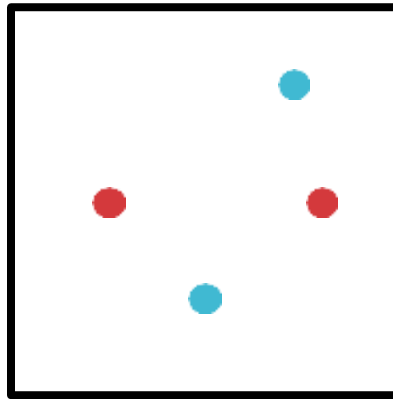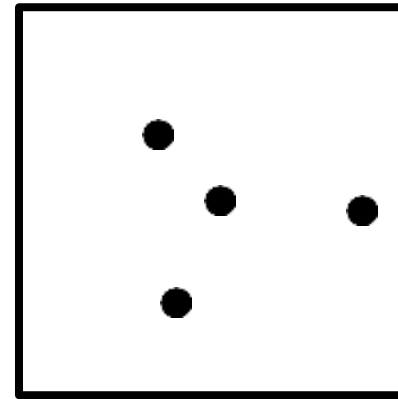
$S_1$

All points on the convex hull

$S_2$

At least one point inside the convex hull

- $x \in \mathbb{R}^2$ and $\mathcal{H} = $ all 2-dimensional linear separators

- What is $VC(\mathcal{H})$?

  - Can $\mathcal{H}$ shatter some set of 1 point?
  - Can $\mathcal{H}$ shatter some set of 2 points?
  - Can $\mathcal{H}$ shatter some set of 3 points?
  - Can $\mathcal{H}$ shatter some set of 4 points?

# VC-Dimension: Example
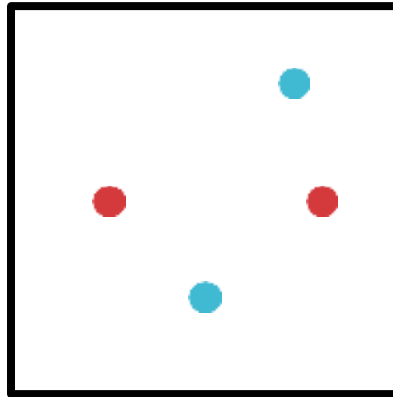


$S_1$

All points on the convex hull

$S_2$

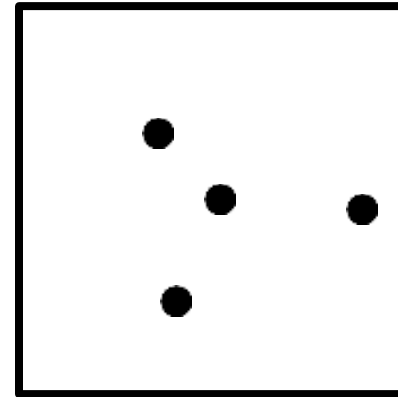At least one point inside the convex hull

# VC-Dimension: Example

- $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators

- What is $VC(\mathcal{H})$?

  - Can $\mathcal{H}$ shatter some set of 1 point?
  - Can $\mathcal{H}$ shatter some set of 2 points?
  - Can $\mathcal{H}$ shatter some set of 3 points?
  - Can $\mathcal{H}$ shatter some set of 4 points?



$S_1$

All points on the
convex hull

$S_2$

At least one point
inside the convex hull

## VC-Dimension: Example

- $x \in \mathbb{R}^2$ and $\mathcal{H}$ = all 2-dimensional linear separators

- What is $VC(\mathcal{H})$?

  - Can $\mathcal{H}$ shatter some set of 1 point?
  - Can $\mathcal{H}$ shatter some set of 2 points?
  - Can $\mathcal{H}$ shatter some set of 3 points?
  - Can $\mathcal{H}$ shatter some set of 4 points?



