

10-701: Introduction to Machine Learning Lecture 21 – Learning Theory (Infinite Case)

Henry Chai

4/3/24

Front Matter

- Announcements
 - Project check-ins due on 4/8 at 11:59 PM
 - **Daniel is on leave and will be for an indeterminate amount of time, please direct all course requests/questions to Henry**

Key Question

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

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...

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

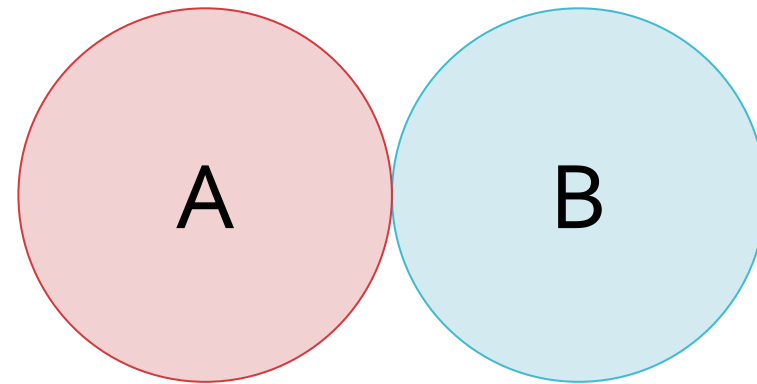
$$R(h) \leq \hat{R}(h) + \sqrt{\frac{1}{2M} \left(\ln(|\mathcal{H}|) + \ln\left(\frac{2}{\delta}\right) \right)}$$

with probability at least $1 - \delta$.

The Union Bound is Bad!

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

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

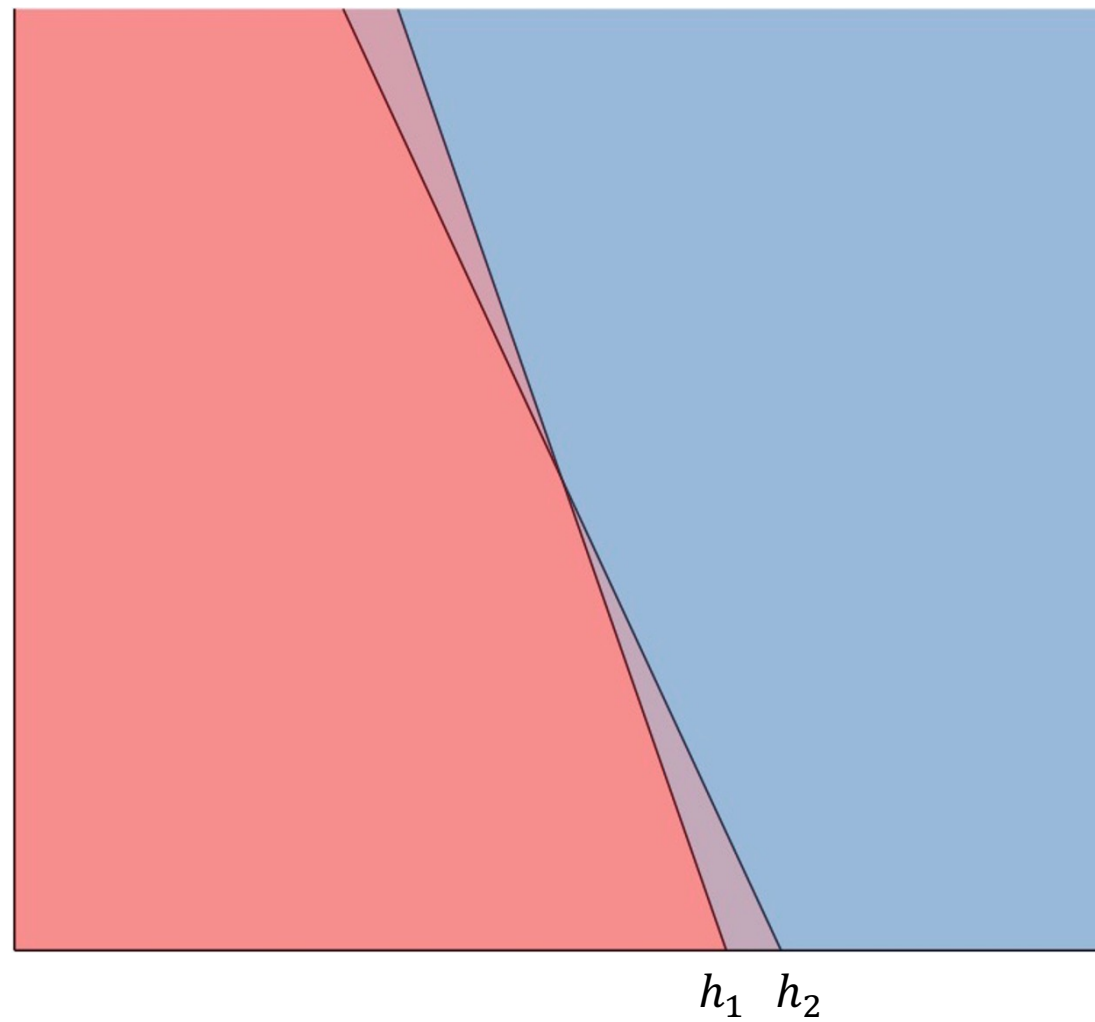


Intuition

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!

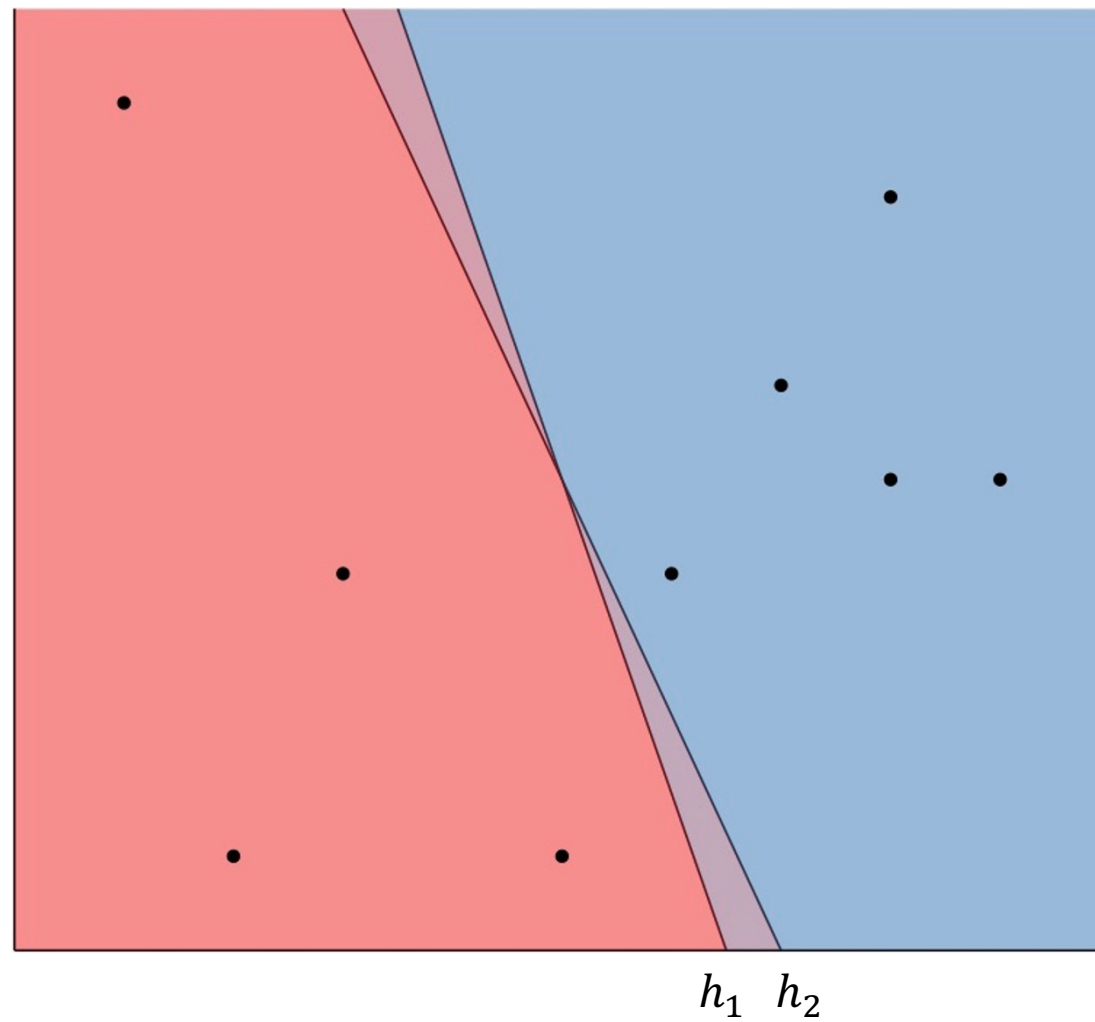


Intuition

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!

