

10-701: Introduction to Machine Learning

Lecture 20 – Learning Theory (Infinite Case)

Hoda Heidari

* Slides adopted from F24 offering of 10701 by Henry Chai.

Statistical Learning Theory Model

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

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

2. Labels are generated from some *unknown* function

$$y^{(n)} = c^*(\mathbf{x}^{(n)})$$

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 Risk (a.k.a. Error)

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

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

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

$$\begin{aligned}\hat{R}(h) &= P_{\mathbf{x} \sim \mathcal{D}}(c^*(\mathbf{x}) \neq h(\mathbf{x})) \\ &= \frac{1}{N} \sum_{n=1}^N \mathbb{1}(c^*(\mathbf{x}^{(n)}) \neq h(\mathbf{x}^{(n)})) \\ &= \frac{1}{N} \sum_{n=1}^N \mathbb{1}(y^{(n)} \neq h(\mathbf{x}^{(n)}))\end{aligned}$$

where $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$ is the training data set with \mathbf{x}^i denoting a point sampled uniformly at random from p^*

Three Hypotheses of Interest

1. The *true function*, c^*

2. The *expected risk minimizer*,

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

3. The *empirical risk minimizer*,

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

Key Question

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

Sample Complexity & PAC Learnability

- A hypothesis class is PAC-learnable if for every $\epsilon, \delta \in (0, 1)$, there exists a sample size $m(\epsilon, \delta)$ polynomial in $1/\epsilon$ and $1/\delta$, such that with m i.i.d. samples from ANY distribution p^* the algorithm outputs a hypothesis whose generalization error is at most ϵ with probability at least $1 - \delta$.

PAC-learnability Results

- 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 $R(h) \leq \epsilon$.

- Making the bound tight (setting the two sides equal to each other) and solving for ϵ gives...

Proof Steps

1. Consider a hypothesis h with $R(h) > \epsilon$. Show that the probability of $\hat{R}(h) = 0$ is bounded.
2. Suppose there are k hypothesis $h \in \mathcal{H}$ with $R(h) > \epsilon$. Use union bound to upper bound the likelihood of at least one of them having $\hat{R}(h) = 0$.
3. Upper bound k with $|\mathcal{H}|$.
4. Set the above upper bound to be less than or equal δ to obtain the statement of the theorem.

Statistical Learning Theory Corollary: Finite, Realizable Case

- 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 $1 - \delta$.

Statistical Learning Theory Corollary: Finite, Agnostic Case

- 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

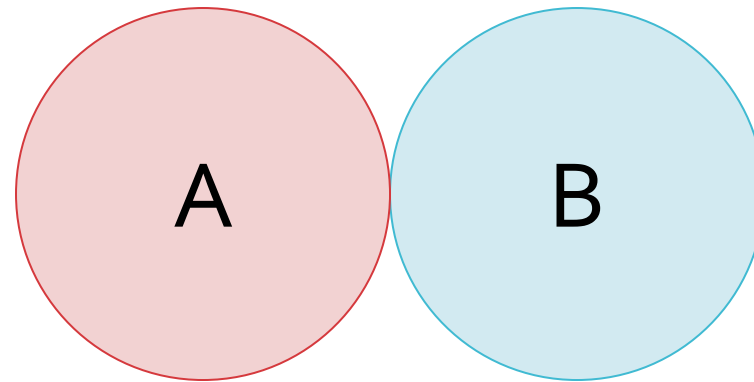
$$R(h) \leq \hat{R}(h) + \sqrt{\frac{1}{2M} \left(\ln(|\mathcal{H}|) + \ln\left(\frac{2}{\delta}\right) \right)}$$

with probability at least $1 - \delta$.

The Union
Bound is not
tight!

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

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

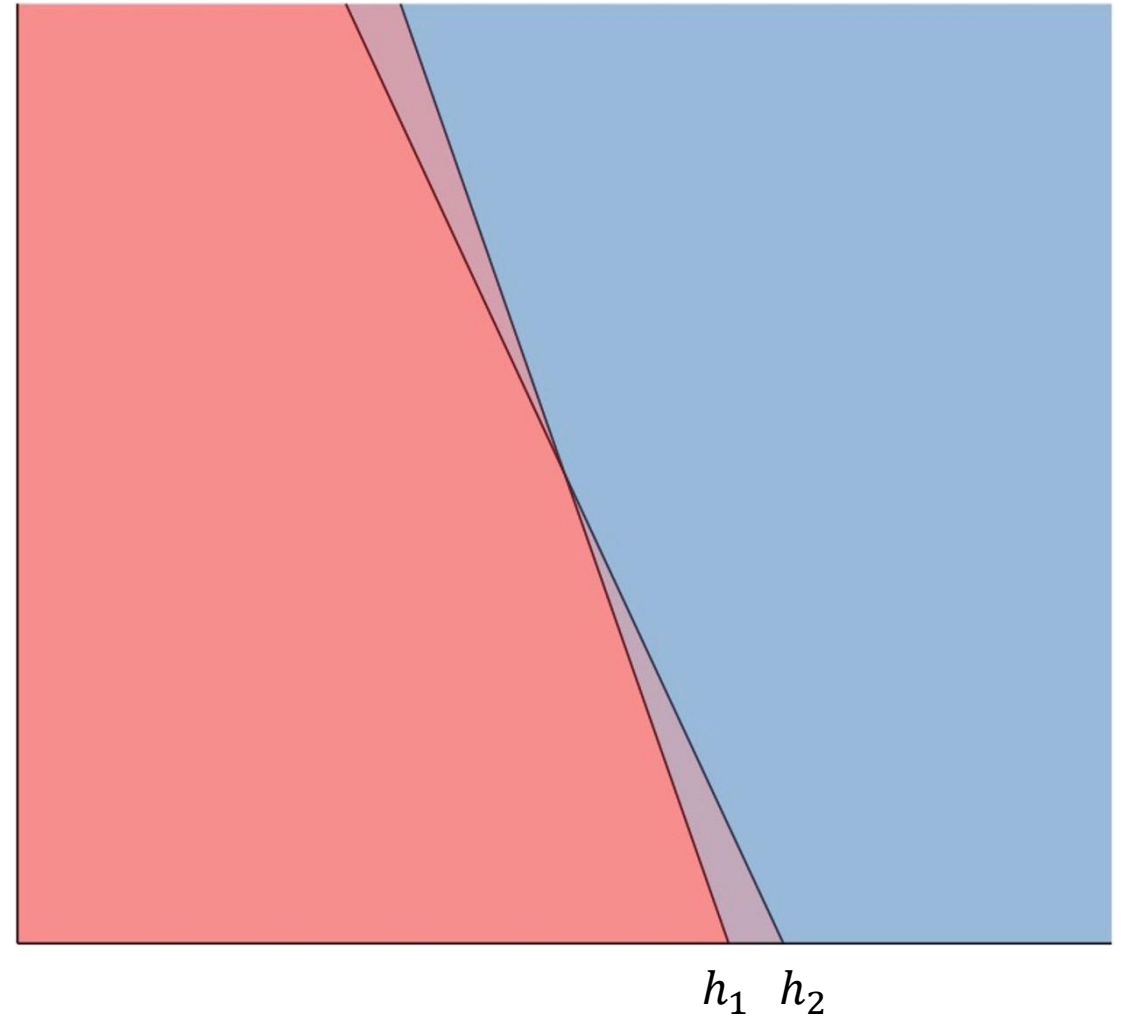


Intuition

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!

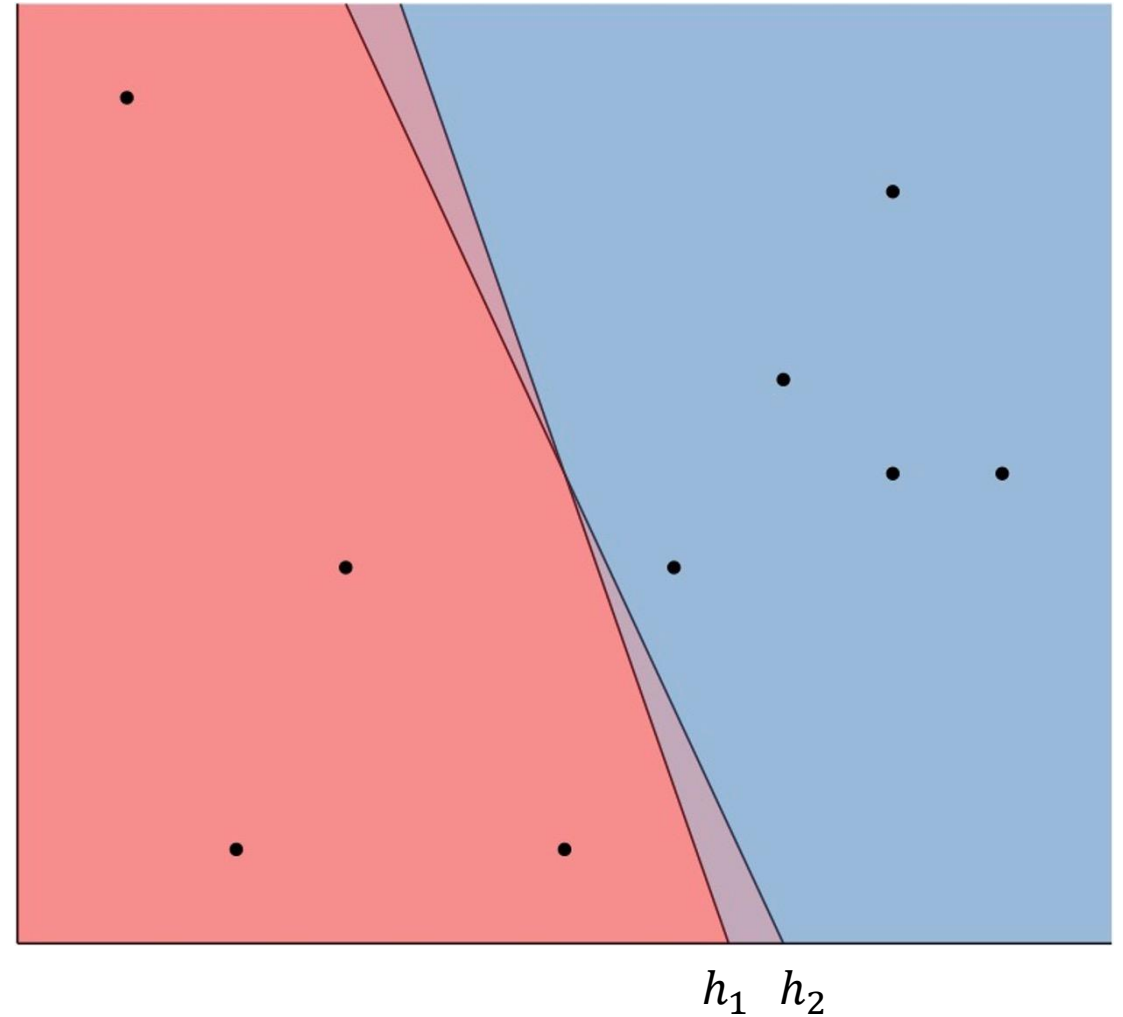


Intuition

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!