$|\mathcal{H}(S_1)| = 14$
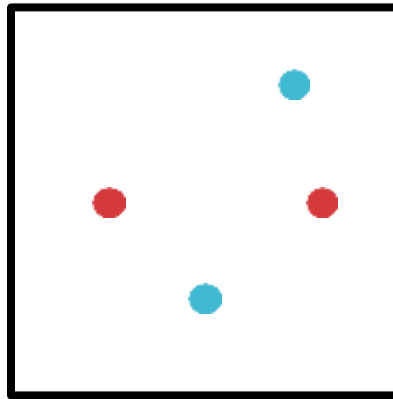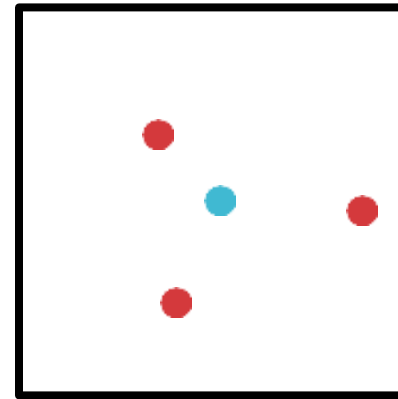
All points on the convex hull

$S_2$

At least one point inside the convex hull

- $x \in \mathbb{R}^2$ and $\mathcal{H}$ = all 2-dimensional linear separators

- What is $VC(\mathcal{H})$?

  - Can $\mathcal{H}$ shatter some set of 1 point?
  - Can $\mathcal{H}$ shatter some set of 2 points?
  - Can $\mathcal{H}$ shatter some set of 3 points?
  - Can $\mathcal{H}$ shatter some set of 4 points?

# VC-Dimension: Example



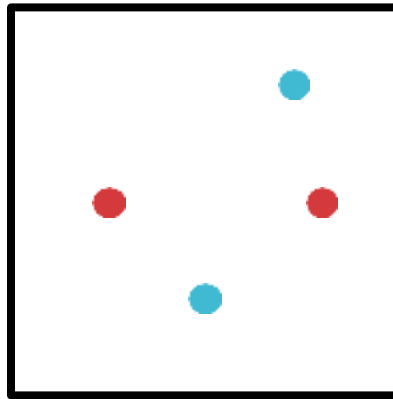$|\mathcal{H}(S_1)| = 14$

All points on the convex hull



$S_2$
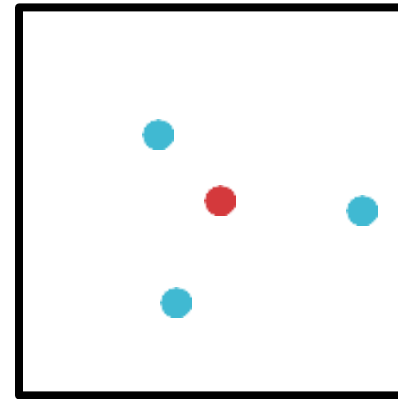
At least one point inside the convex hull

# VC-Dimension: Example

- $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators

- What is $VC(\mathcal{H})$?

  - Can $\mathcal{H}$ shatter some set of 1 point?
  - Can $\mathcal{H}$ shatter some set of 2 points?
  - Can $\mathcal{H}$ shatter some set of 3 points?
  - Can $\mathcal{H}$ shatter some set of 4 points?



$|\mathcal{H}(S_1)| = 14$

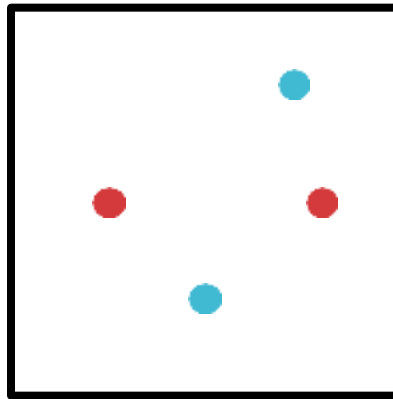All points on the convex hull
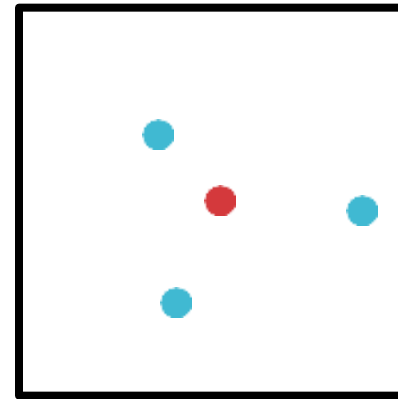
$S_2$

At least one point inside the convex hull

# VC-Dimension: Example

- $x \in \mathbb{R}^2$ and $\mathcal{H} = $ all 2-dimensional linear separators

- What is $VC(\mathcal{H})$?