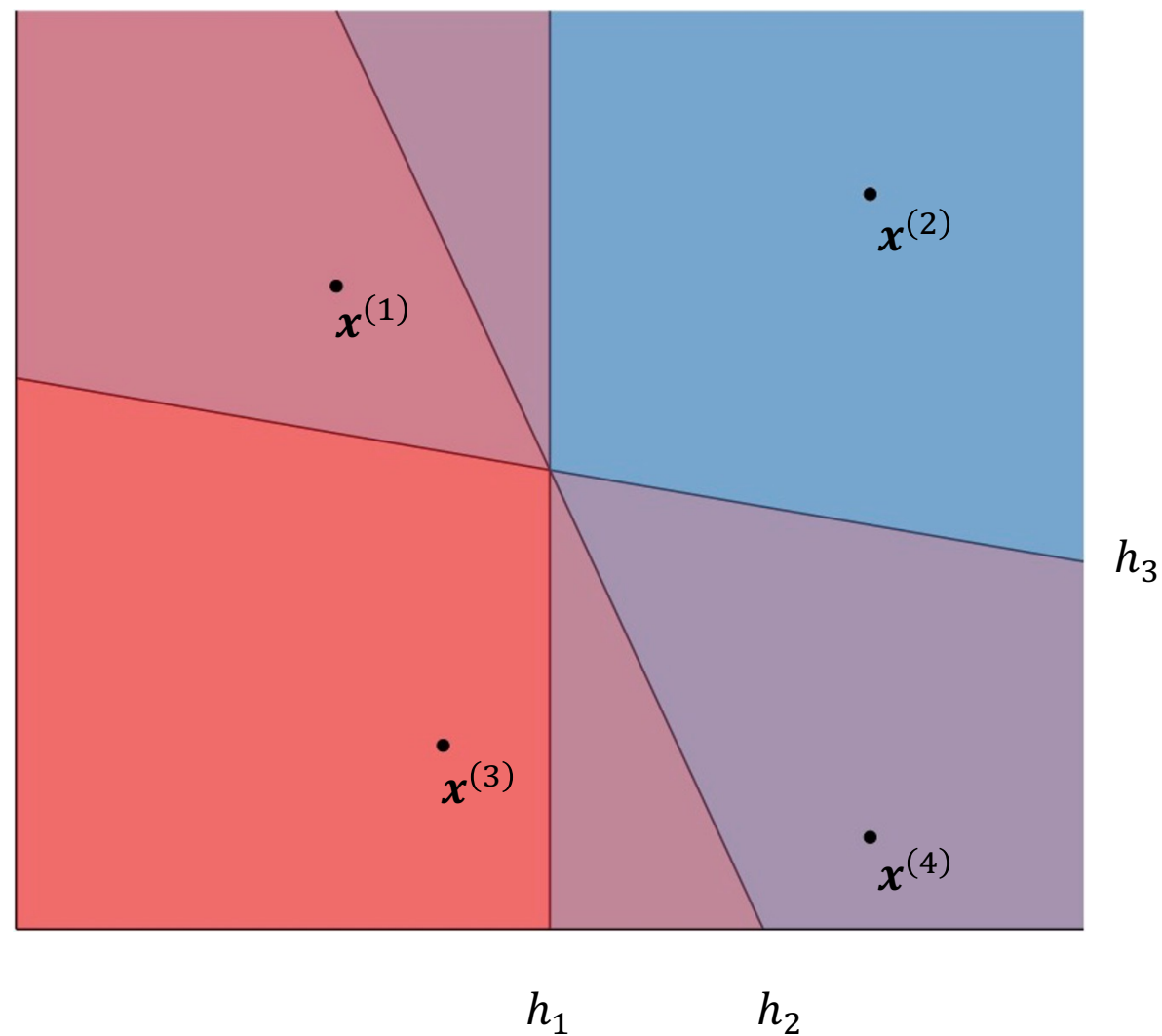


Labellings

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

Example: Labellings

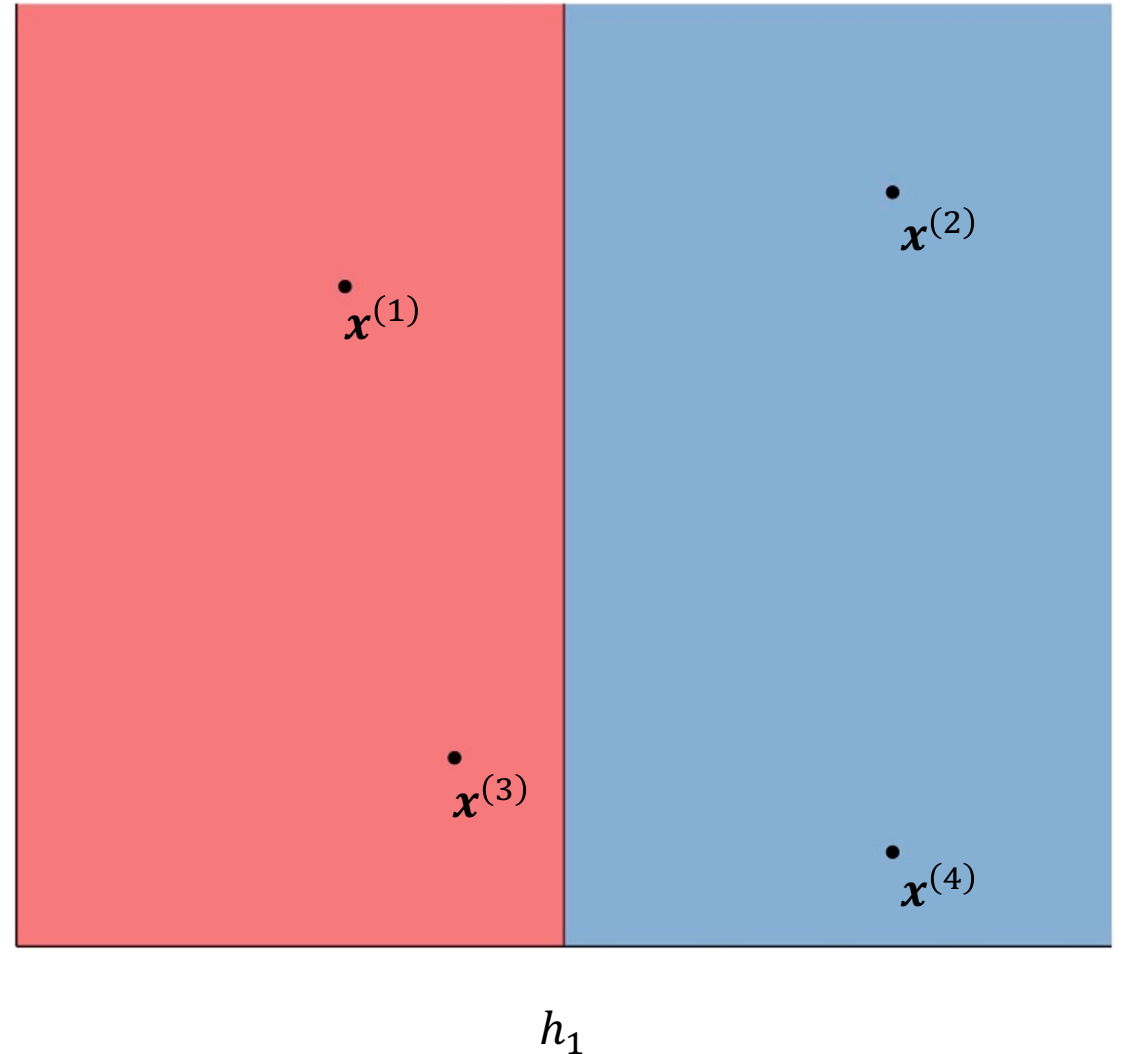
$$\mathcal{H} = \{h_1, h_2, h_3\}$$



Example: Labellings

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

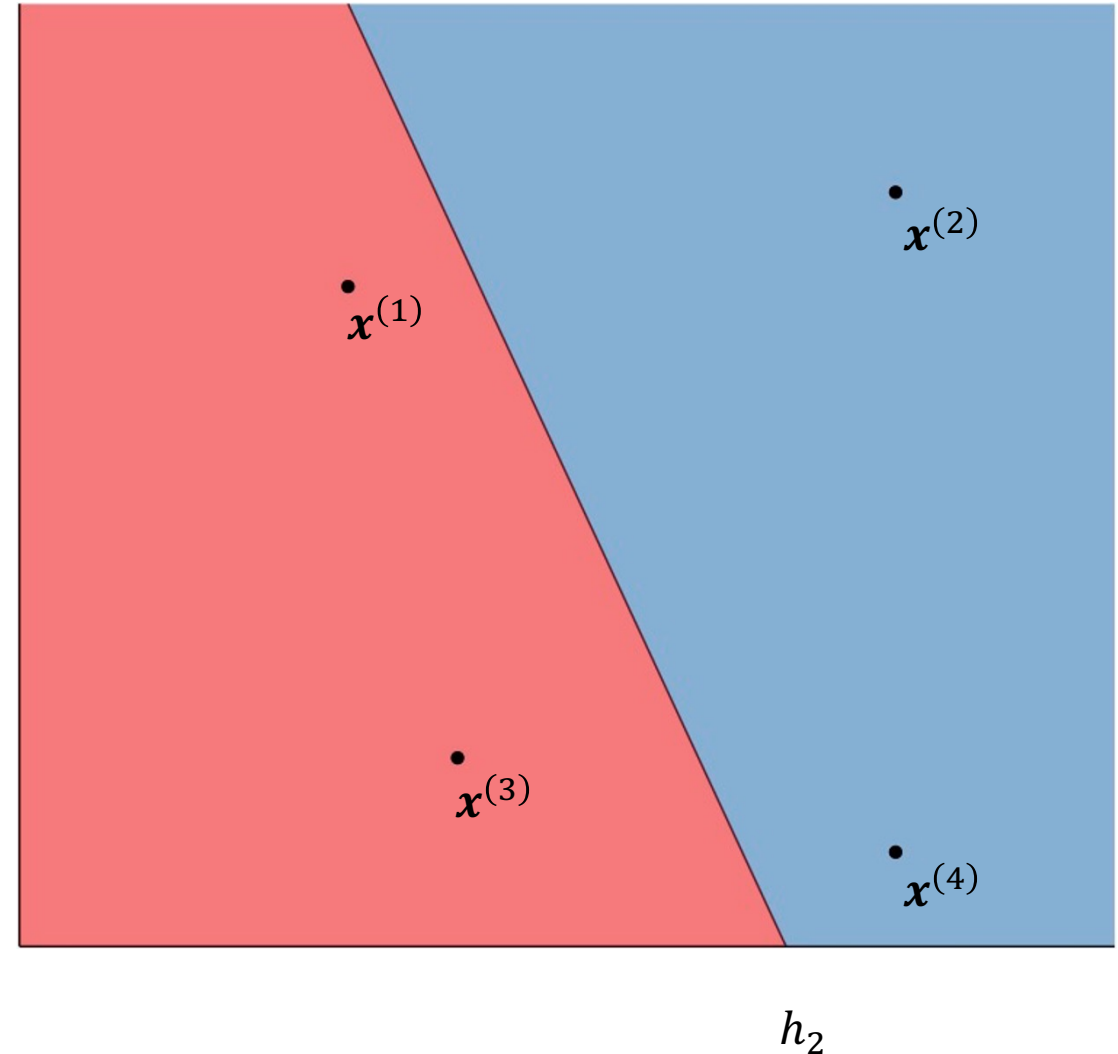
$$\begin{aligned} & (h_1(\mathbf{x}^{(1)}), h_1(\mathbf{x}^{(2)}), h_1(\mathbf{x}^{(3)}), h_1(\mathbf{x}^{(4)})) \\ &= (-1, +1, -1, +1) \end{aligned}$$



Example: Labellings

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

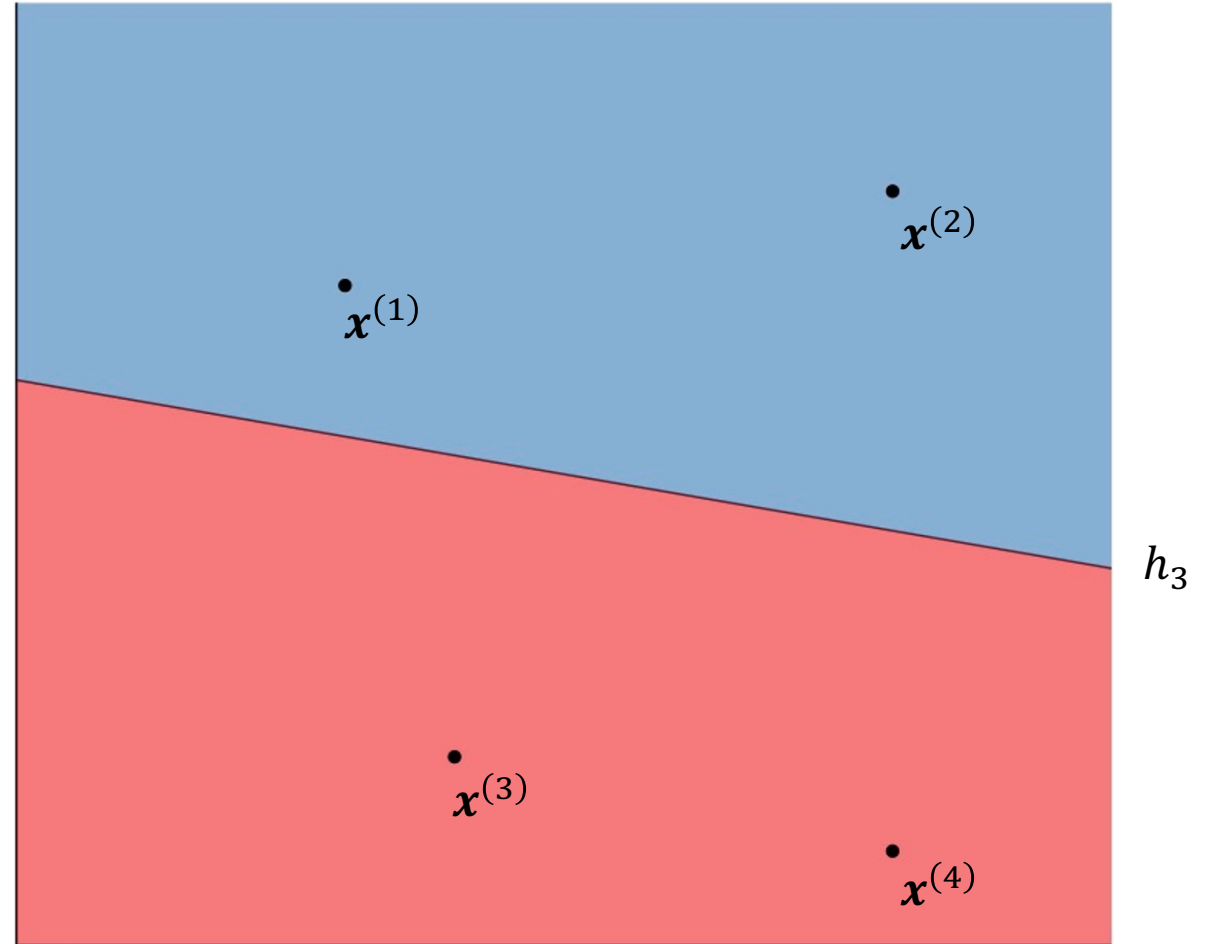
$$\begin{aligned} & \left(h_2(\mathbf{x}^{(1)}), h_2(\mathbf{x}^{(2)}), h_2(\mathbf{x}^{(3)}), h_2(\mathbf{x}^{(4)}) \right) \\ & = (-1, +1, -1, +1) \end{aligned}$$