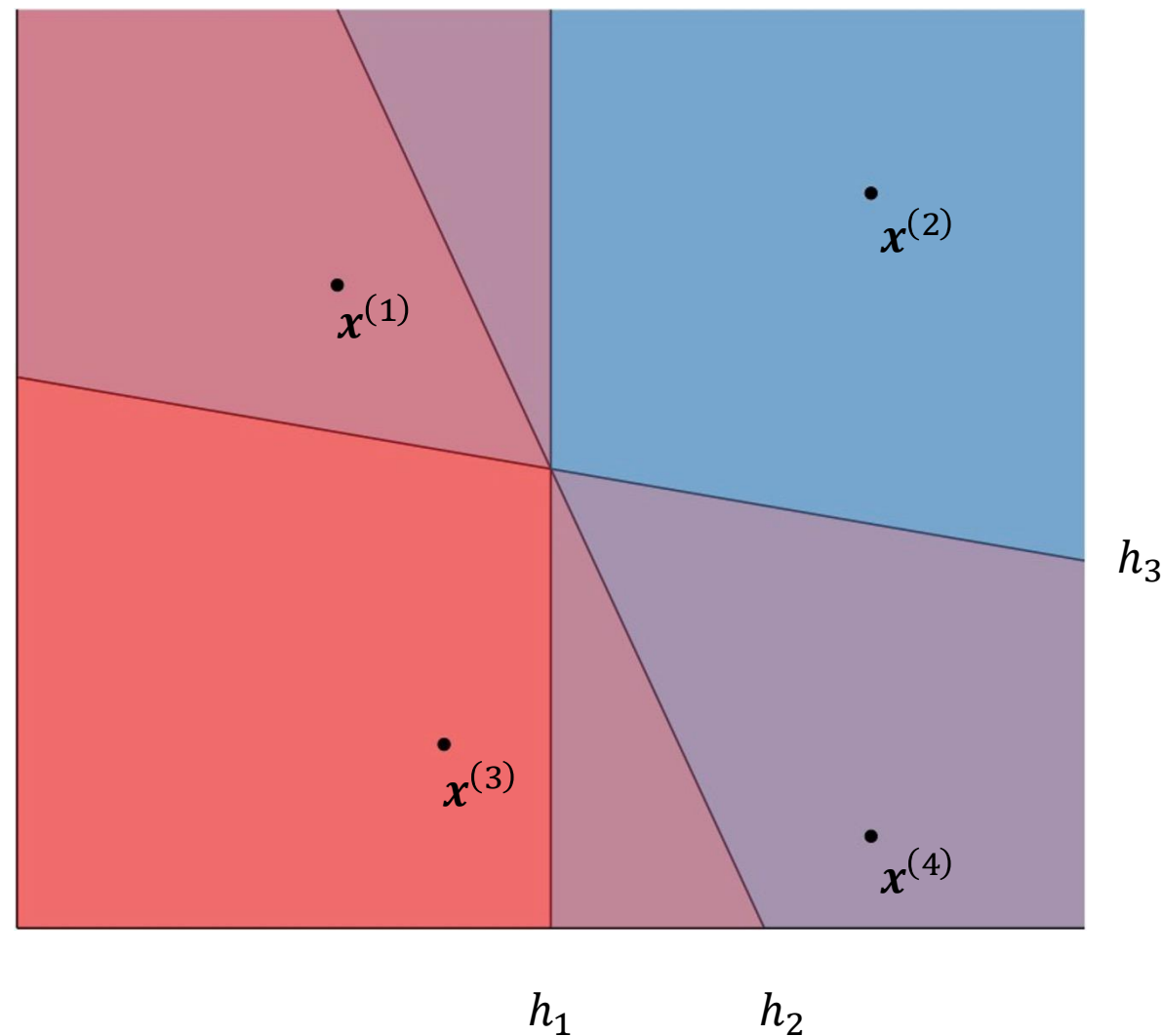
Labellings

- Given some finite set of data points $S = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)})$ and some hypothesis $h \in \mathcal{H}$, applying h to each point in S results in a labelling
 - $(h(\mathbf{x}^{(1)}), \dots, h(\mathbf{x}^{(M)}))$ is a vector of M +1's and -1's
- Insight: given $S = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)})$, 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

$$L_{\mathcal{H}}(S) = \left\{ \left(h(\mathbf{x}^{(1)}), \dots, h(\mathbf{x}^{(M)}) \right) \mid h \in \mathcal{H} \right\}$$

Example: Labellings

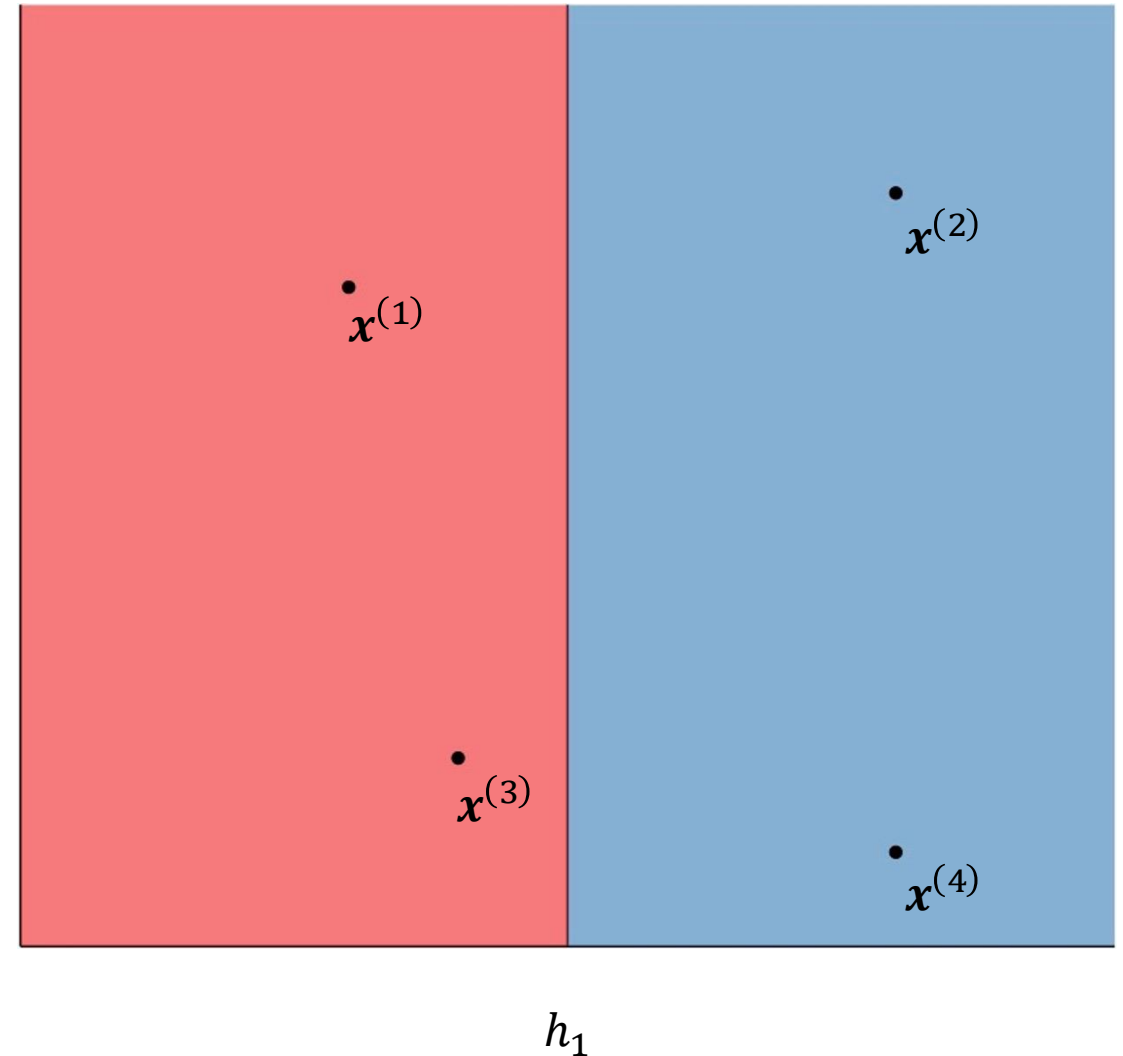
$$\mathcal{H} = \{h_1, h_2, h_3\}$$



Example: Labellings

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

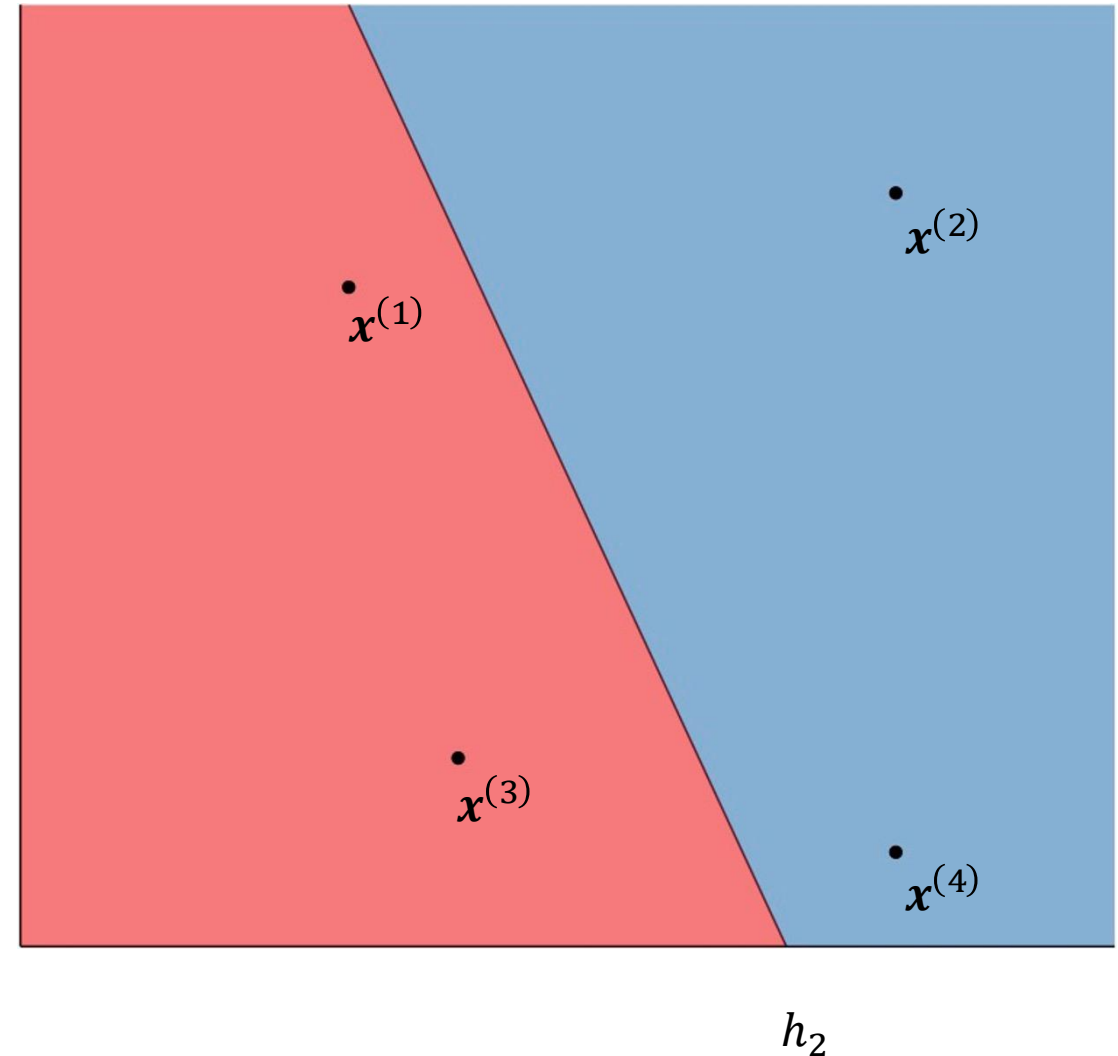
$$\begin{aligned} & \left(h_1(\mathbf{x}^{(1)}), h_1(\mathbf{x}^{(2)}), h_1(\mathbf{x}^{(3)}), h_1(\mathbf{x}^{(4)}) \right) \\ &= (-1, +1, -1, +1) \end{aligned}$$