  - Can $\mathcal{H}$ shatter some set of 1 point?
  - Can $\mathcal{H}$ shatter some set of 2 points?
  - Can $\mathcal{H}$ shatter some set of 3 points?
  - Can $\mathcal{H}$ shatter some set of 4 points?

$|\mathcal{H}(S_1)| = 14$

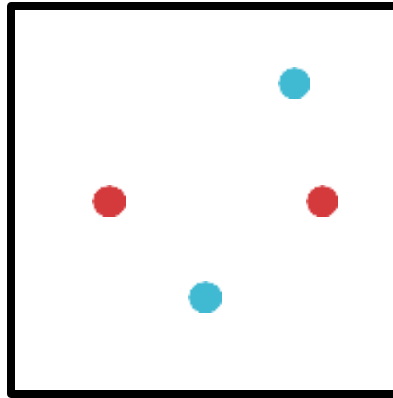All points on the
convex hull

$|\mathcal{H}(S_2)| = 14$

At least one point
inside the convex hull

# VC-Dimension: Example

- $x \in \mathbb{R}^2$ and $\mathcal{H} = $ all 2-dimensional linear separators

- $VC(\mathcal{H}) = 3$

  - Can $\mathcal{H}$ shatter some set of 1 point?
  - Can $\mathcal{H}$ shatter some set of 2 points?
  - Can $\mathcal{H}$ shatter some set of 3 points?
  - Can $\mathcal{H}$ shatter some set of 4 points?



$|\mathcal{H}(S_1)| = 14$

All points on the convex hull



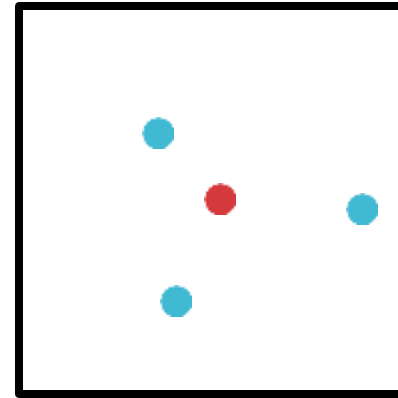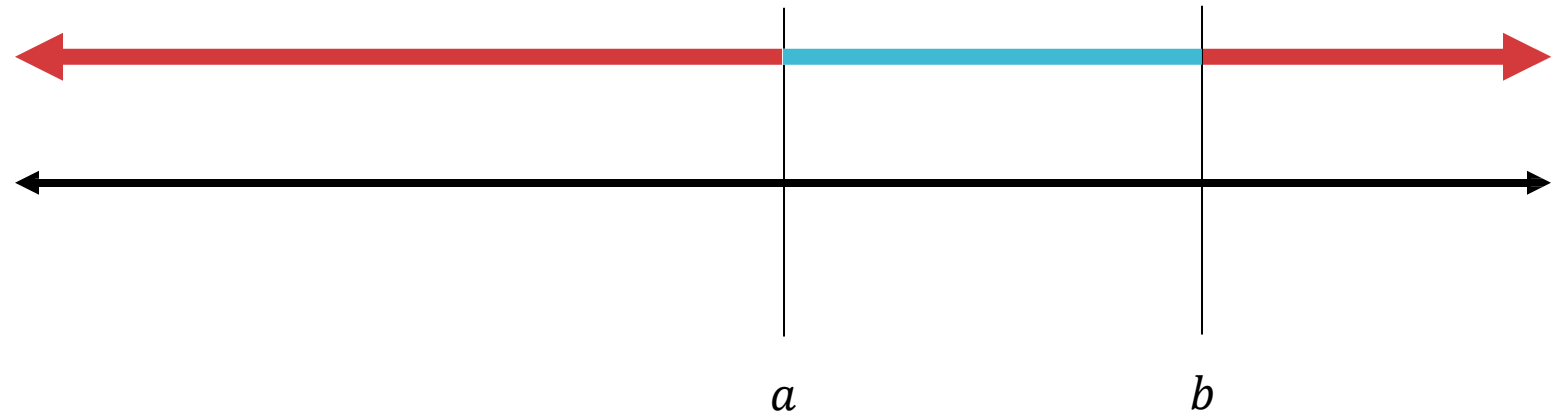$|\mathcal{H}(S_2)| = 14$