Example: Labellings

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

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

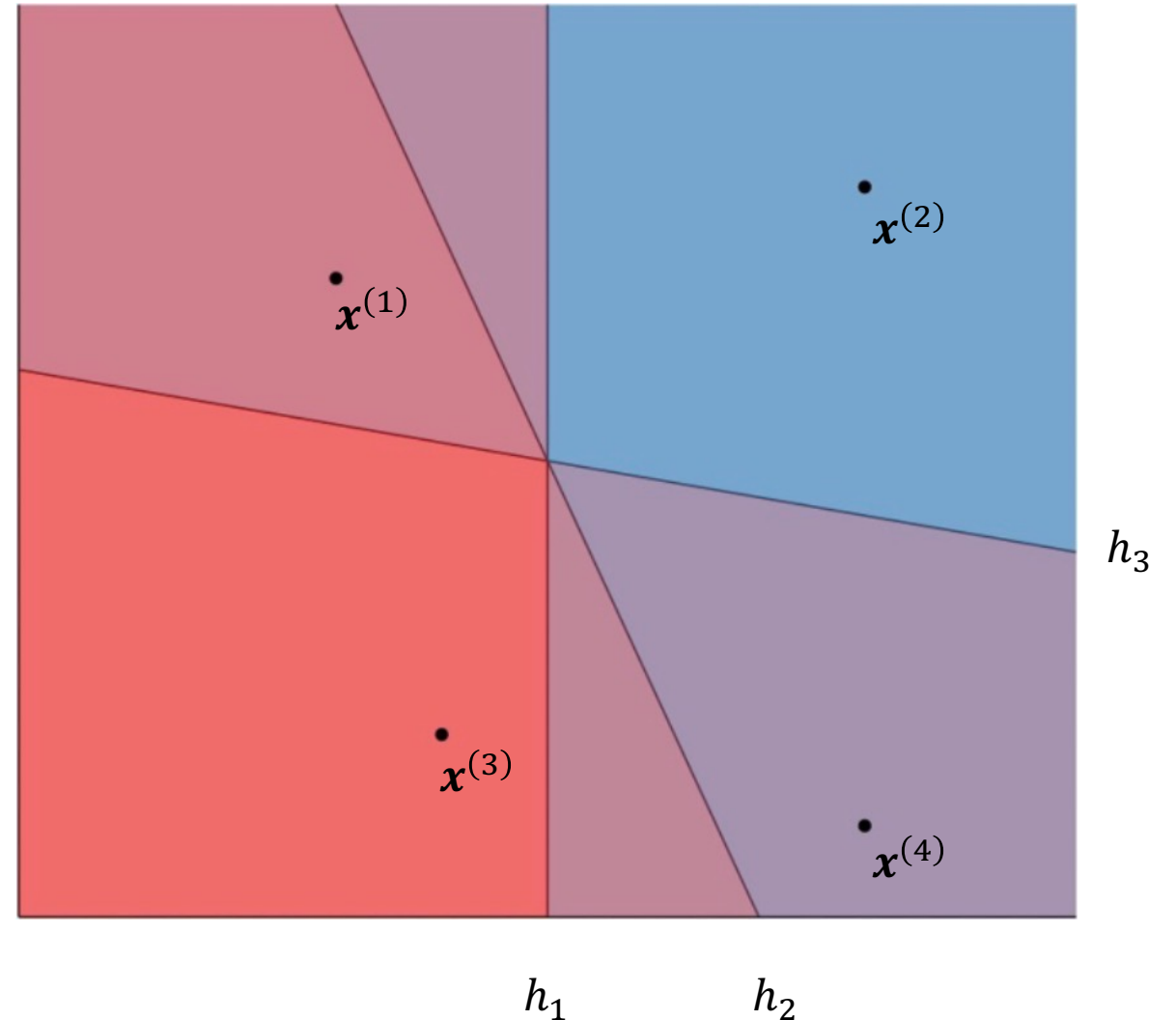


Example: Labellings

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

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

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

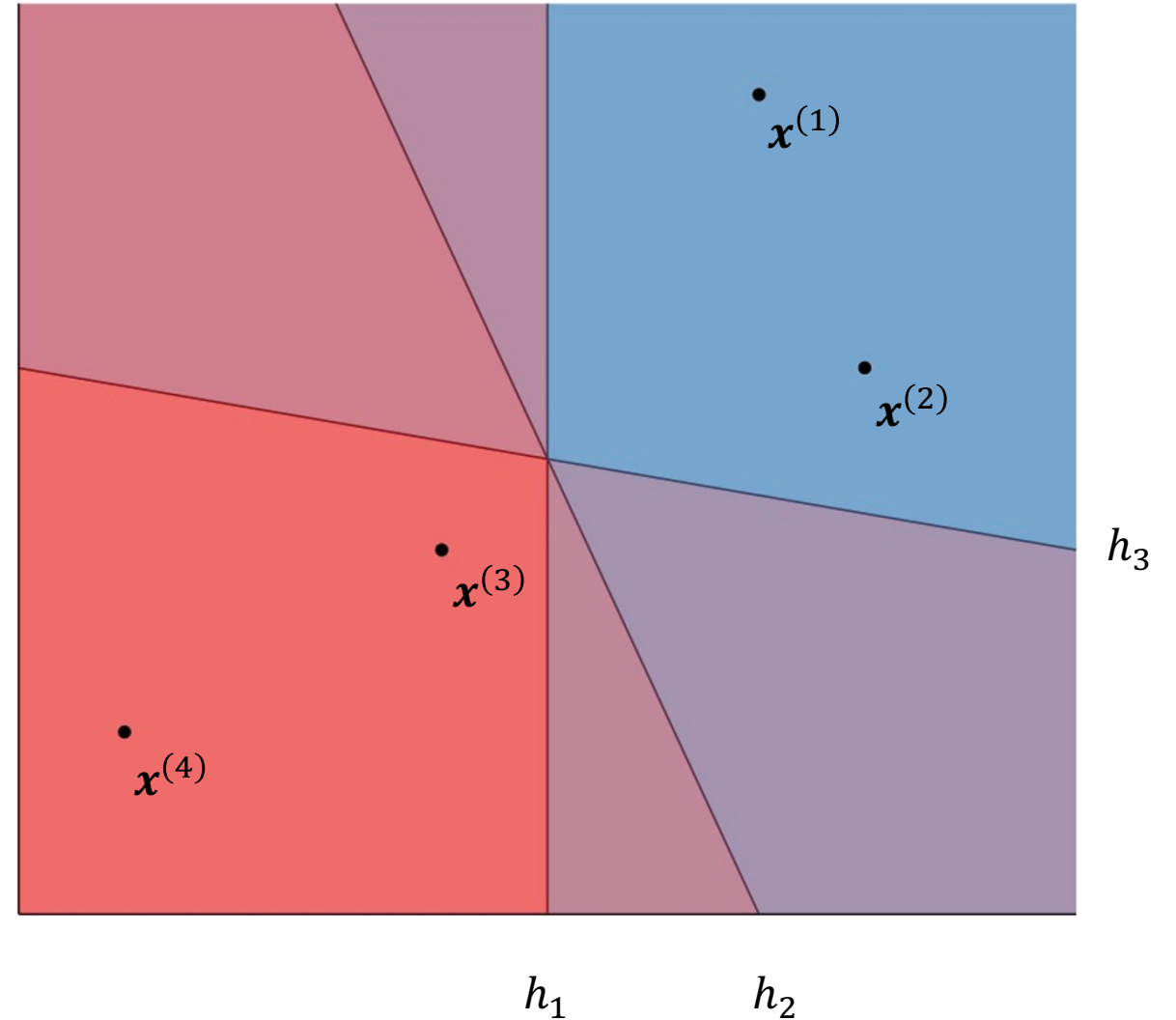


Example: Labellings

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

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

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



Growth Function

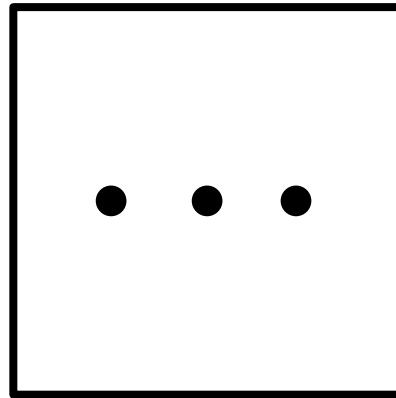
- The growth function of \mathcal{H} is the maximum number of distinct labellings \mathcal{H} can induce on **any** set of M data points:

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

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

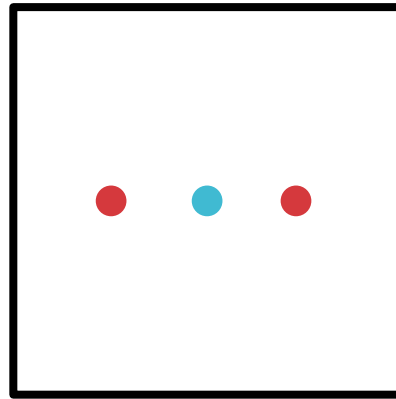
Growth Function: Example

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



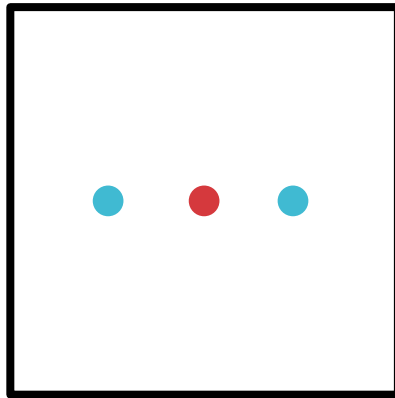
Growth Function: Example

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



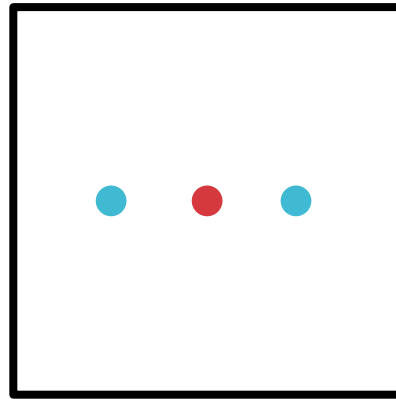
Growth Function: Example

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

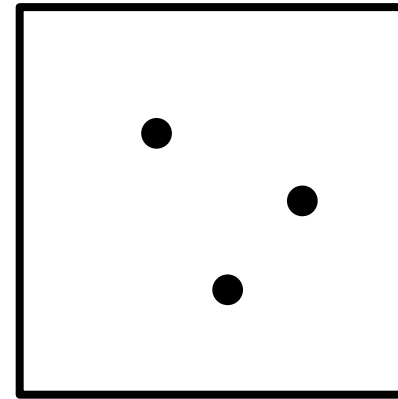


Growth Function: Example

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



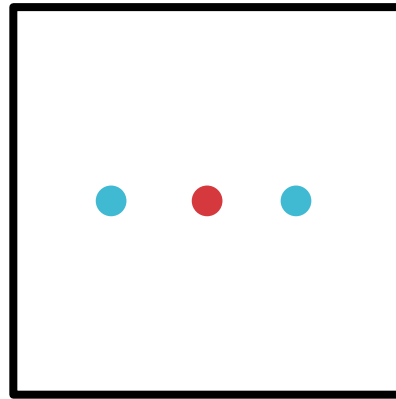
$$|\mathcal{H}(S_1)| = 6$$



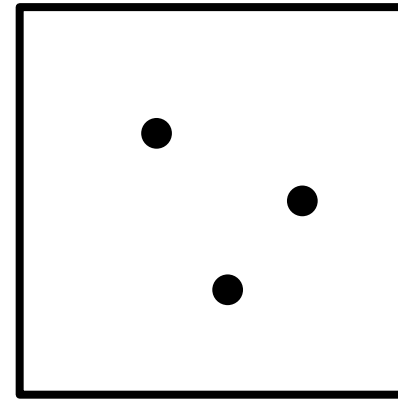
$$|\mathcal{H}(S_2)| = 8$$

Growth Function: Example

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



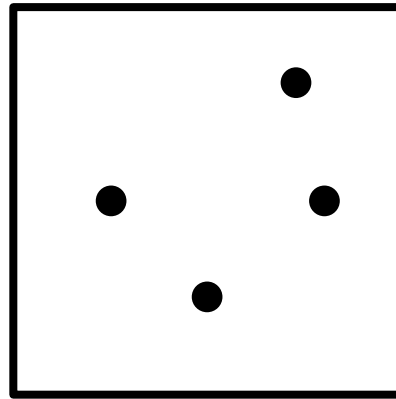
$$|\mathcal{H}(S_1)| = 6$$



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

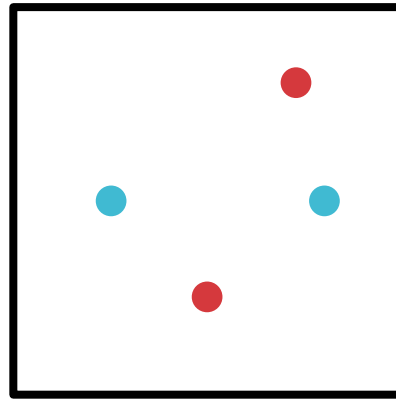
Growth Function: Example

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



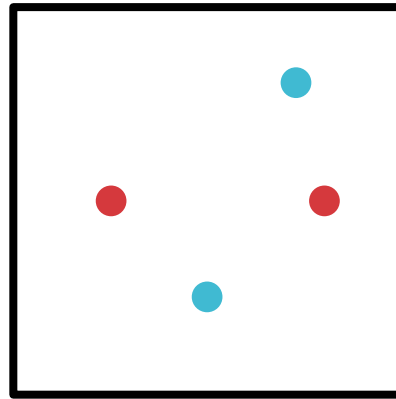
Growth Function: Example

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



Growth Function: Example

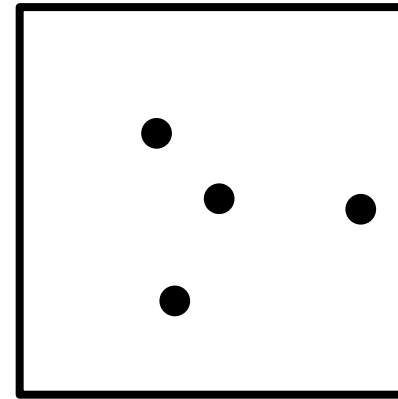
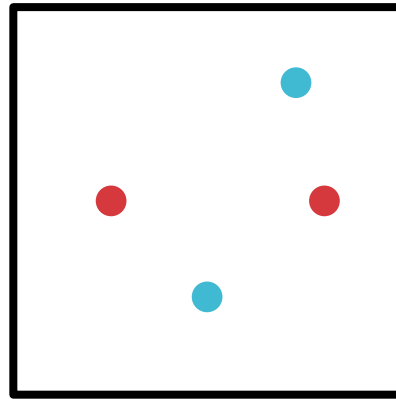
- $\mathbf{x}^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $g_{\mathcal{H}}(4)$?



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

Growth Function: Example

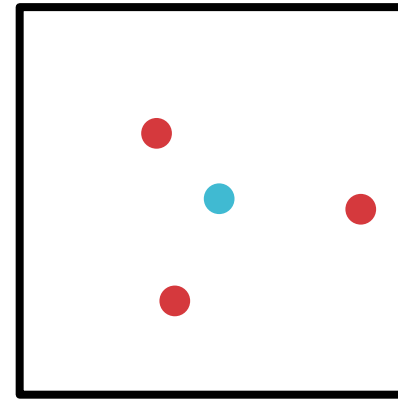
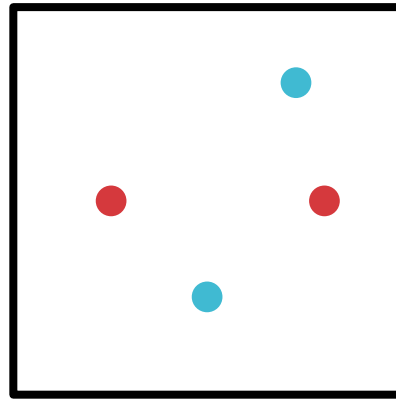
- $\mathbf{x}^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $g_{\mathcal{H}}(4)$?