Example: Labellings

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

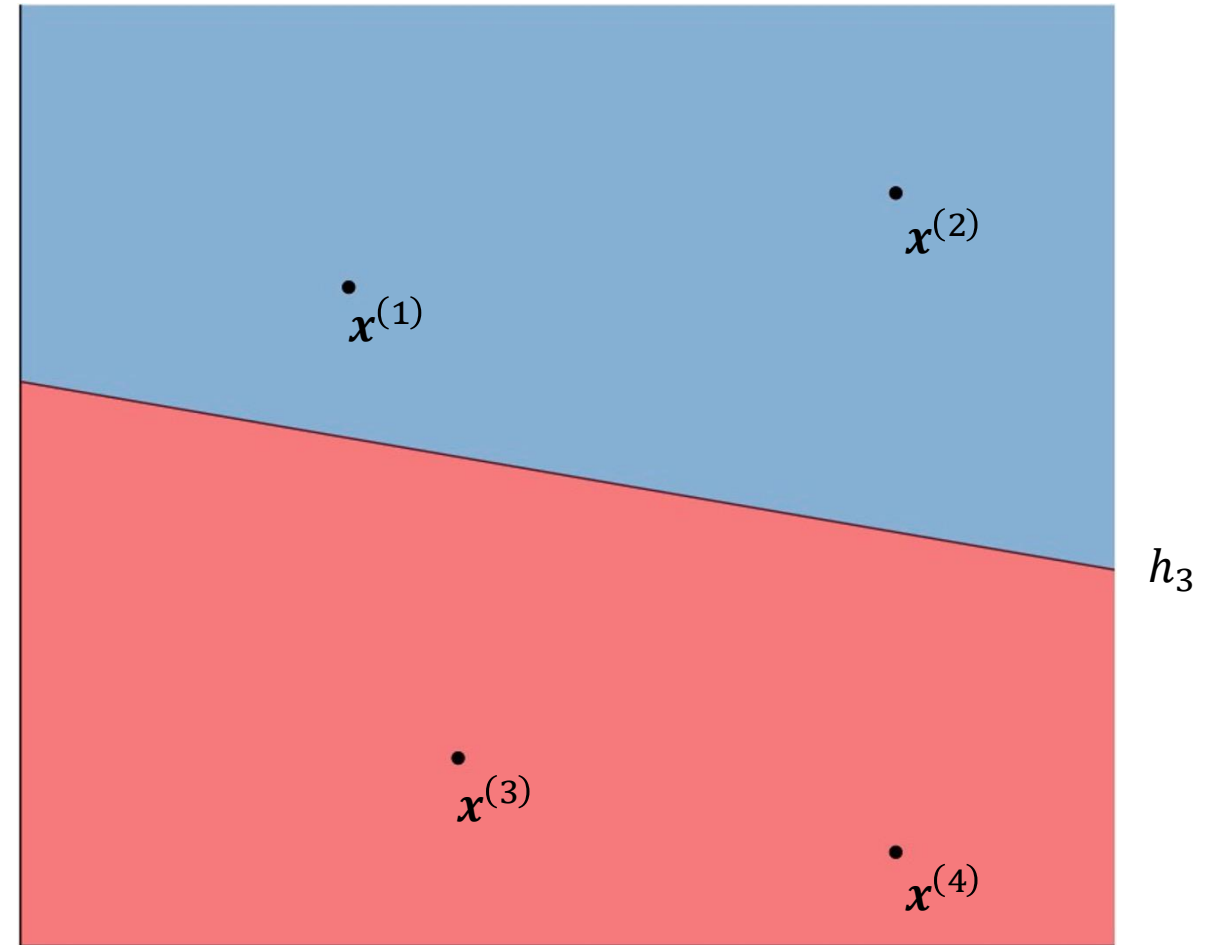
$$\begin{aligned} & \left(h_2(\mathbf{x}^{(1)}), h_2(\mathbf{x}^{(2)}), h_2(\mathbf{x}^{(3)}), h_2(\mathbf{x}^{(4)}) \right) \\ &= (-1, +1, -1, +1) \end{aligned}$$



Example: Labellings

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

$$\begin{aligned} & \left(h_3(\mathbf{x}^{(1)}), h_3(\mathbf{x}^{(2)}), h_3(\mathbf{x}^{(3)}), h_3(\mathbf{x}^{(4)}) \right) \\ &= (+1, +1, -1, -1) \end{aligned}$$