At least one point inside the convex hull

# VC-Dimension: Example

- $x \in \mathbb{R}^d$ and $\mathcal{H} =$ all $d$-dimensional linear separators

- $VC(\mathcal{H}) = d + 1$

- $x \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive intervals
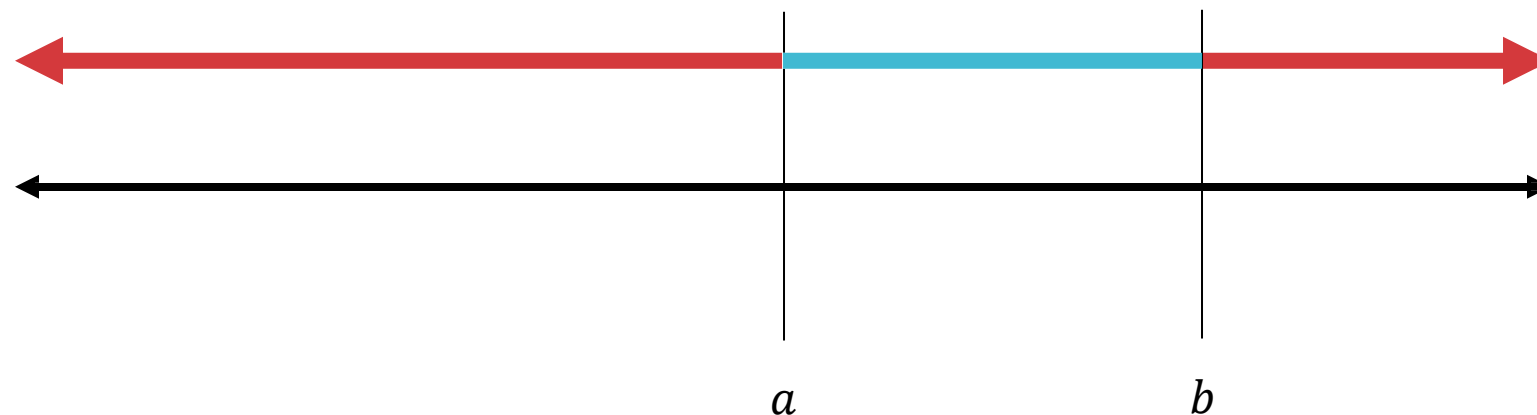
# VC-Dimension: Example



$a$          $b$

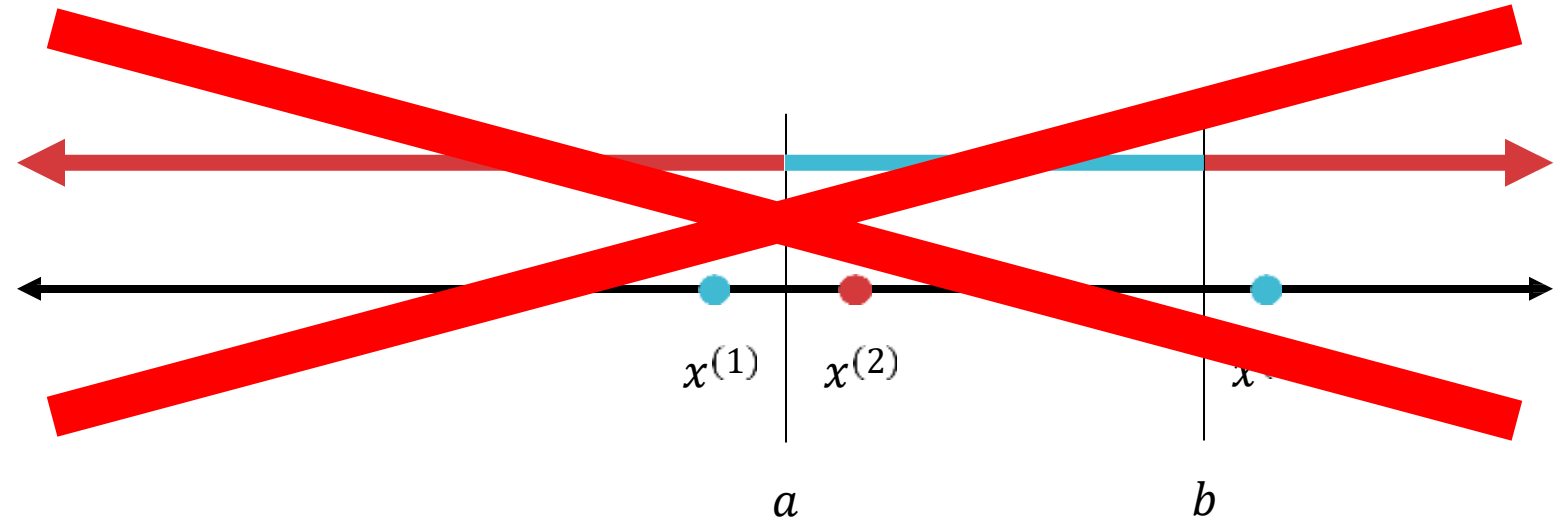# Poll Question 1:

What is $VC(\mathcal{H})$?

A. 0
B. 1
C. 1.5 **(TOXIC)**
D. 2
E. 3

- $x \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive intervals



$a$  $b$

VC-Dimension: Example

- $x \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive intervals

$x^{(1)}$  $x^{(2)}$  $x$

$a$  $b$

- $VC(\mathcal{H}) = 2$

## Theorem 3: Vapnik-Chervonenkis (VC)-Bound

- Infinite, realizable case: for any hypothesis set $\mathcal{H}$ and distribution $p^*$, if the number of labelled training data points satisfies

$$M = O\left(\frac{1}{\epsilon}\left(VC(\mathcal{H})\log\left(\frac{1}{\epsilon}\right) + \log\left(\left(\frac{1}{\delta}\right)\right)\right)\right)$$

then with probability at least $1 - \delta$, all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \epsilon$

## Statistical Learning Theory Corollary 3

- Infinite, realizable case: for any hypothesis set $\mathcal{H}$ and distribution $p^*$, given a training data set $S$ s.t. $|S| = M$, all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have

$$R(h) \leq O\left(\frac{1}{M}\left(VC(\mathcal{H})\log\left(\frac{M}{VC(\mathcal{H})}\right) + \log\left(\frac{1}{\delta}\right)\right)\right)$$

with probability at least $1 - \delta$.

# Theorem 4: Vapnik-Chervonenkis (VC)-Bound

- Infinite, agnostic case: for any hypothesis set $\mathcal{H}$ and distribution $p^*$, if the number of labelled training data points satisfies

$$M = O\left(\frac{1}{\epsilon^2}\left(VC(\mathcal{H}) + \log\left(\frac{1}{\delta}\right)\right)\right)$$

then with probability at least $1 - \delta$, all $h \in \mathcal{H}$ have

$$\left|R(h) - \hat{R}(h)\right| \leq \epsilon$$

# Statistical Learning Theory Corollary 4

- Infinite, agnostic case: for any hypothesis set $\mathcal{H}$ and distribution $p^*$, given a training data set $S$ s.t. $|S| = M$, all $h \in \mathcal{H}$ have

$$R(h) \leq \hat{R}(h) + O\left(\sqrt{\frac{1}{M}\left(VC(\mathcal{H}) + \log\left(\frac{1}{\delta}\right)\right)}\right)$$

with probability at least $1 - \delta$.

# Approximation Generalization Tradeoff

How well does $h$ generalize?

$$R(h) \leq \hat{R}(h) + O\left(\sqrt{\frac{1}{M}\left(VC(\mathcal{H}) + \log\left(\frac{1}{\delta}\right)\right)}\right)$$

How well does $h$ approximate $c*$?

# Approximation Generalization Tradeoff

Increases as $VC(\mathcal{H})$ increases

$$R(h) \leq \hat{R}(h) + O\left(\sqrt{\frac{1}{M}\left(VC(\mathcal{H}) + \log\left(\frac{1}{\delta}\right)\right)}\right)$$

Decreases as $VC(\mathcal{H})$ increases

**Learning Theory Learning Objectives**

You should be able to...

- Identify the properties of a learning setting and assumptions required to ensure low generalization error
- Distinguish true error, train error, test error
- Define PAC and explain what it means to be approximately correct and what occurs with high probability
- Apply sample complexity bounds to real-world machine learning examples