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

Growth Function: Example

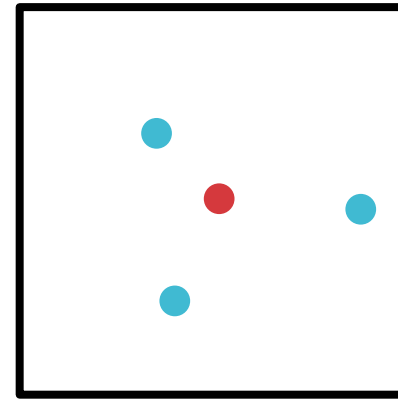
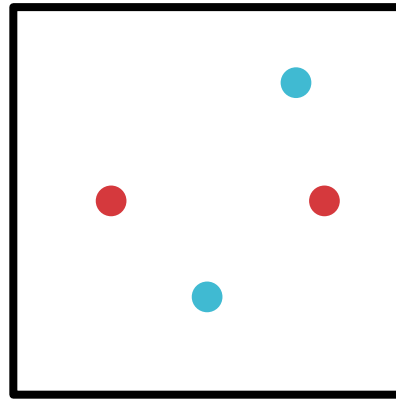
- $\mathbf{x}^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $g_{\mathcal{H}}(4)$?



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

Growth Function: Example

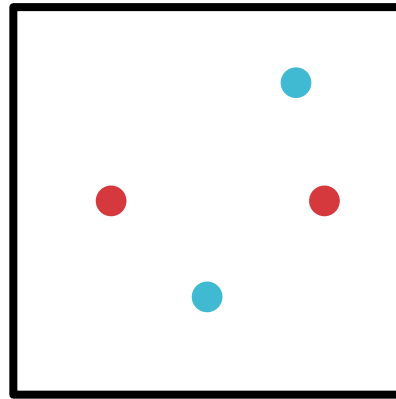
- $\mathbf{x}^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $g_{\mathcal{H}}(4)$?



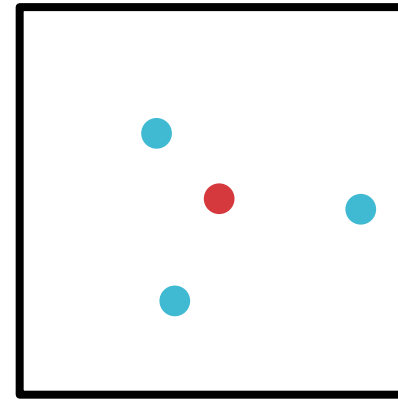
$$|\mathcal{H}(S_1)| = 14$$

Growth Function: Example

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



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



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

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 $R(h) \geq \epsilon$ have $\hat{R}(h) > 0$

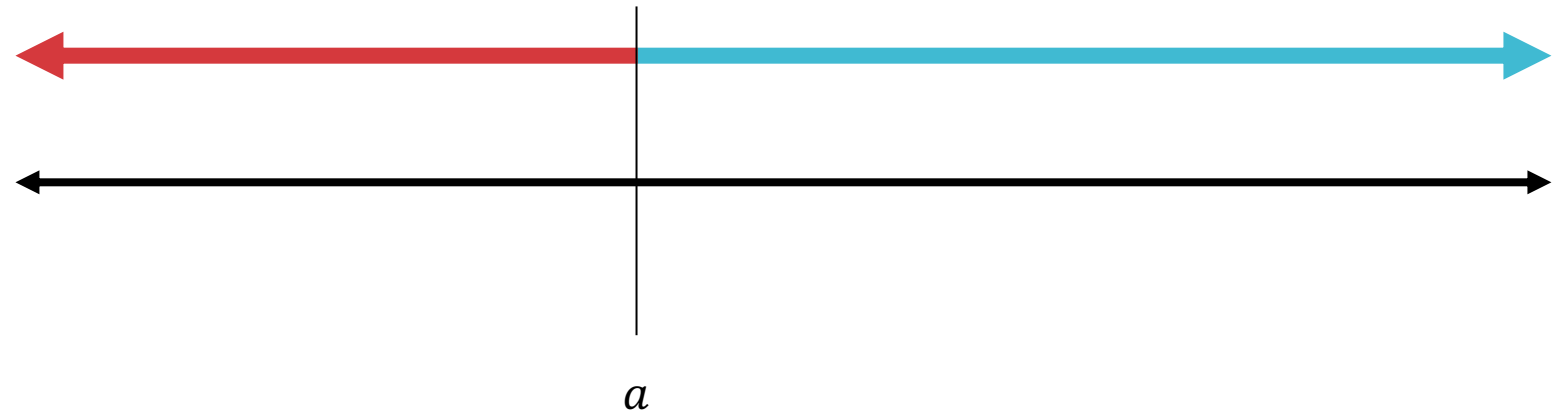
- M appears on both sides of the inequality...

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

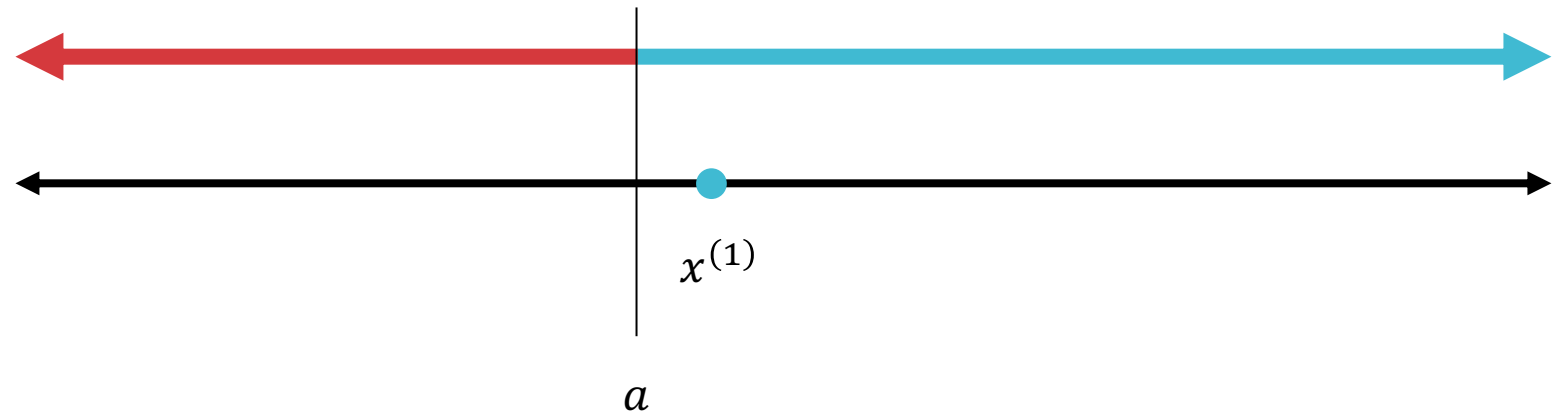
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

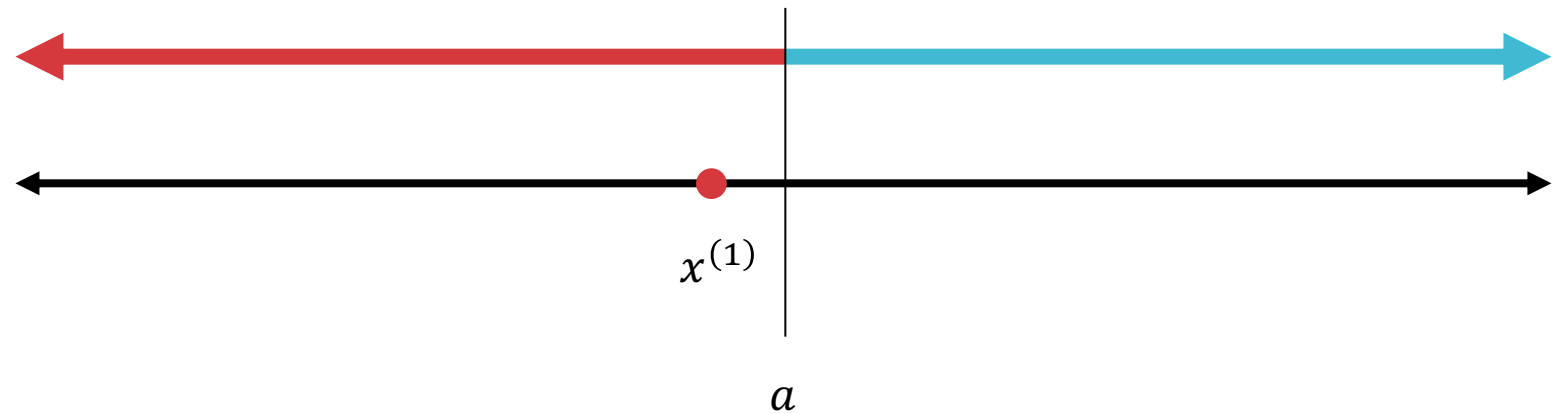
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

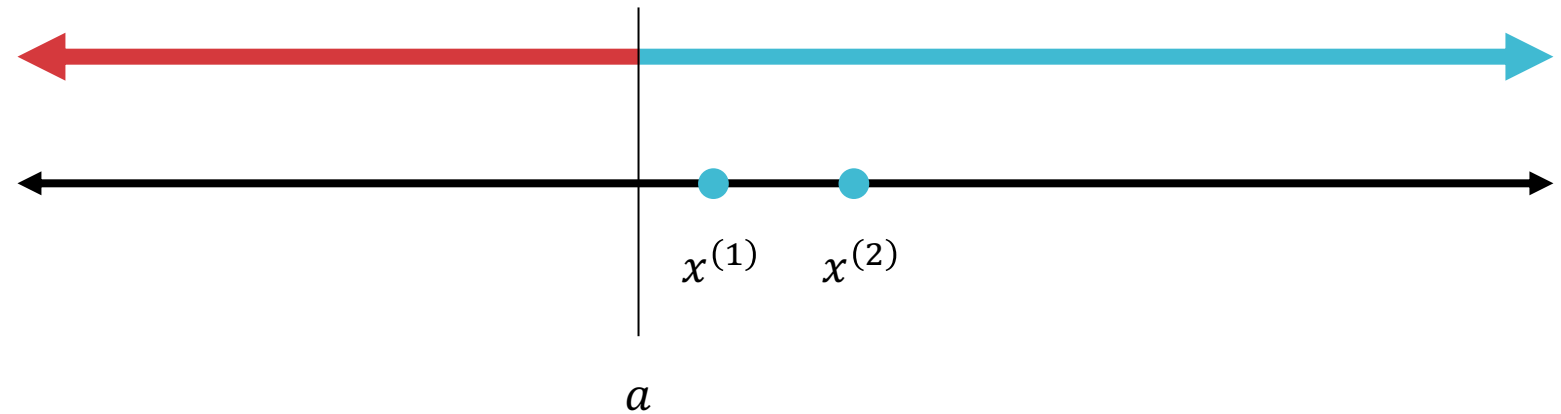
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

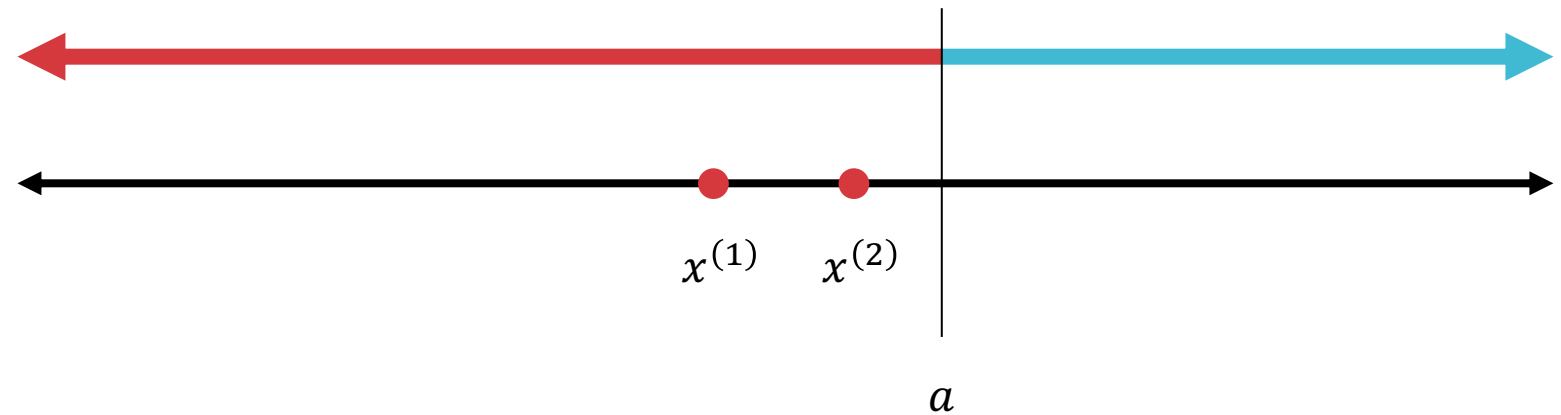
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