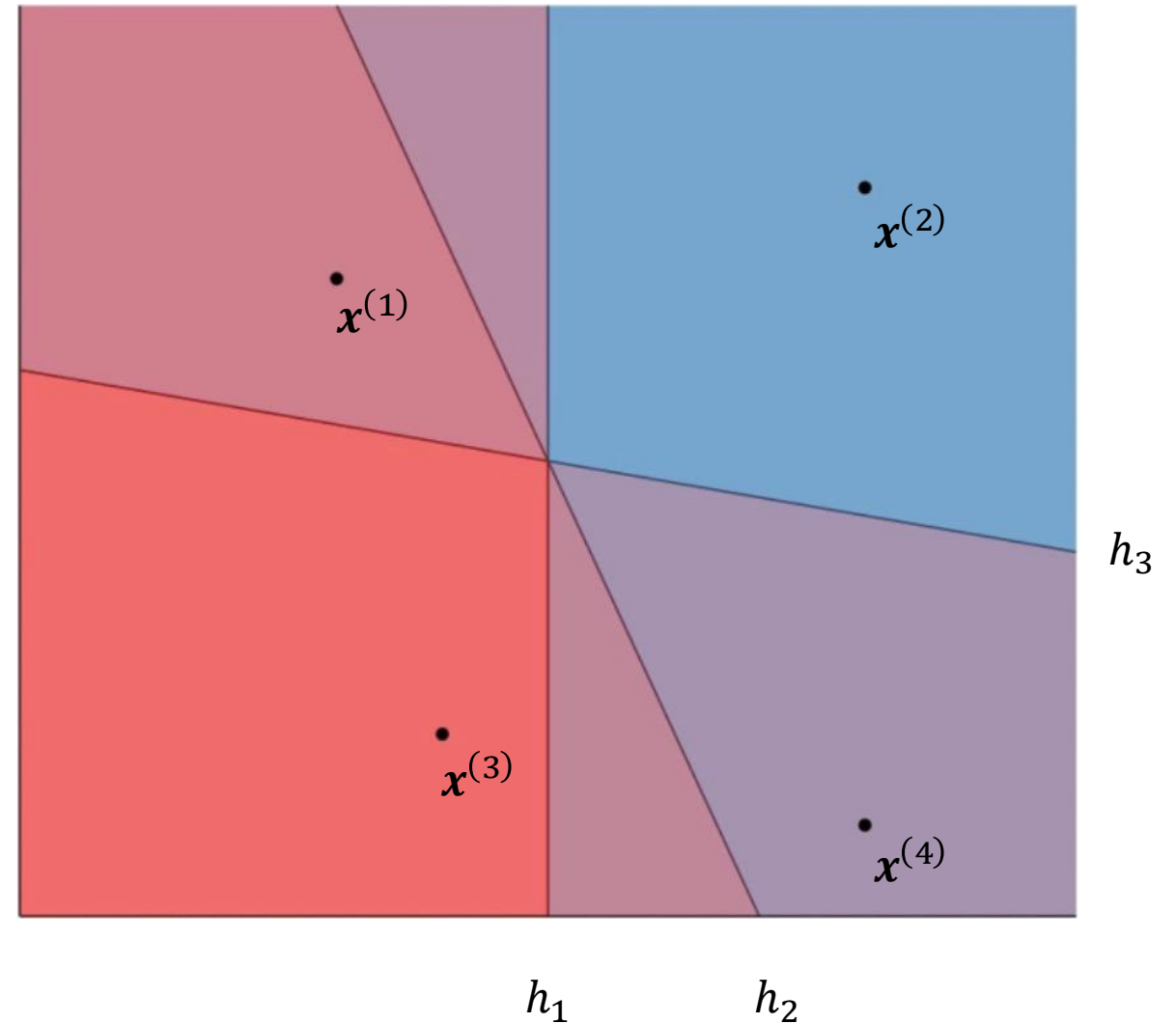


Example: Labellings

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

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

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

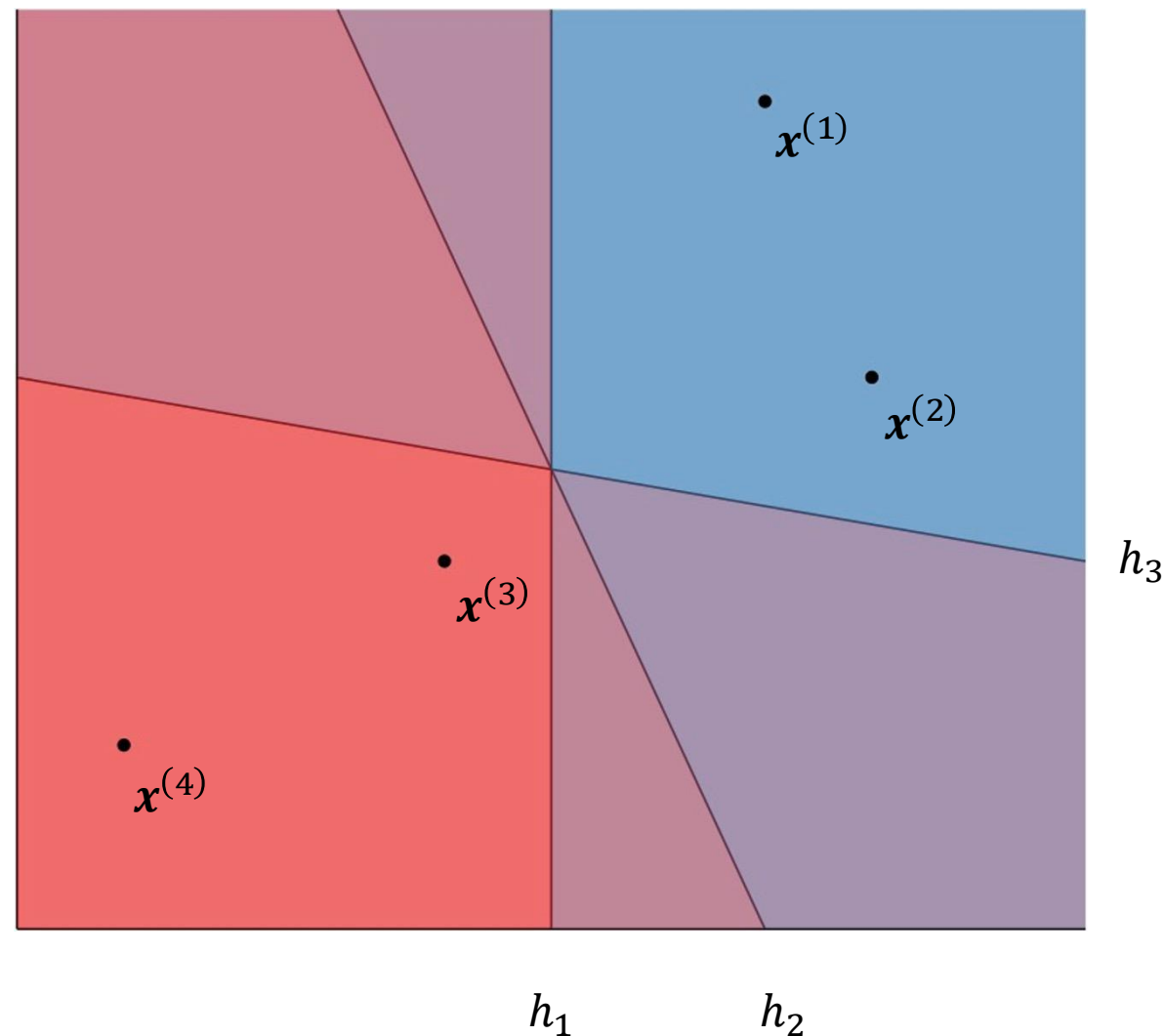


Example: Labellings

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

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

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



Growth Function

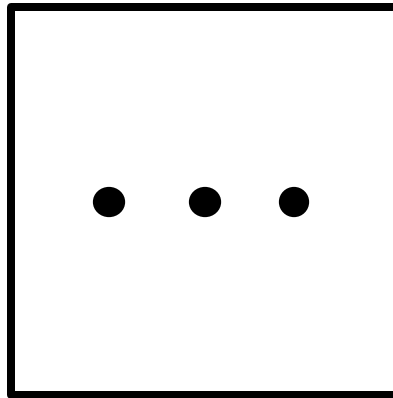
- The growth function of \mathcal{H} is the maximum number of distinct labellings \mathcal{H} can induce on **any** set of M data points:

$$g_{\mathcal{H}}(M) = \max_{S: |S|=M} |L_{\mathcal{H}}(S)|$$

- $g_{\mathcal{H}}(M) \leq 2^M \forall \mathcal{H}$ and M
- \mathcal{H} shatters S if $|L_{\mathcal{H}}(S)| = 2^M$
- If $\exists S$ s.t. $|S| = M$ and \mathcal{H} shatters S , then $g_{\mathcal{H}}(M) = 2^M$

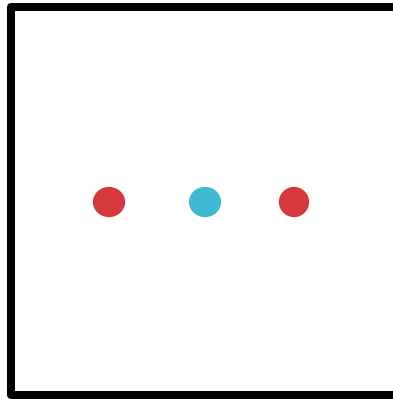
Growth Function: Example

- $\mathbf{x}^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $g_{\mathcal{H}}(3)$?



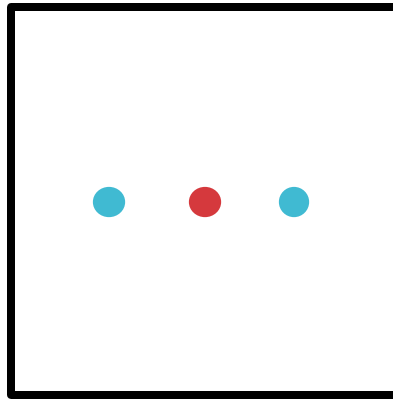
Growth Function: Example

- $\mathbf{x}^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $g_{\mathcal{H}}(3)$?



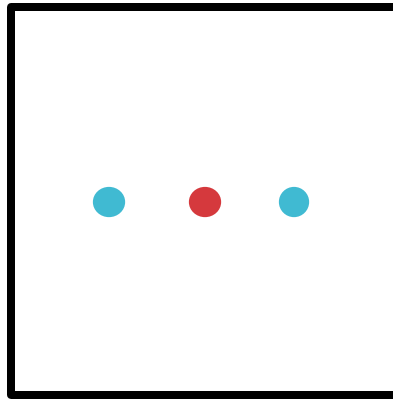
Growth Function: Example

- $\mathbf{x}^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $g_{\mathcal{H}}(3)$?

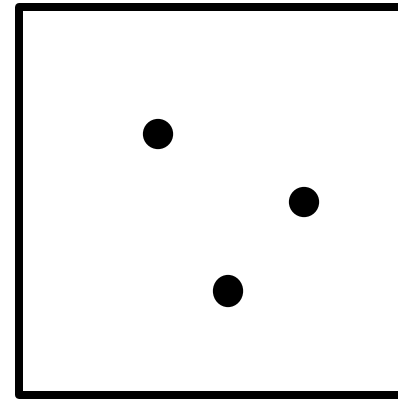


Growth Function: Example

- $\mathbf{x}^{(m)} \in \mathbb{R}^2$ and \mathcal{H} = all 2-dimensional linear separators
- What is $g_{\mathcal{H}}(3)$?



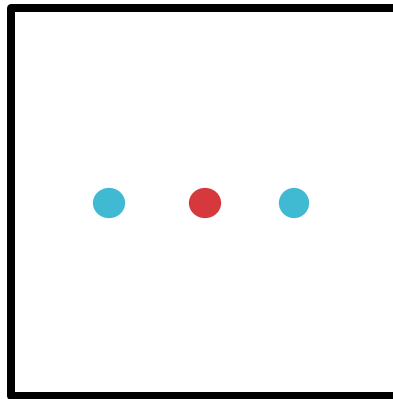
$$|\mathcal{H}(S_1)| = 6$$



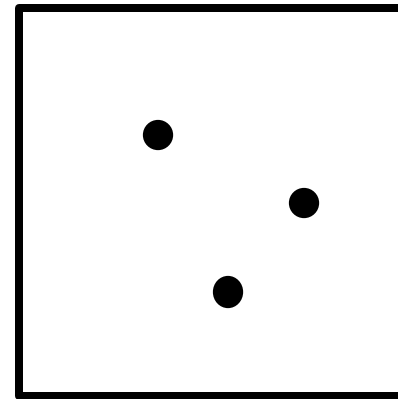
$$|\mathcal{H}(S_2)| = 8$$

Growth Function: Example

- $\mathbf{x}^{(m)} \in \mathbb{R}^2$ and \mathcal{H} = all 2-dimensional linear separators
- $g_{\mathcal{H}}(3) = 8 = 2^3$



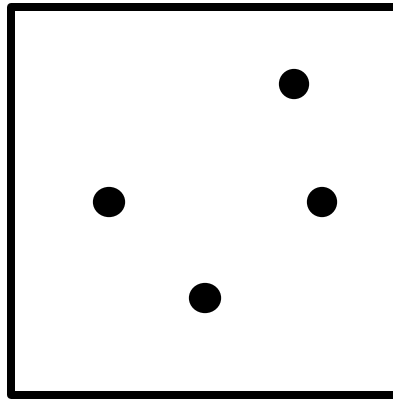
$$|\mathcal{H}(S_1)| = 6$$



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

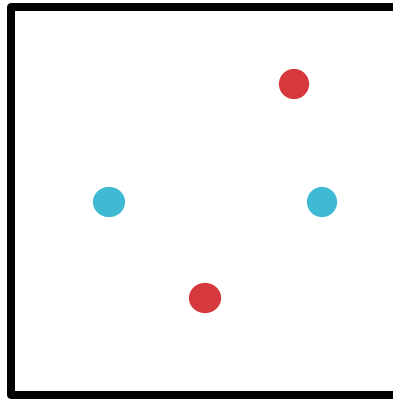
Growth Function: Example

- $\mathbf{x}^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $g_{\mathcal{H}}(4)$?



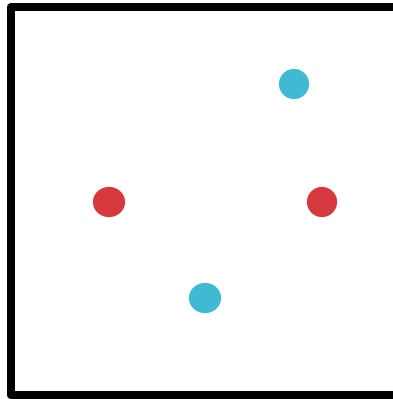
Growth Function: Example

- $\mathbf{x}^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $g_{\mathcal{H}}(4)$?



Growth Function: Example

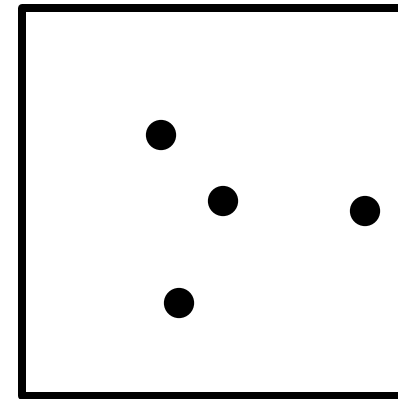
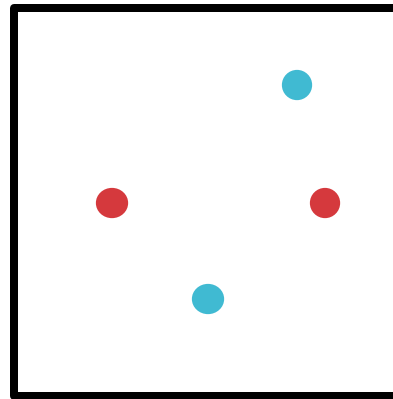
- $\mathbf{x}^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $g_{\mathcal{H}}(4)$?



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

Growth Function: Example

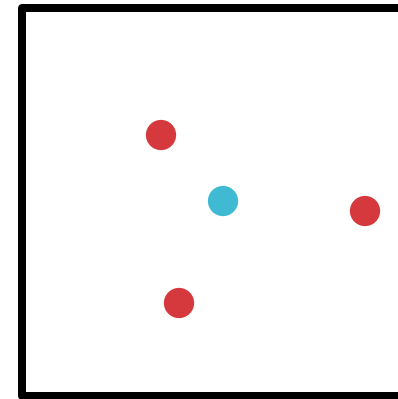
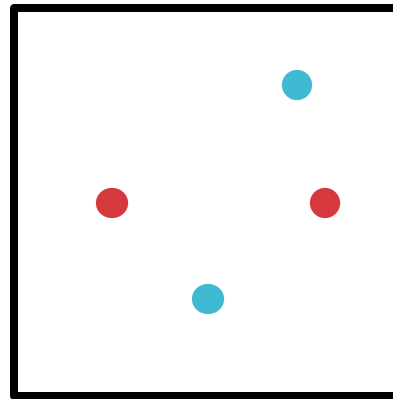
- $\mathbf{x}^{(m)} \in \mathbb{R}^2$ and \mathcal{H} = all 2-dimensional linear separators
- What is $g_{\mathcal{H}}(4)$?



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

Growth Function: Example

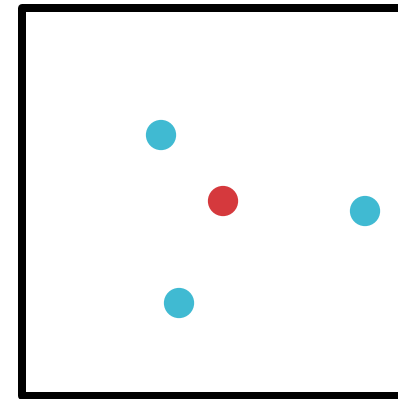
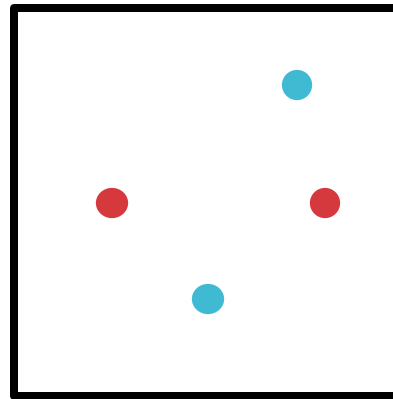
- $\mathbf{x}^{(m)} \in \mathbb{R}^2$ and \mathcal{H} = all 2-dimensional linear separators
- What is $g_{\mathcal{H}}(4)$?



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

Growth Function: Example

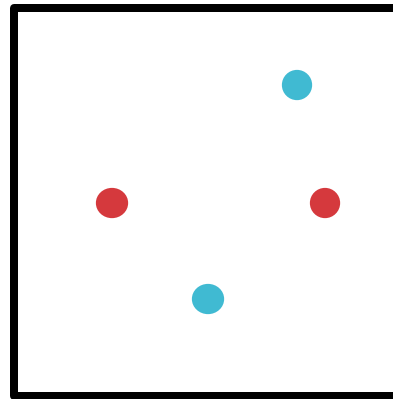
- $\mathbf{x}^{(m)} \in \mathbb{R}^2$ and \mathcal{H} = all 2-dimensional linear separators
- What is $g_{\mathcal{H}}(4)$?



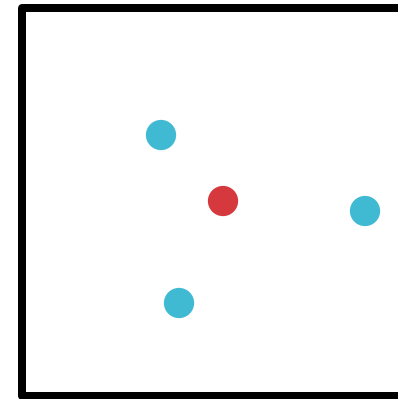
$$|\mathcal{H}(S_1)| = 14$$

Growth Function: Example

- $\mathbf{x}^{(m)} \in \mathbb{R}^2$ and \mathcal{H} = all 2-dimensional linear separators
- $g_{\mathcal{H}}(4) = 14 < 2^4$



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



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

Growth Function: Example

- $\mathbf{x}^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $g_{\mathcal{H}}(5)$?

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

$$M \geq \frac{2}{\epsilon} \left(\log_2(2g_{\mathcal{H}}(2M)) + \log_2\left(\frac{1}{\delta}\right) \right)$$

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

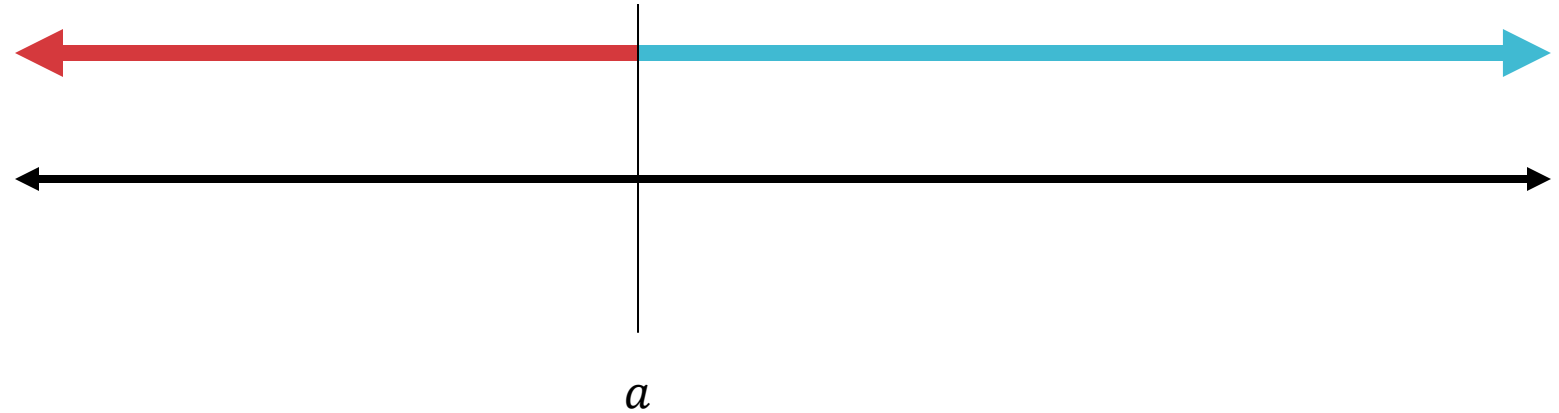
M appears on both sides of the inequality...

Theorem 3: Vapnik- Chervonenkis (VC)-Dimension

- $d_{VC}(\mathcal{H})$ = the largest value of M s.t. $g_{\mathcal{H}}(M) = 2^M$, i.e., the greatest number of data points that can be shattered by \mathcal{H}
 - If \mathcal{H} can shatter arbitrarily large finite sets, then $d_{VC}(\mathcal{H}) = \infty$
 - $g_{\mathcal{H}}(M) = O(M^{d_{VC}(\mathcal{H})})$ (Sauer-Shelah lemma)
- To prove that $d_{VC}(\mathcal{H}) = C$, you need to show
 1. \exists some set of C data points that \mathcal{H} can shatter and
 2. \nexists a set of $C + 1$ data points that \mathcal{H} can shatter

VC-Dimension: Example

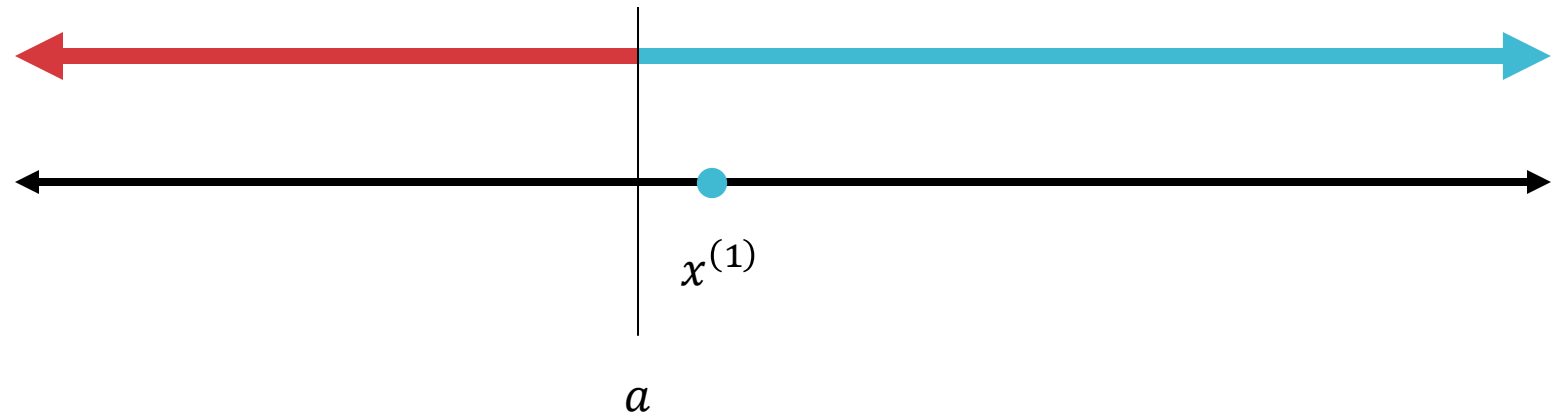
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

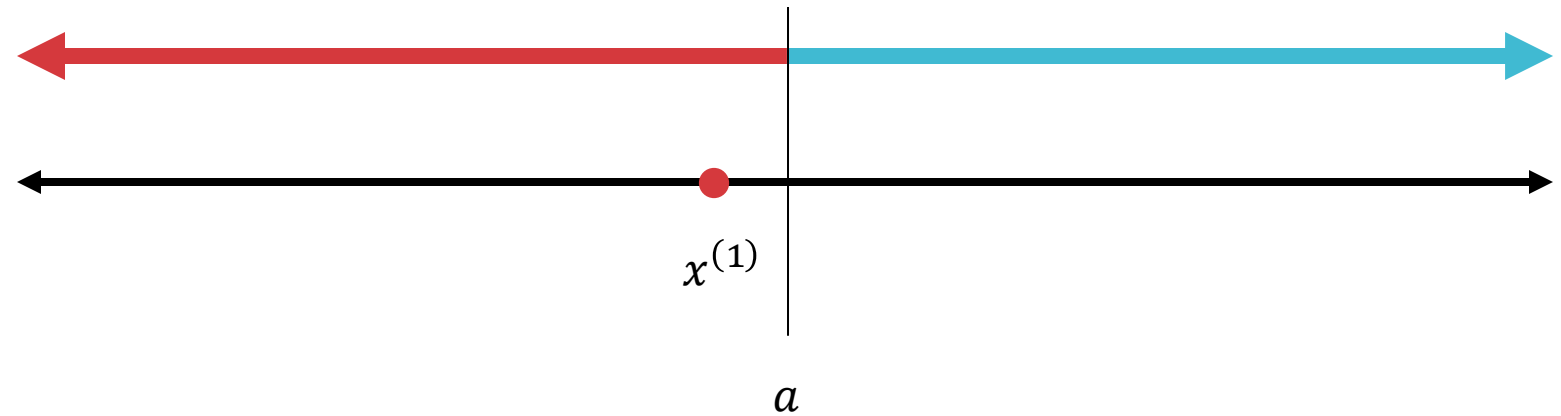
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