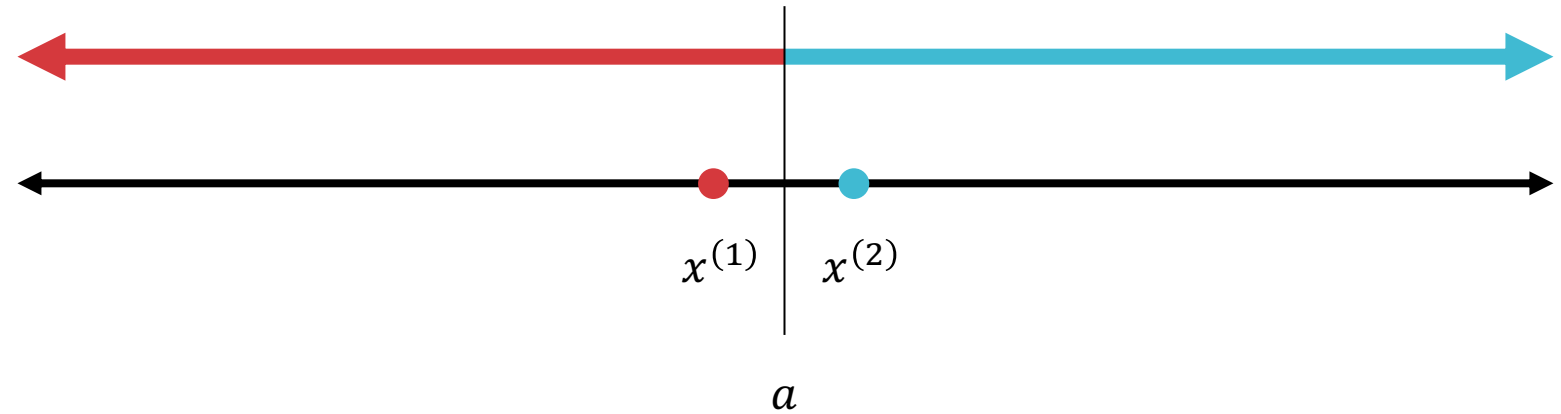
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

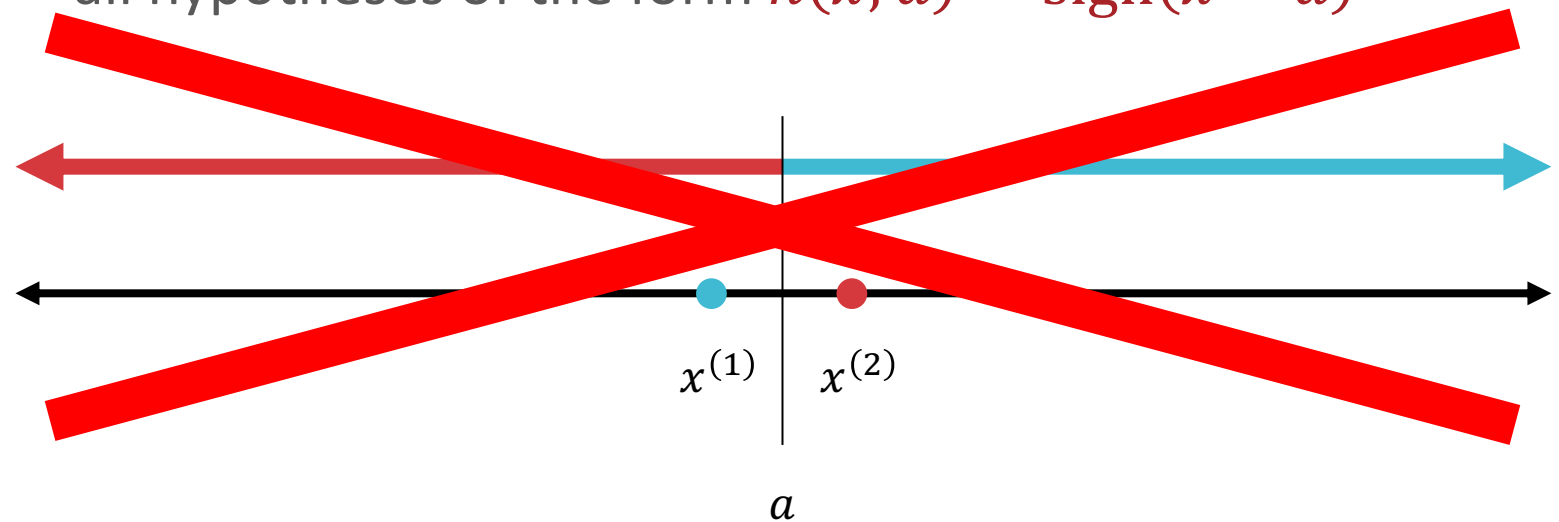
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

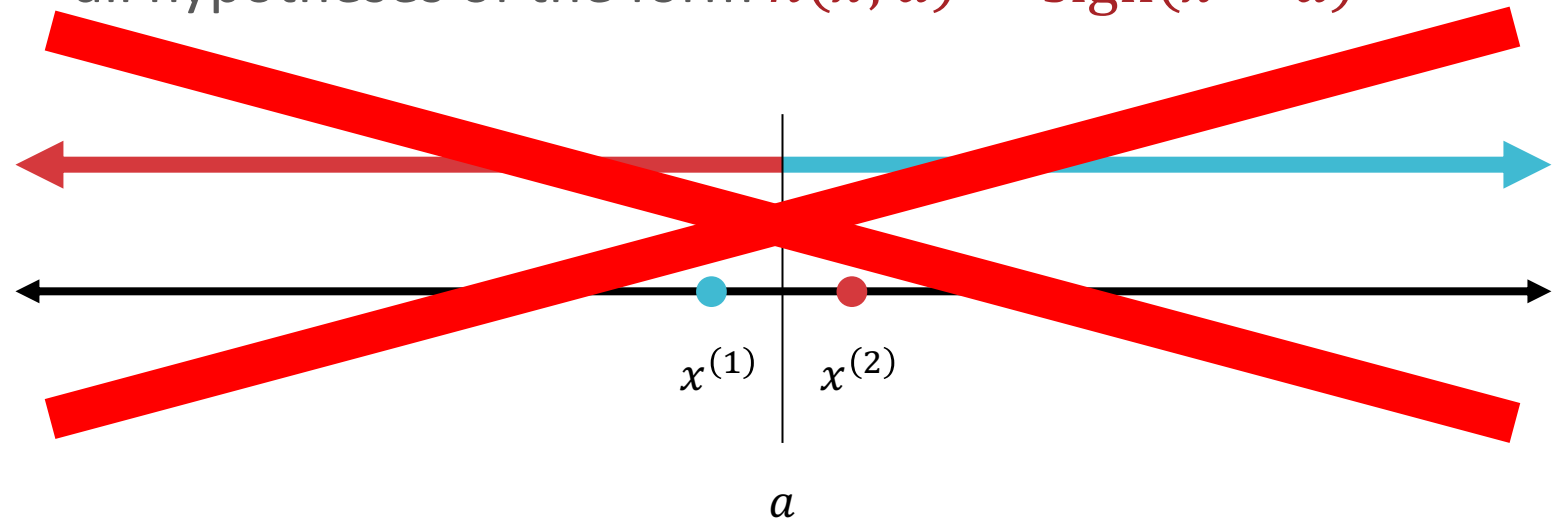
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

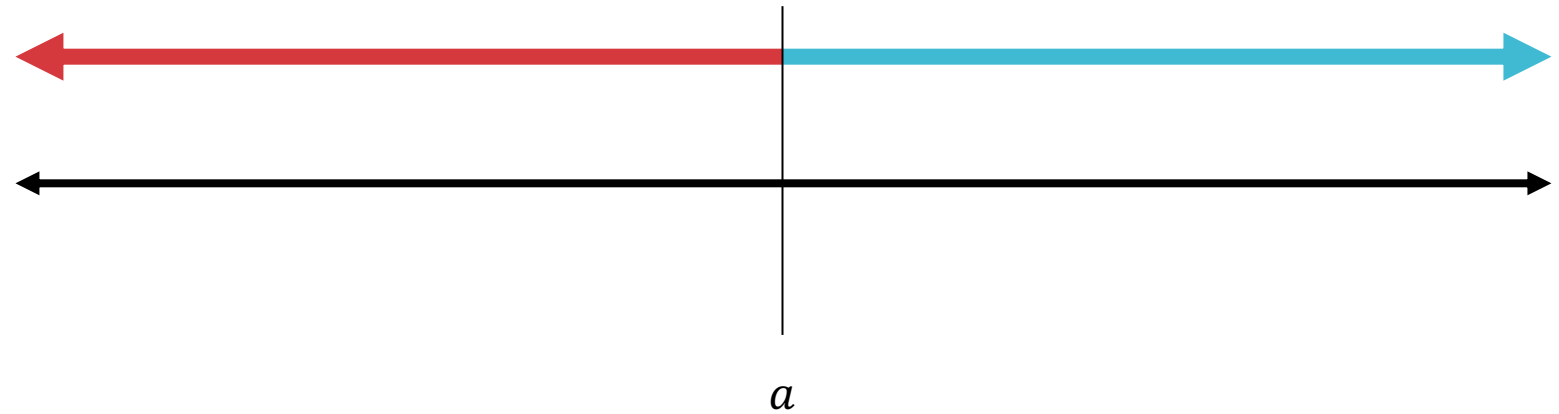
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

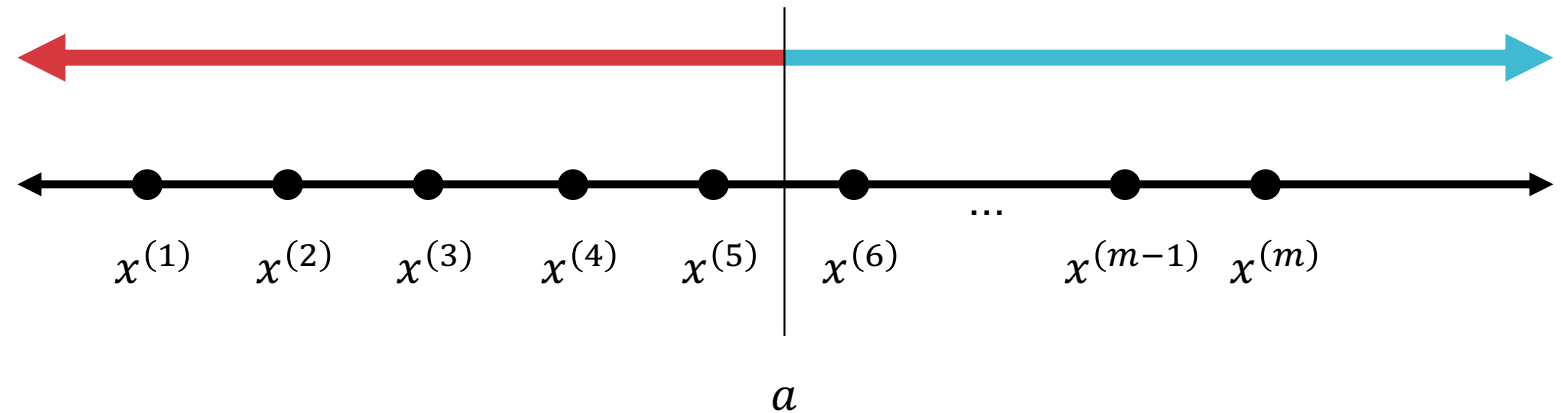
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

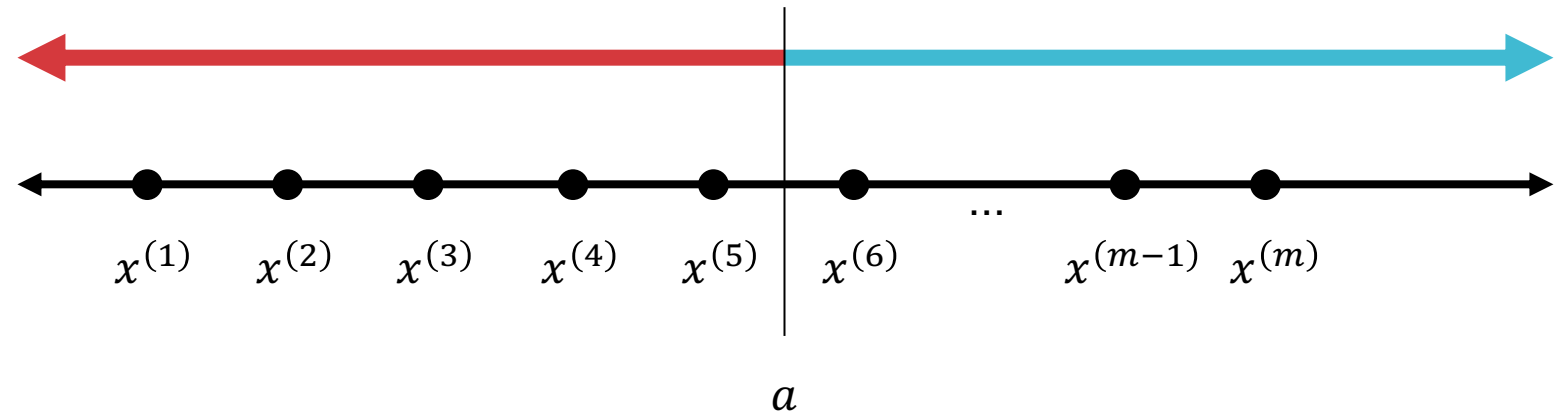
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

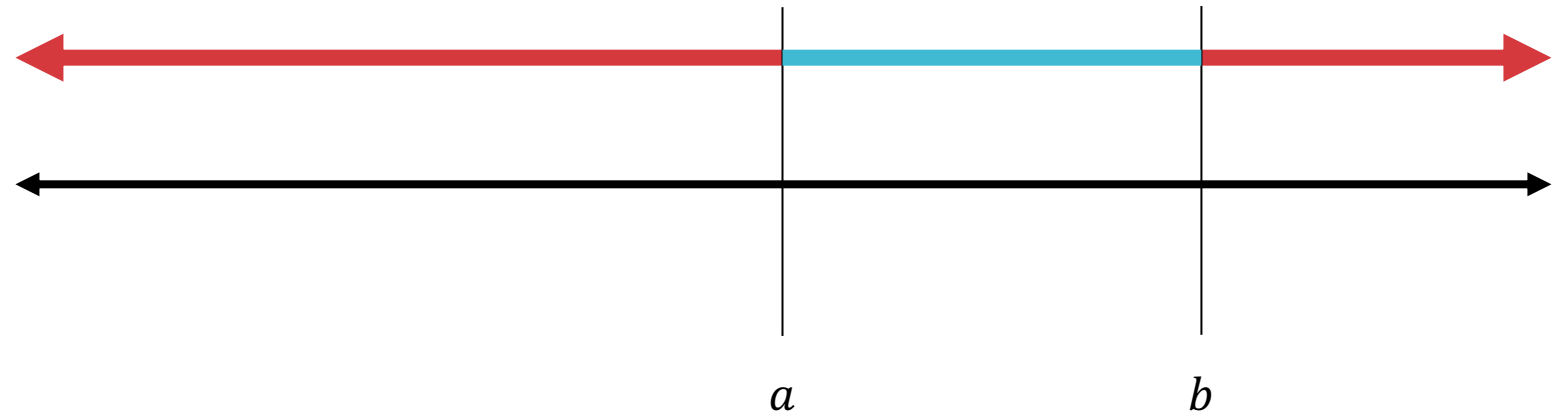
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \text{sign}(x - a)$



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

VC-Dimension: Example

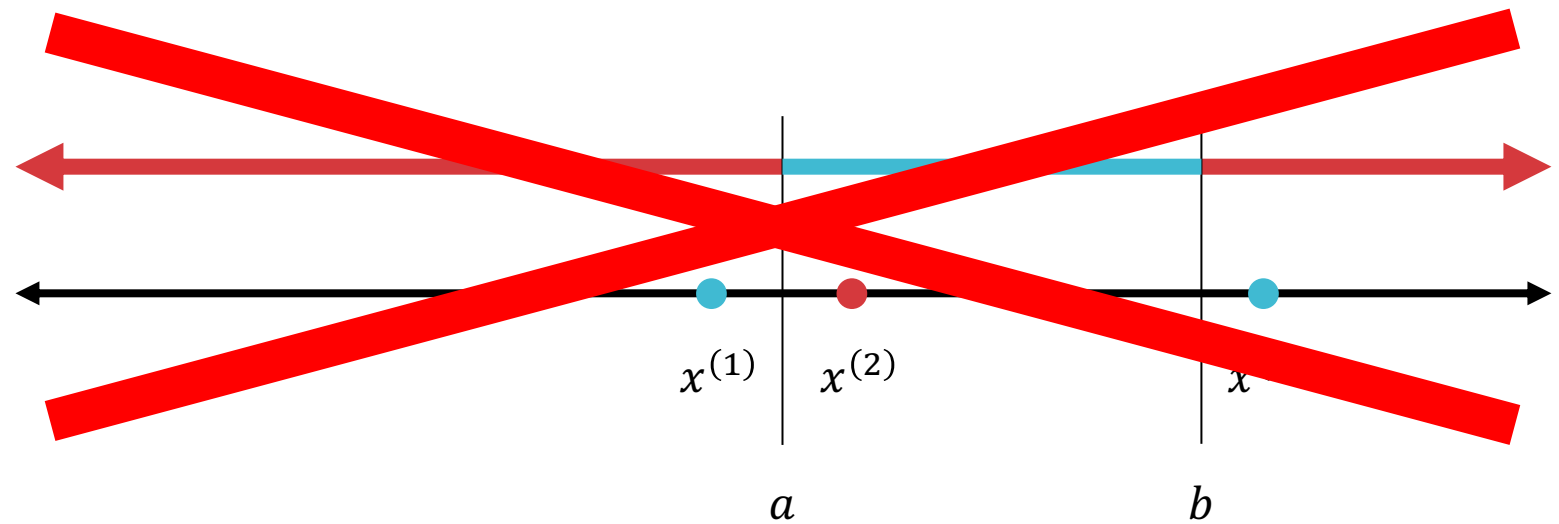
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive intervals



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

VC-Dimension: Example

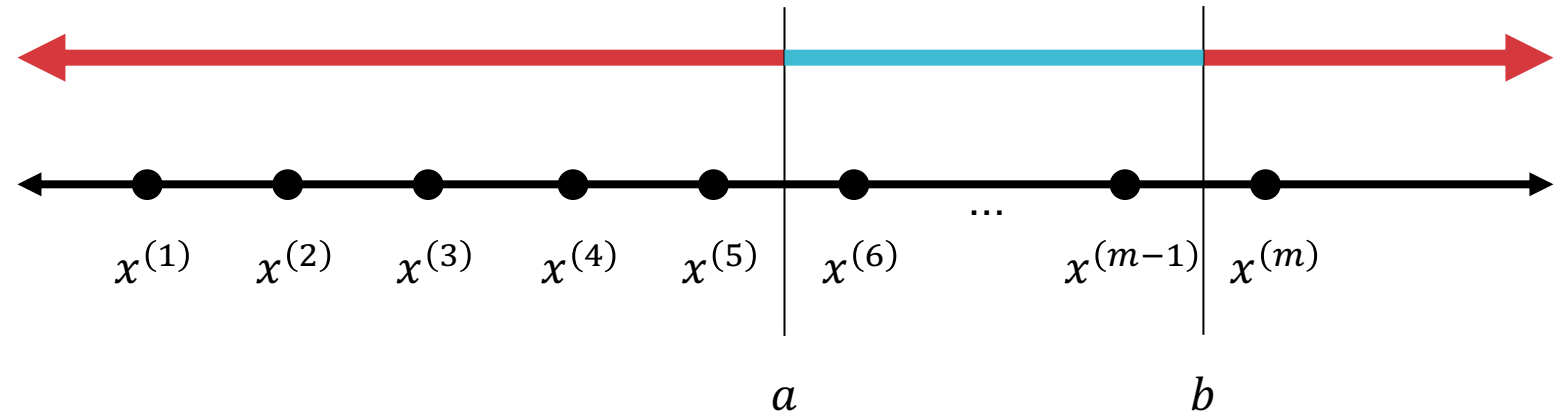
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive intervals



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

VC-Dimension: Example

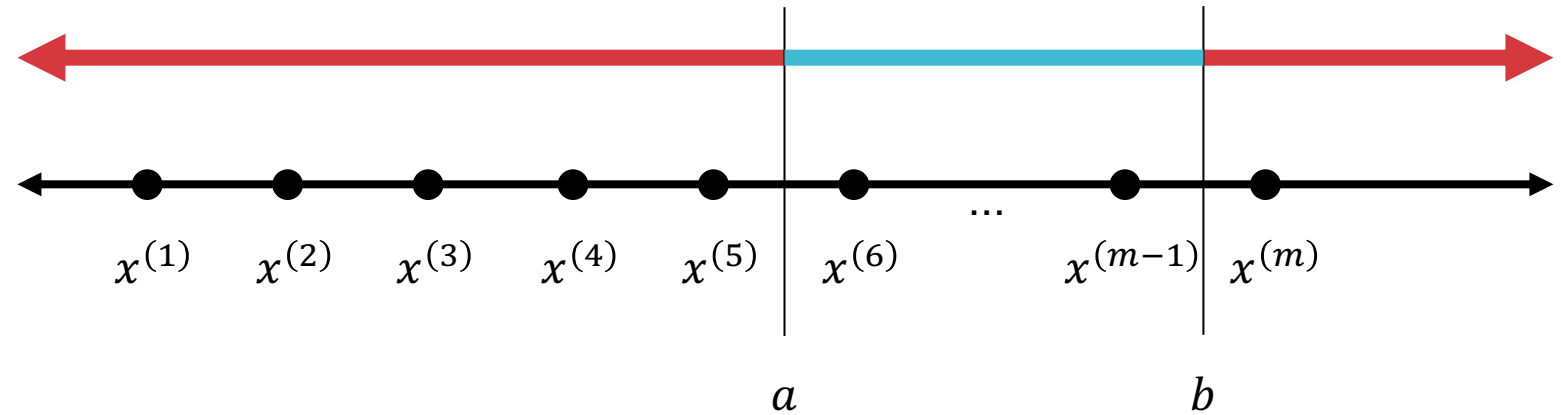
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive intervals



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

VC-Dimension: Example

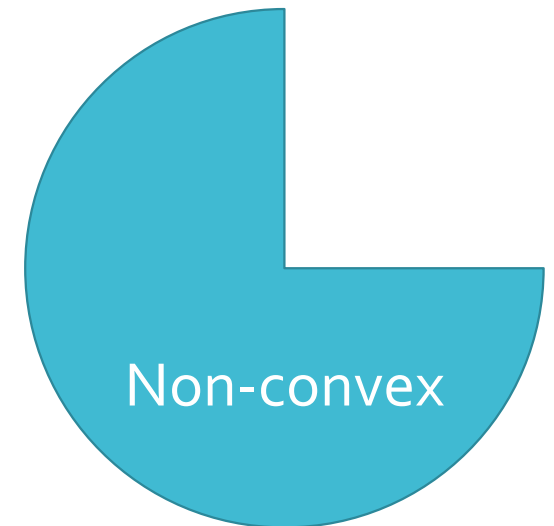
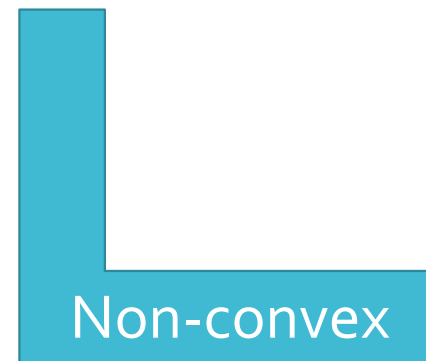
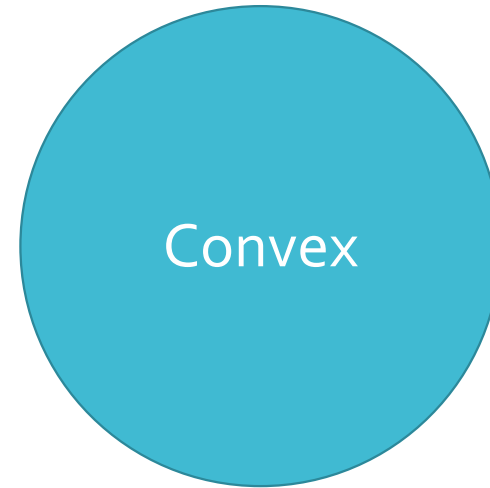
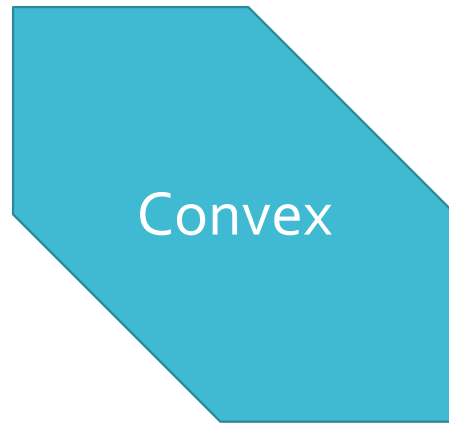
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive intervals



- $d_{VC}(\mathcal{H}) = 2$ and $g_{\mathcal{H}}(m) = \binom{m+1}{2} + 1 = O(m^2)$