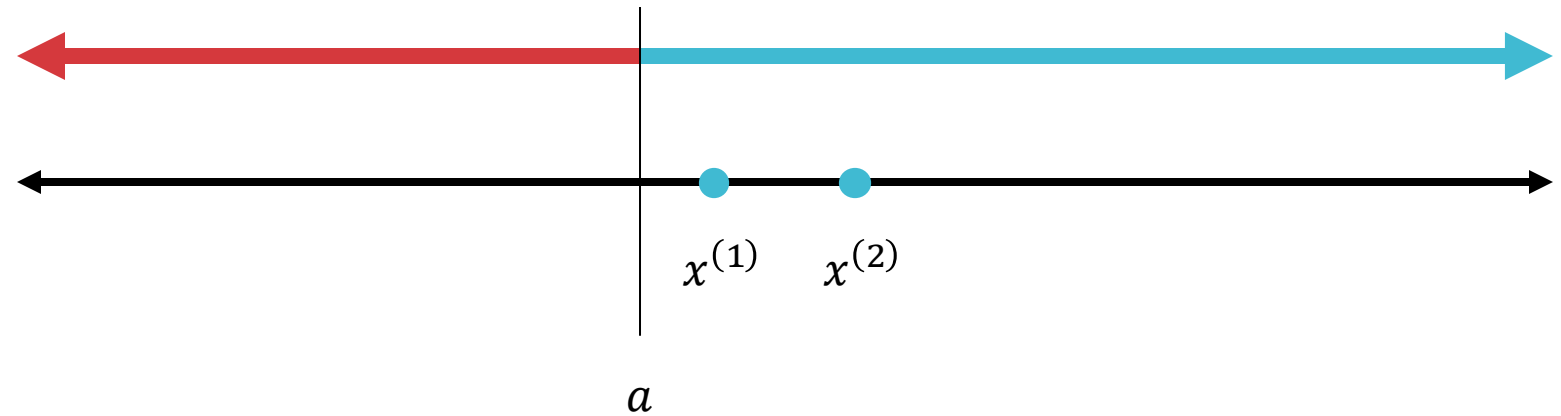
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

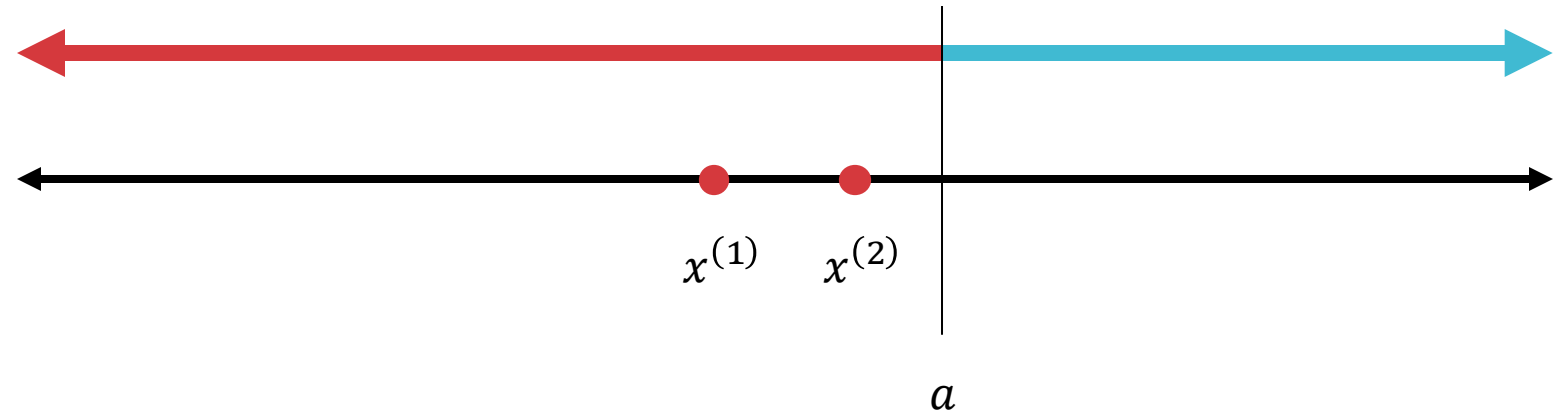
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

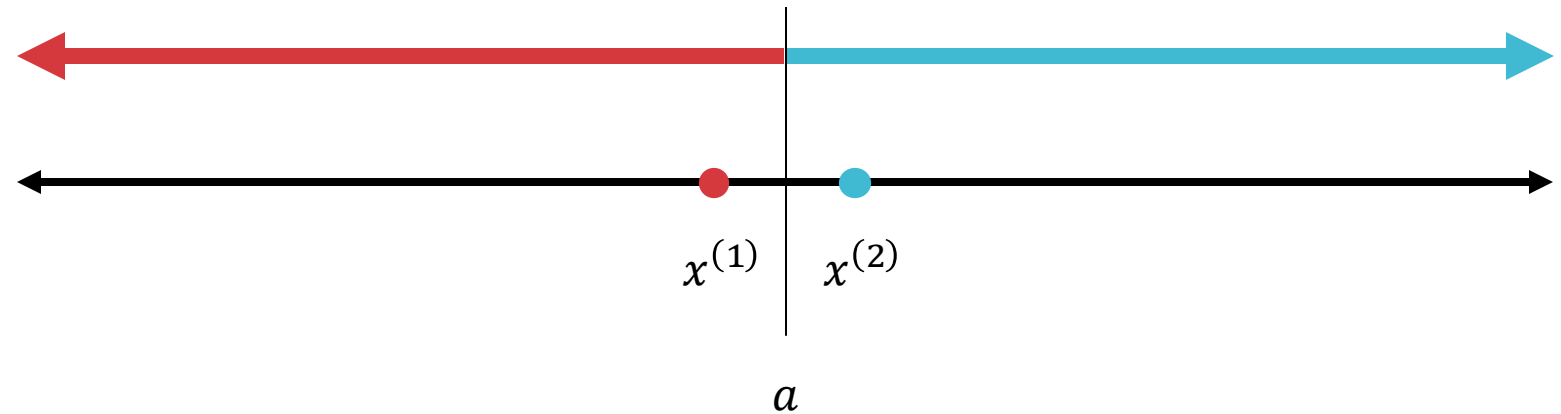
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

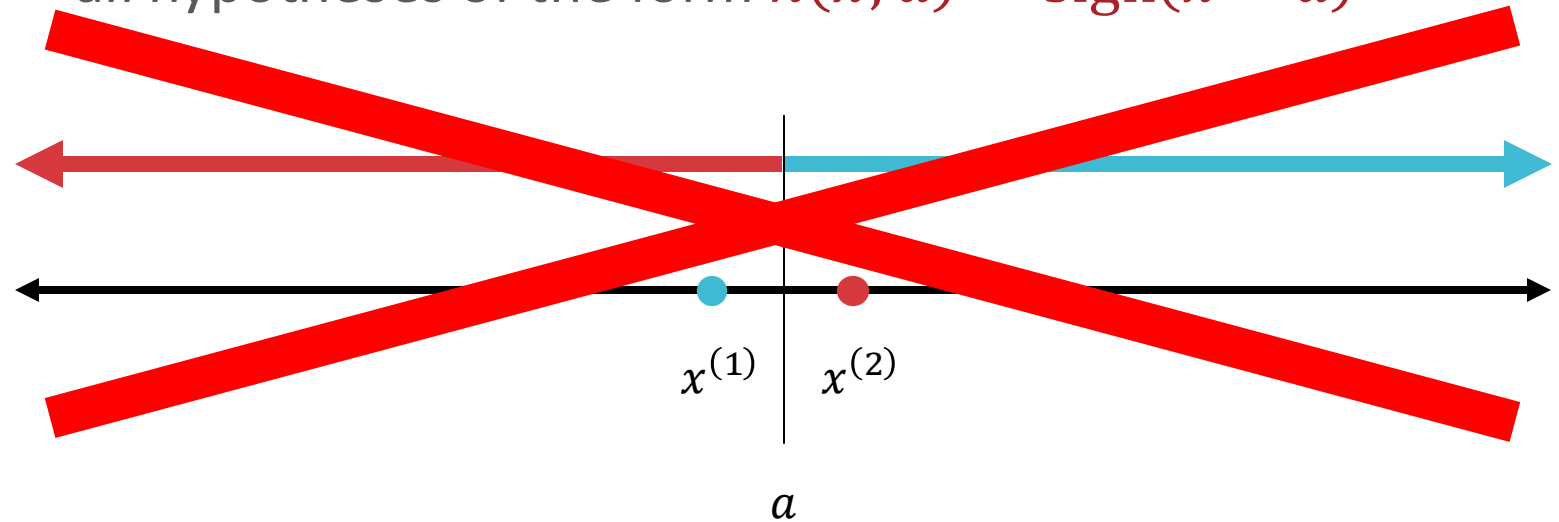
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

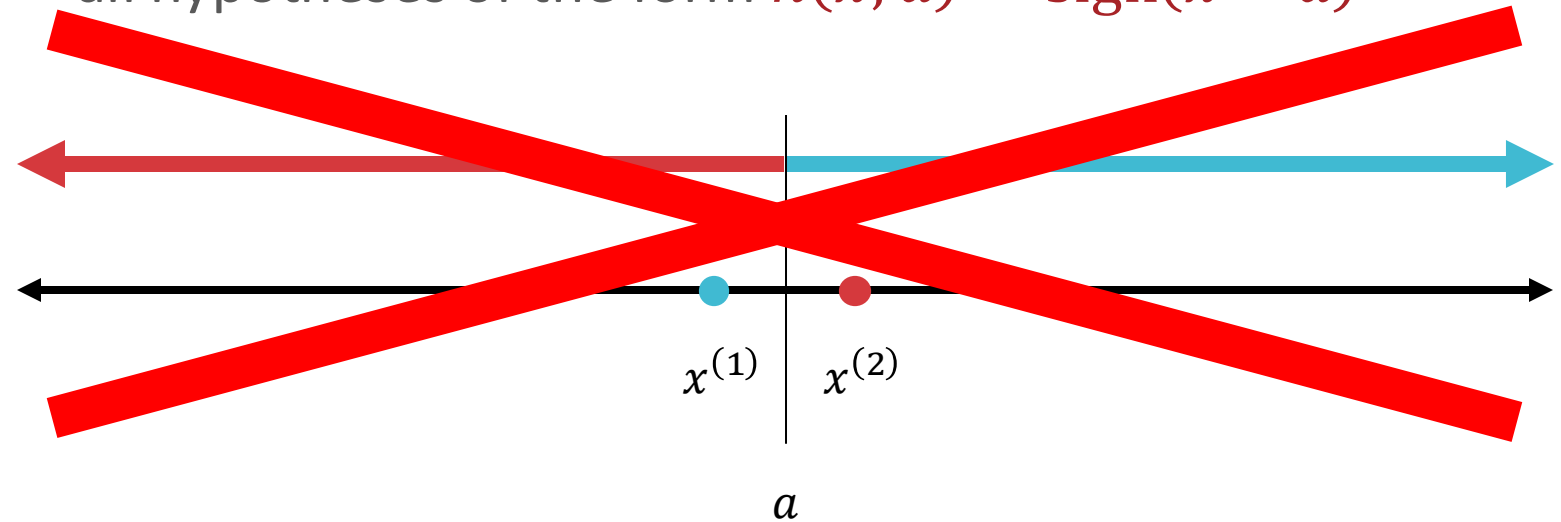
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