Growth Function: Example

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

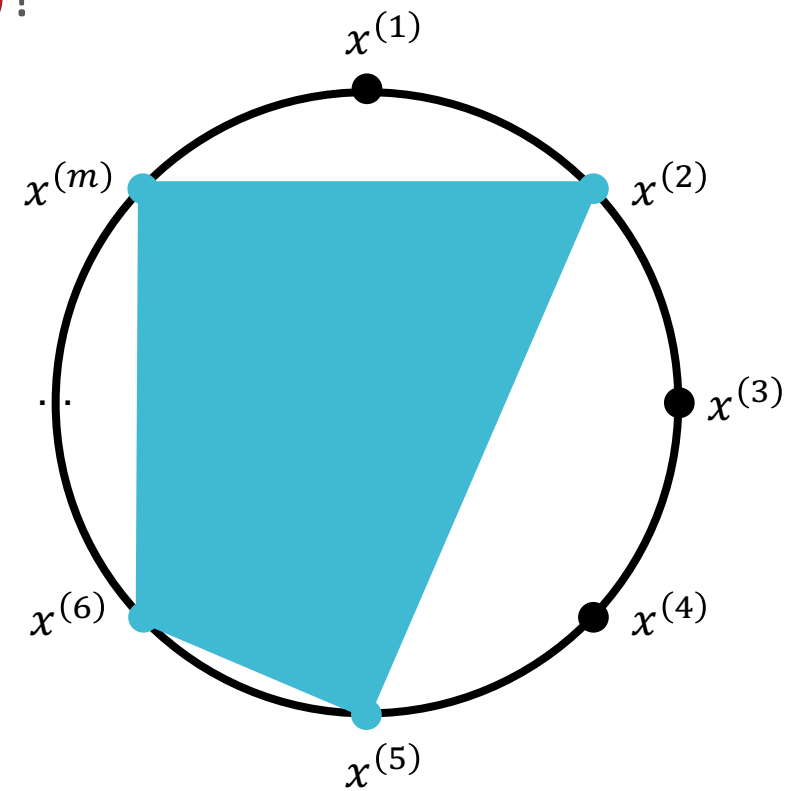


Growth Function: Example

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

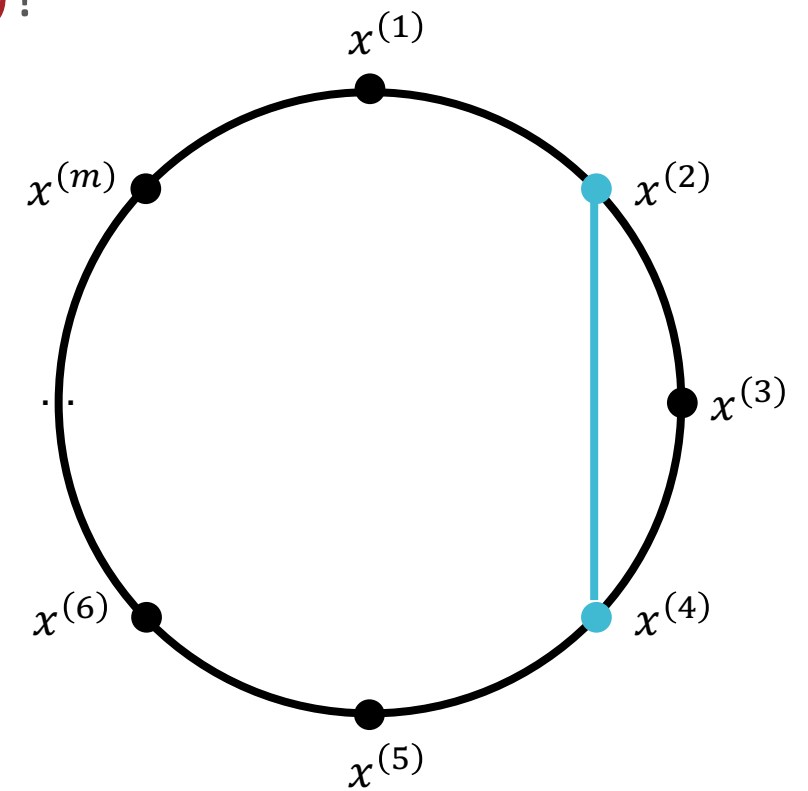
Growth Function: Example

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



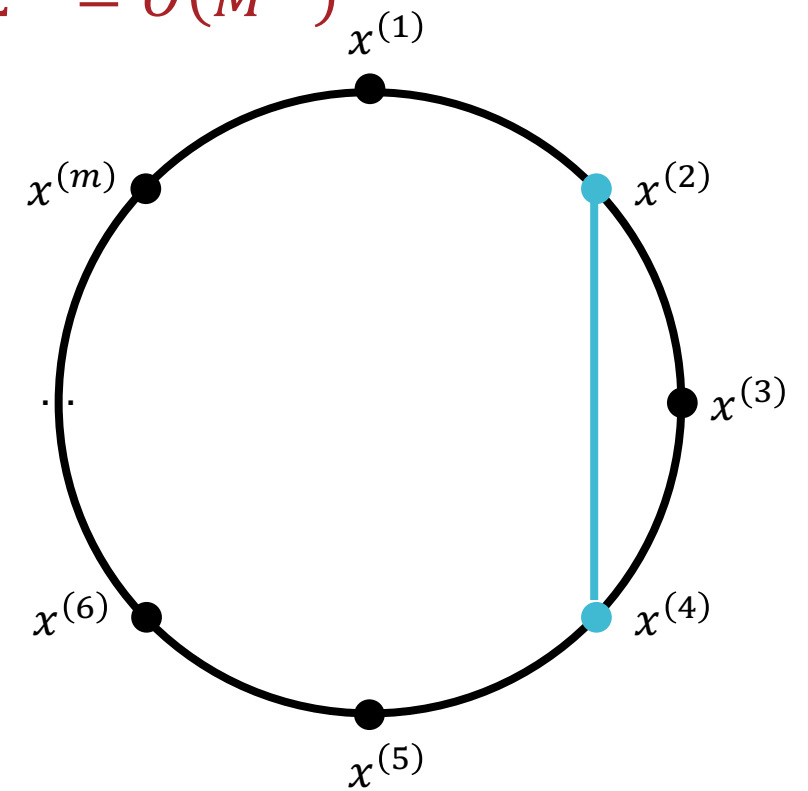
Growth Function: Example

- $x^{(m)} \in \mathbb{R}^2$ 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
approximate c^* ?

Approximation Generalization Tradeoff

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

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

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

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 $d_{VC}(\mathcal{H})$

Bias-Variance Tradeoff

- Assume a regression task with squared error and let $h_S \in \mathcal{H}$ = the hypothesis trained on training data S
- $err_D(h_S) = \mathbb{E}_{\mathbf{x} \sim D} \left[(h_S(\mathbf{x}) - c^*(\mathbf{x}))^2 \right]$
- $\mathbb{E}_S[err_D(h_S)] = \mathbb{E}_S \left[\mathbb{E}_{\mathbf{x} \sim D} \left[(h_S(\mathbf{x}) - c^*(\mathbf{x}))^2 \right] \right]$
 $= \mathbb{E}_{\mathbf{x} \sim D} \left[\mathbb{E}_S \left[(h_S(\mathbf{x}) - c^*(\mathbf{x}))^2 \right] \right]$
 $= \mathbb{E}_{\mathbf{x} \sim D} \left[\mathbb{E}_S [h_S(\mathbf{x})^2 - 2h_S(\mathbf{x})c^*(\mathbf{x}) + c^*(\mathbf{x})^2] \right]$
 $= \mathbb{E}_{\mathbf{x} \sim D} \left[\mathbb{E}_S [h_S(\mathbf{x})^2] - 2\bar{h}(\mathbf{x})c^*(\mathbf{x}) + c^*(\mathbf{x})^2 \right]$
- where $\bar{h}(\vec{x}) = \mathbb{E}_S[h_S(\mathbf{x})] \approx \frac{1}{k} \sum_{i=1}^k h_{S_i}(\mathbf{x})$

Bias-Variance Tradeoff

- Assume a regression task with squared error and let $h_S \in \mathcal{H}$ = the hypothesis trained on training data S
- $err_D(h_S) = \mathbb{E}_{\mathbf{x} \sim D} \left[(h_S(\mathbf{x}) - c^*(\mathbf{x}))^2 \right]$
- $$\begin{aligned} \mathbb{E}_S[err_D(h_S)] &= \mathbb{E}_{\mathbf{x} \sim D} \left[\mathbb{E}_S[h_S(\mathbf{x})^2] - 2\bar{h}(\mathbf{x})c^*(\mathbf{x}) + c^*(\mathbf{x})^2 \right] \\ &= \mathbb{E}_{\mathbf{x} \sim D} \left[\mathbb{E}_S[h_S(\mathbf{x})^2] - \bar{h}(\mathbf{x})^2 \right. \\ &\quad \left. + \bar{h}(\mathbf{x})^2 - 2\bar{h}(\mathbf{x})c^*(\mathbf{x}) + c^*(\mathbf{x})^2 \right] \\ &= \mathbb{E}_{\mathbf{x} \sim D} \left[\mathbb{E}_S[h_S(\mathbf{x})^2 - \bar{h}(\mathbf{x})^2] + \left(\bar{h}(\mathbf{x}) - c^*(\mathbf{x}) \right)^2 \right] \\ &= \mathbb{E}_{\mathbf{x} \sim D} \left[\text{Variance of } h_S(\mathbf{x}) + \text{Bias of } \bar{h}(\mathbf{x}) \right] \end{aligned}$$

Bias-Variance Tradeoff

How much does h change if the training data set changes?

$$\mathbb{E}_S[\text{err}_D(h_S)] = \mathbb{E}_{\mathbf{x} \sim D} \left[\underbrace{\mathbb{E}_S[h_S(\mathbf{x})^2 - \bar{h}(\mathbf{x})^2]}_{\text{Variance}} + \underbrace{(\bar{h}(\mathbf{x}) - c^*(\mathbf{x}))^2}_{\text{Bias}} \right]$$

How well on average does h approximate c^* ?

Bias-Variance Tradeoff

How well could h approximate anything?

$$\mathbb{E}_S[\text{err}_D(h_S)] = \mathbb{E}_{\mathbf{x} \sim D} \left[\underbrace{\mathbb{E}_S[h_S(\mathbf{x})^2 - \bar{h}(\mathbf{x})^2]}_{\text{How well could } h \text{ approximate anything?}} + \underbrace{(\bar{h}(\mathbf{x}) - c^*(\mathbf{x}))^2}_{\text{How well on average does } h \text{ approximate } c^*?} \right]$$

How well on average does h approximate c^* ?

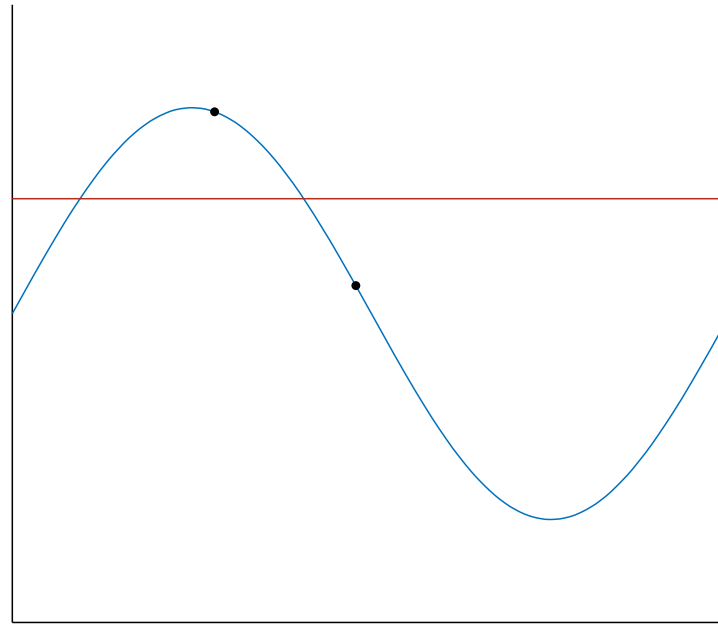
Bias-Variance Tradeoff

$$\mathbb{E}_S[\text{err}_D(h_S)] = \mathbb{E}_{\mathbf{x} \sim D} \left[\underbrace{\mathbb{E}_S[h_S(\mathbf{x})^2 - \bar{h}(\mathbf{x})^2]}_{\text{Increases as } \mathcal{H} \text{ becomes more complex}} + \underbrace{(\bar{h}(\mathbf{x}) - c^*(\mathbf{x}))^2}_{\text{Decreases as } \mathcal{H} \text{ becomes more complex}} \right]$$

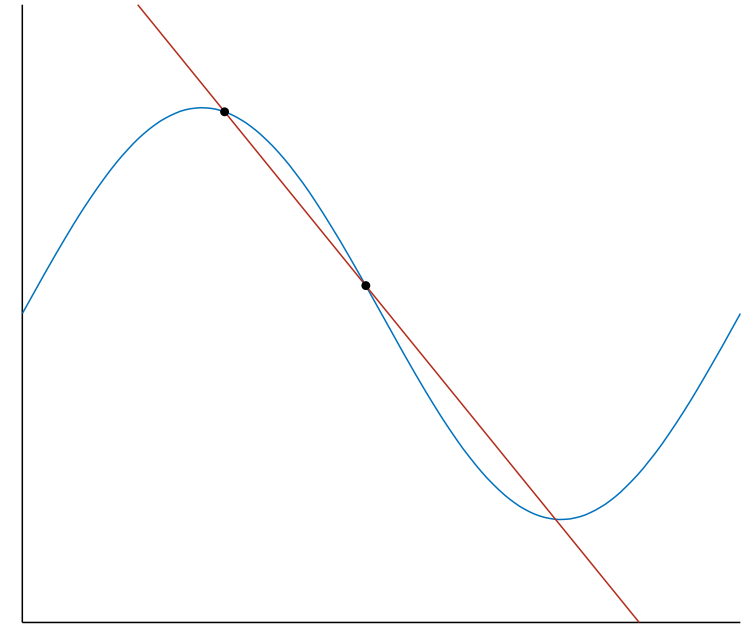
Bias-Variance Tradeoff: Example

- $x^{(i)} \in \mathbb{R}$ and $D = \text{Uniform}(0, 2\pi)$
- $c^* = \sin(\cdot)$, i.e., $y = \sin(x)$
- $N = 2 \rightarrow \mathcal{D} = \{(x^{(1)}, \sin(x^{(1)})), (x^{(2)}, \sin(x^{(2)}))\}$
- $\mathcal{H}_0 = \{h : h(x) = b\}$ and $\mathcal{H}_1 = \{h : h(x) = ax + b\}$