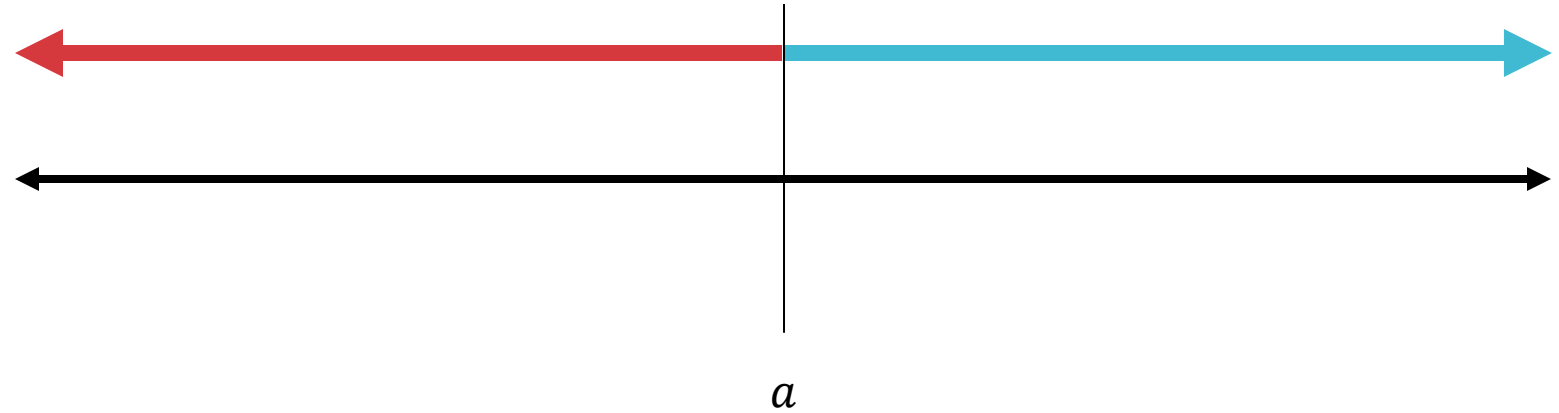
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

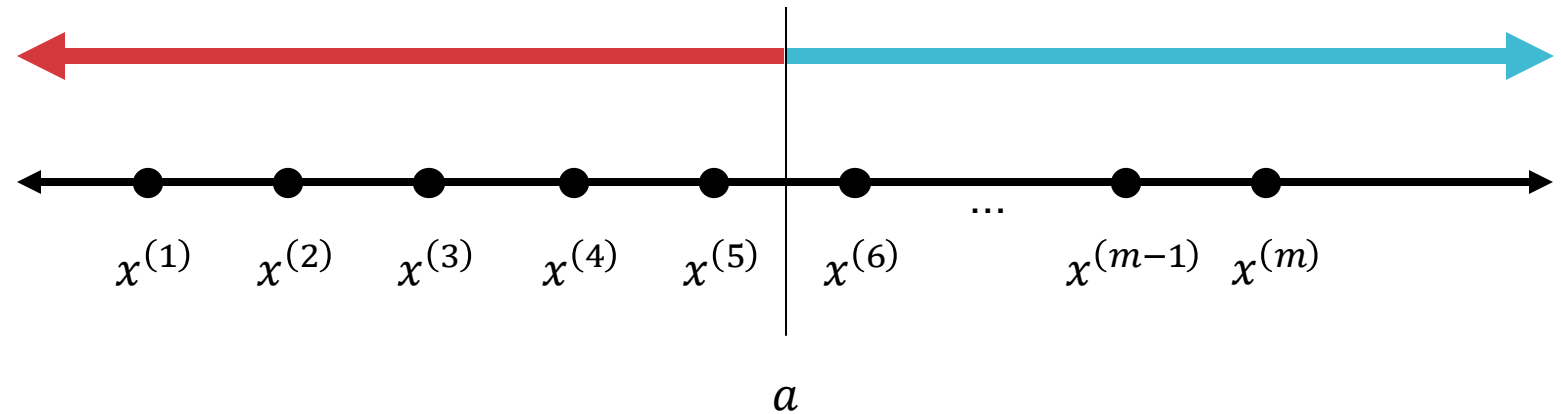
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

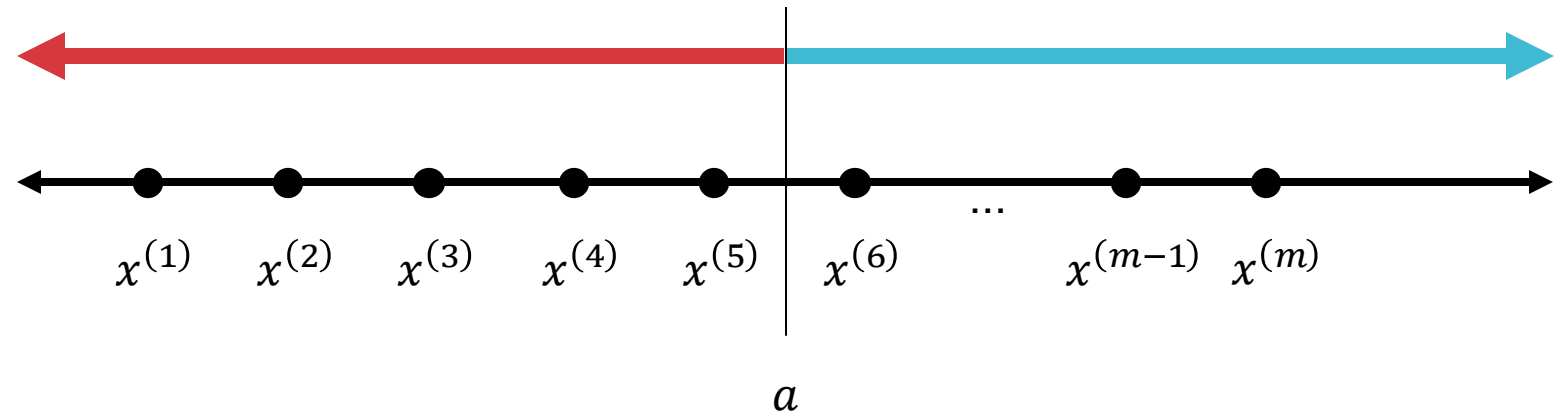
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

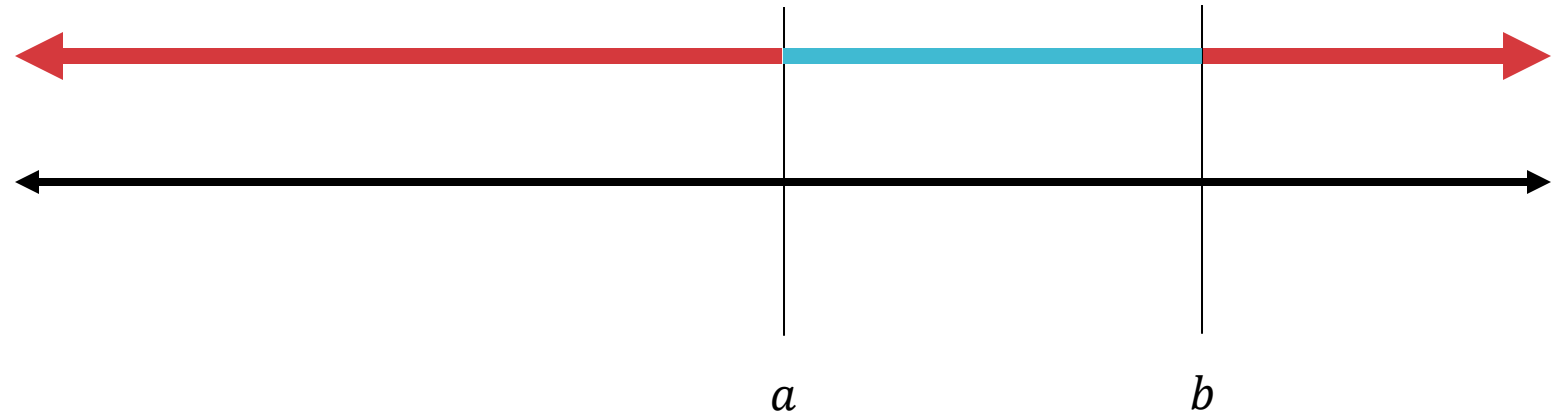
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



- $g_{\mathcal{H}}(m) = m + 1 = O(m^1)$

VC-Dimension: In-class Poll

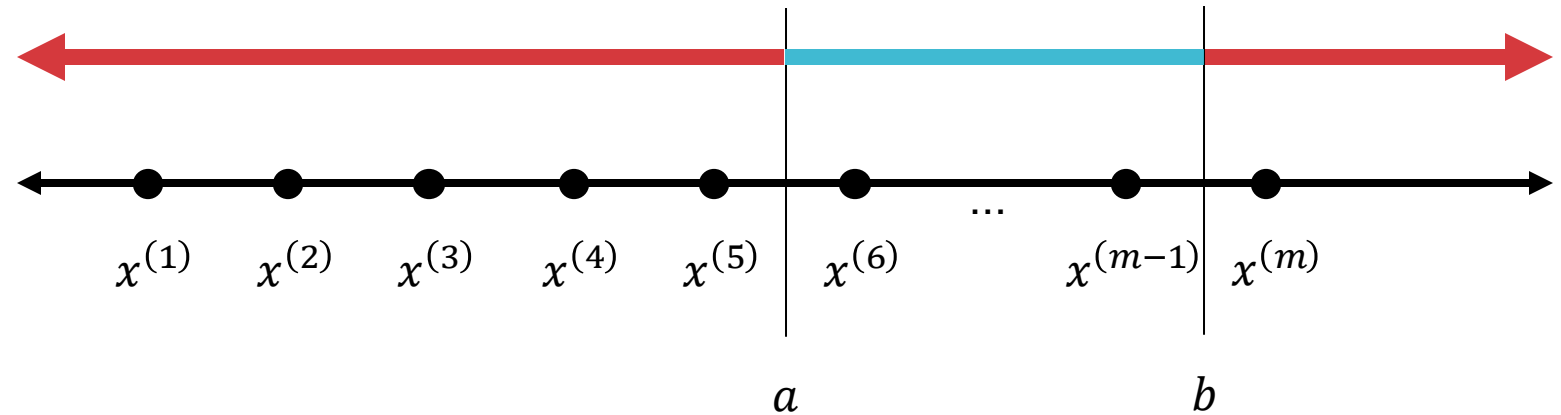
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive intervals



- What are $d_{VC}(\mathcal{H})$ and $g_{\mathcal{H}}(m)$?
 - 1 and $m+1$
 - 2 and $m+1$
 - 2 and $\frac{1}{2}(m^2 + m + 4)$
 - 2 and $\frac{1}{2}(m^2 + m + 2)$

VC-Dimension: Example

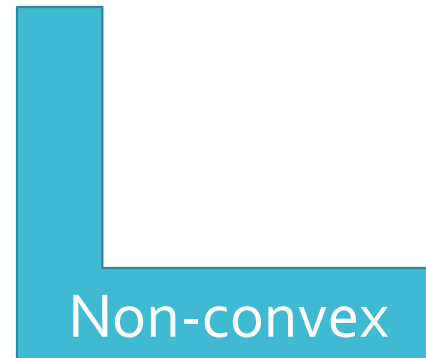
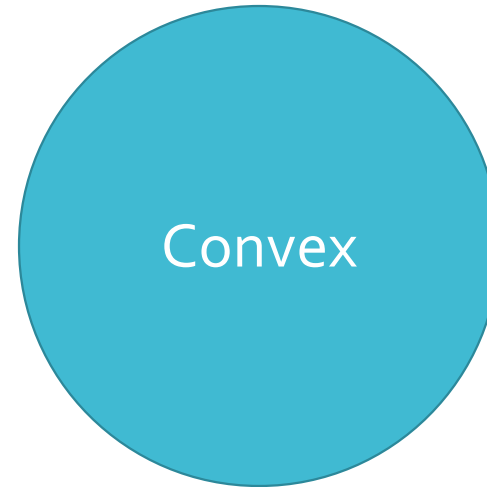
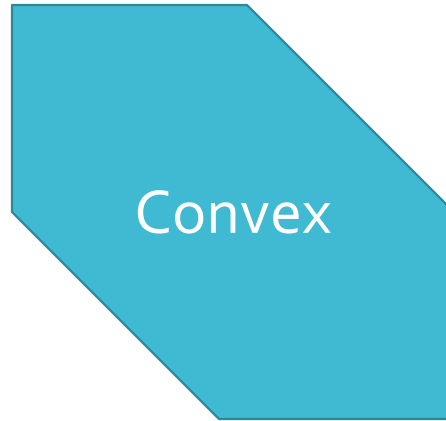
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive intervals



- What are $d_{VC}(\mathcal{H})$ and $g_{\mathcal{H}}(m)$?

Growth Function: Example

- $x^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional positive convex sets

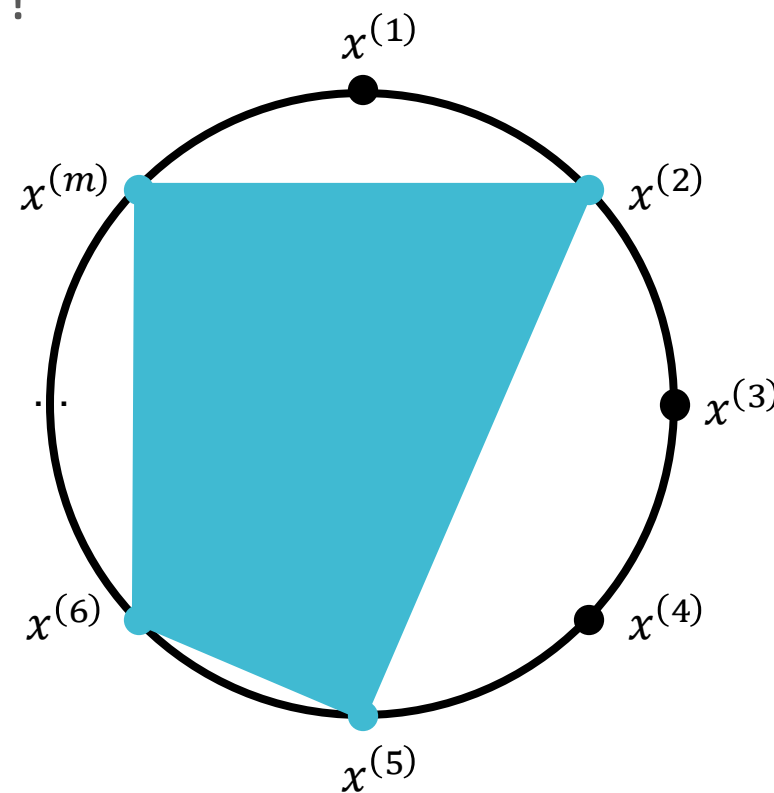


Growth Function: Example

- $x^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional positive convex sets
- What are $d_{VC}(\mathcal{H})$ and $g_{\mathcal{H}}(M)$?

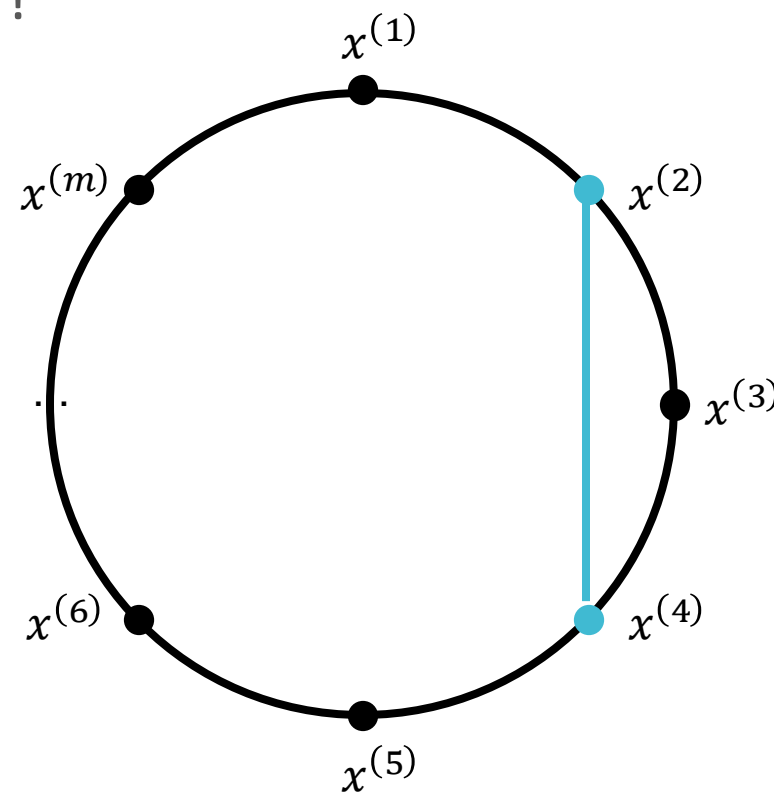
Growth Function: Example

- $x^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional positive convex sets
- What are $d_{VC}(\mathcal{H})$ and $g_{\mathcal{H}}(M)$?



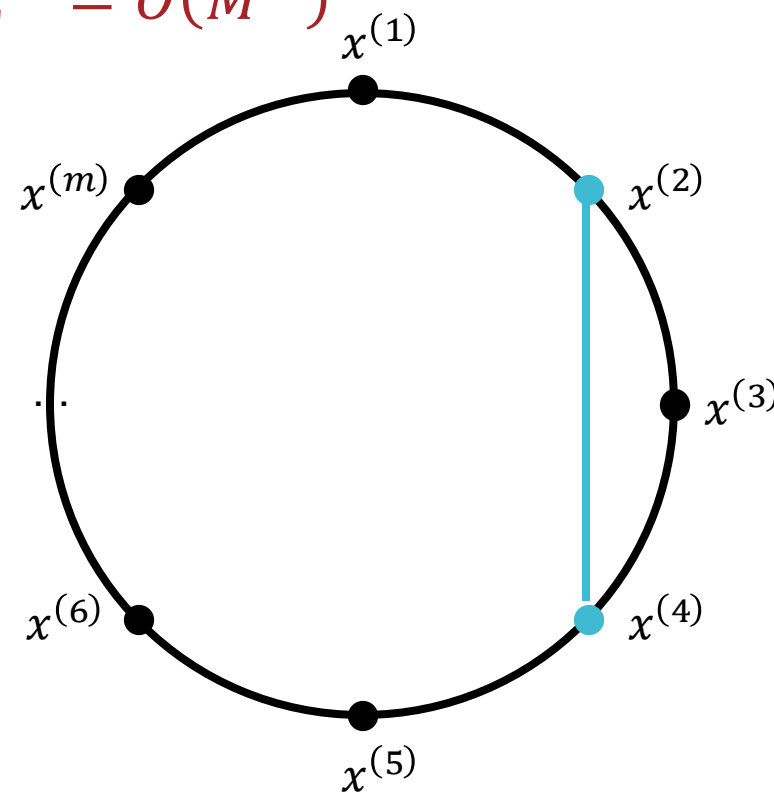
Growth Function: Example

- $x^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional positive convex sets
- What are $d_{VC}(\mathcal{H})$ and $g_{\mathcal{H}}(M)$?



Growth Function: Example

- $x^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional positive convex sets
- $d_{VC}(\mathcal{H}) = \infty$ and $g_{\mathcal{H}}(M) = 2^M = O(M^\infty)$



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

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

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

- 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

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

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

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

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

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

Statistical Learning Theory Corollary

- 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

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

with probability at least $1 - \delta$.

Approximation Generalization Tradeoff

How well does
 h generalize?

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

How well does h
fit the data?

Approximation Generalization Tradeoff

$$R(h) \leq \underbrace{\hat{R}(h)}_{\text{Decreases as } d_{VC}(\mathcal{H}) \text{ increases}} + O\left(\underbrace{\sqrt{\frac{1}{M} \left(d_{VC}(\mathcal{H}) + \log\left(\frac{1}{\delta}\right) \right)}}_{\text{Increases as } d_{VC}(\mathcal{H}) \text{ increases}}\right)$$

Key Takeaways

- For infinite hypothesis sets, use the VC-dimension (or the growth function) as a measure of complexity
 - Computing $d_{VC}(\mathcal{H})$ and $g_{\mathcal{H}}(M)$
 - Connection between VC-dimension and the growth function (Sauer-Shelah lemma)
 - Sample complexity and statistical learning theory style bounds using $d_{VC}(\mathcal{H})$