Bias-Variance Tradeoff: Example

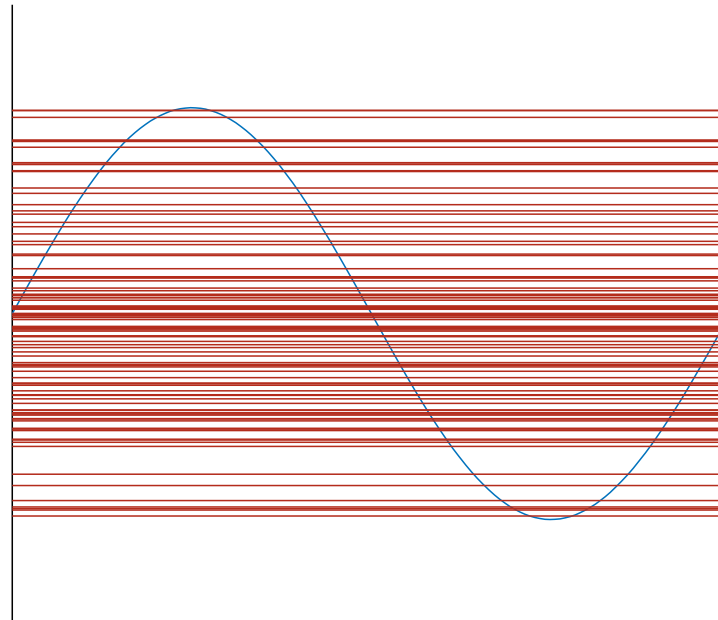


\mathcal{H}_0

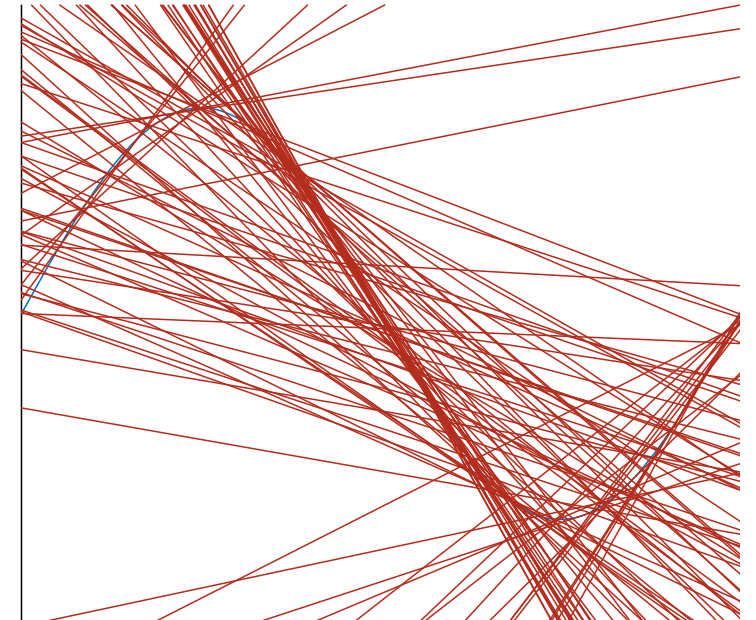


\mathcal{H}_1

Bias-Variance Tradeoff: Example

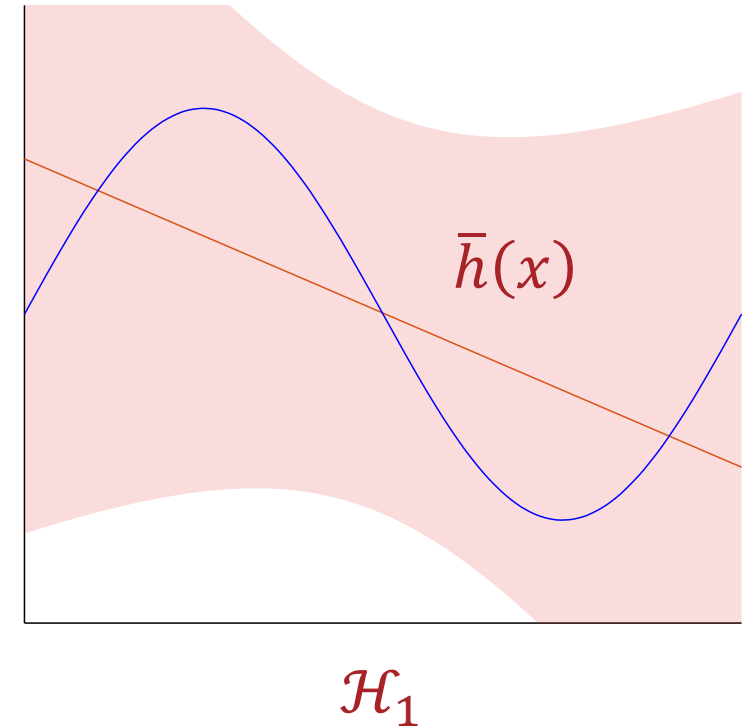
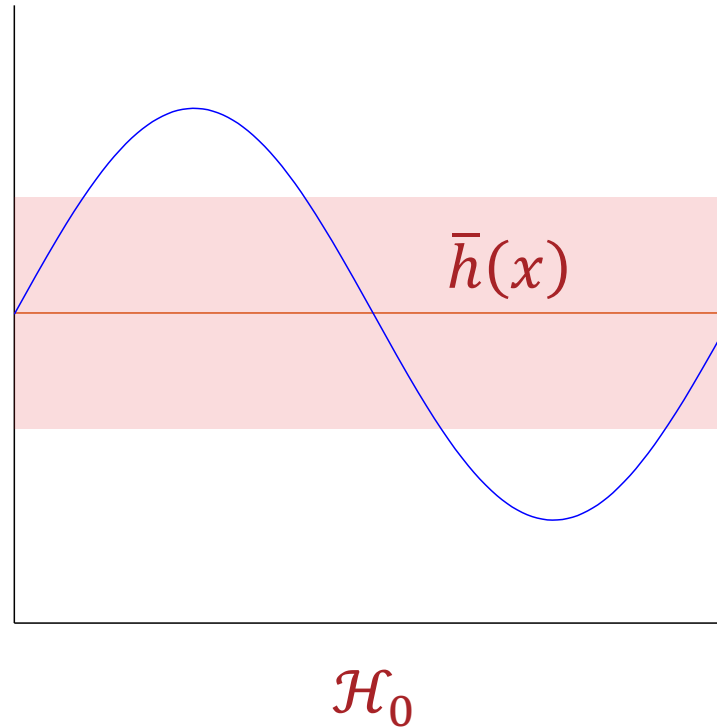


\mathcal{H}_0

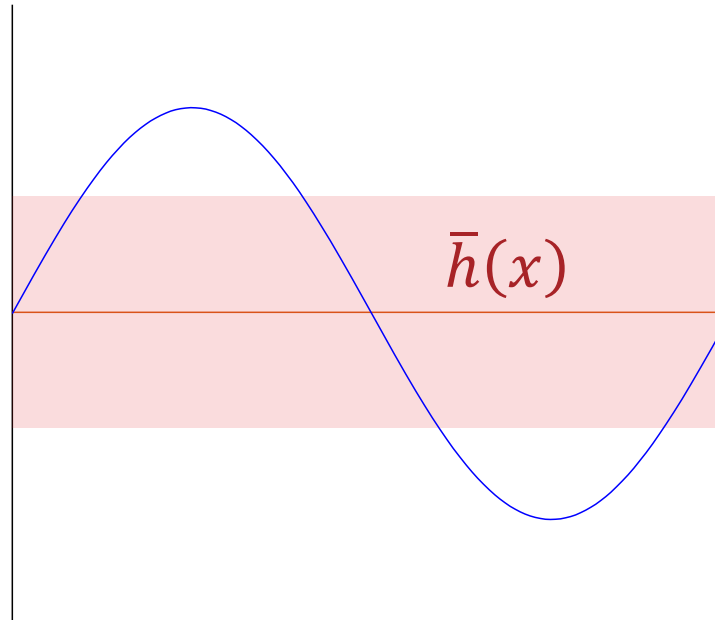


\mathcal{H}_1

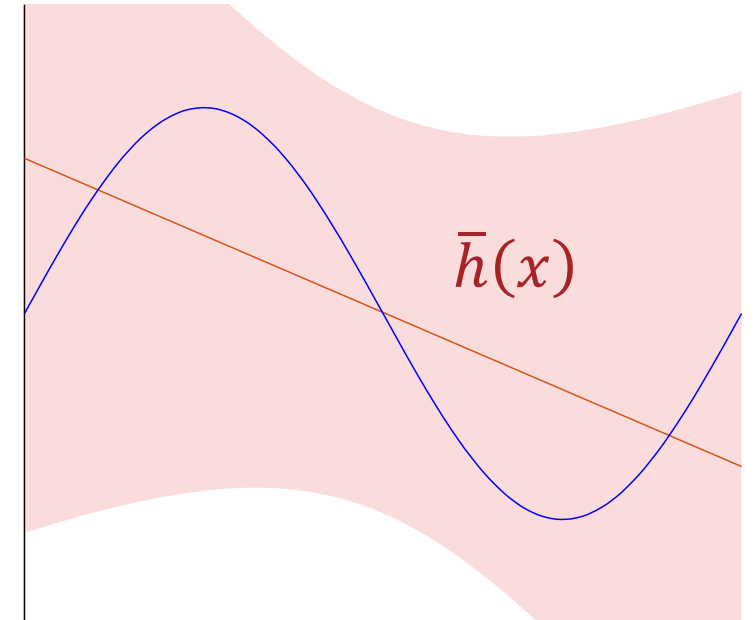
Bias-Variance Tradeoff: Example



Bias-Variance Tradeoff: Example

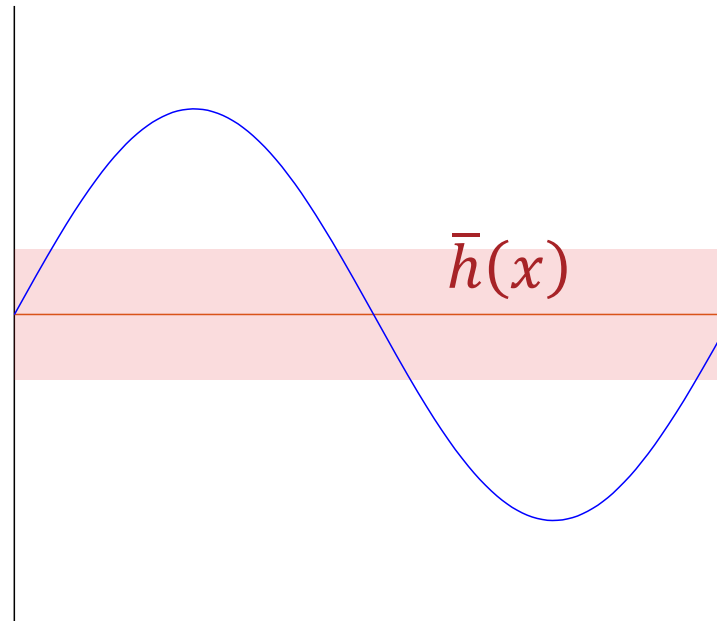


Bias of $\bar{h}(x) \approx 0.50$
Variance of $h_S(x) \approx 0.25$
 $\mathbb{E}_S[err_D(h_S)] \approx 0.75$

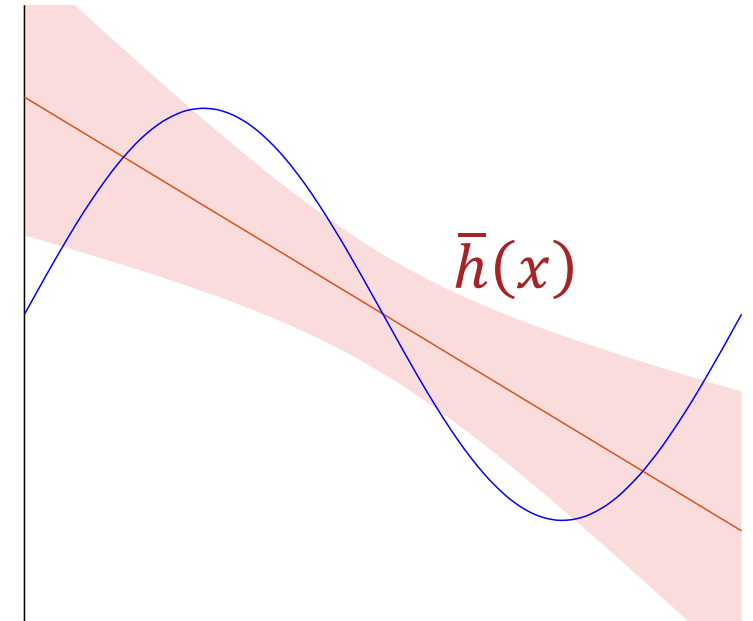


Bias of $\bar{h}(x) \approx 0.21$
Variance of $h_S(x) \approx 1.74$
 $\mathbb{E}_S[err_D(h_S)] \approx 1.95$

Bias-Variance Tradeoff: Example ($N = 5$)



Bias of $\bar{h}(x) \approx 0.50$
Variance of $h_S(x) \approx 0.10$
 $\mathbb{E}_S[err_D(h_S)] \approx 0.60$



Bias of $\bar{h}(x) \approx 0.21$
Variance of $h_S(x) \approx 0.21$
 $\mathbb{E}_S[err_D(h_S)] \approx 0.42$