# 10-301/601: Introduction to Machine Learning Lecture 17 – Learning Theory (Infinite Case)

Henry Chai

7/5/23

# Front Matter

- Announcements
  - PA4 released 6/15, due 7/13 at 11:59 PM
    - You still have one week from this Thursday!
  - Quiz 6: Deep Learning & Learning Theory on 7/11

- Recommended Readings
  - Mitchell, Chapter 7.4

# Recall: 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 $R(h) \leq \epsilon$

- 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

- 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 $1 - \delta$.

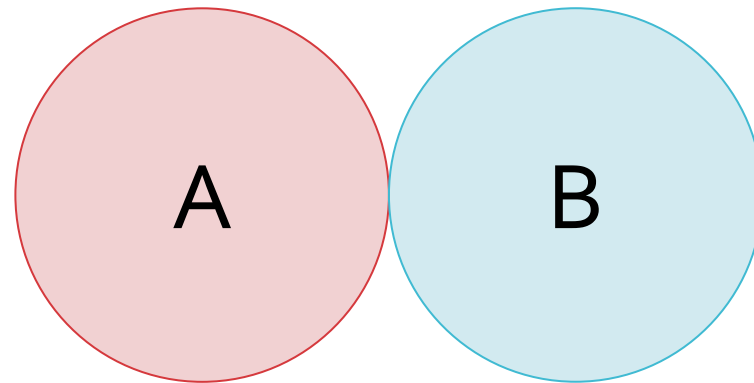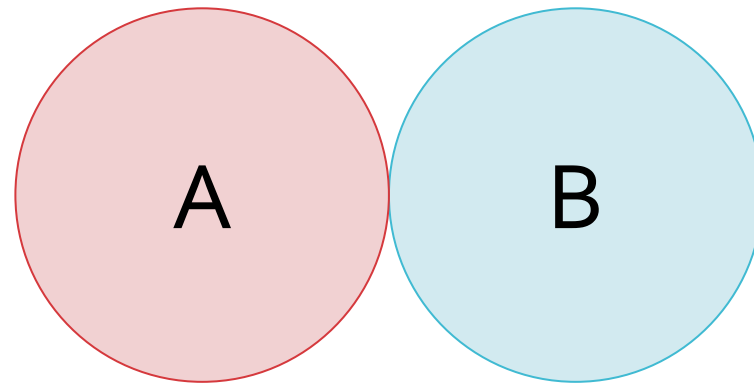# What happens when $|\mathcal{H}| = \infty$?

- 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…

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

# The Union Bound is Bad!

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

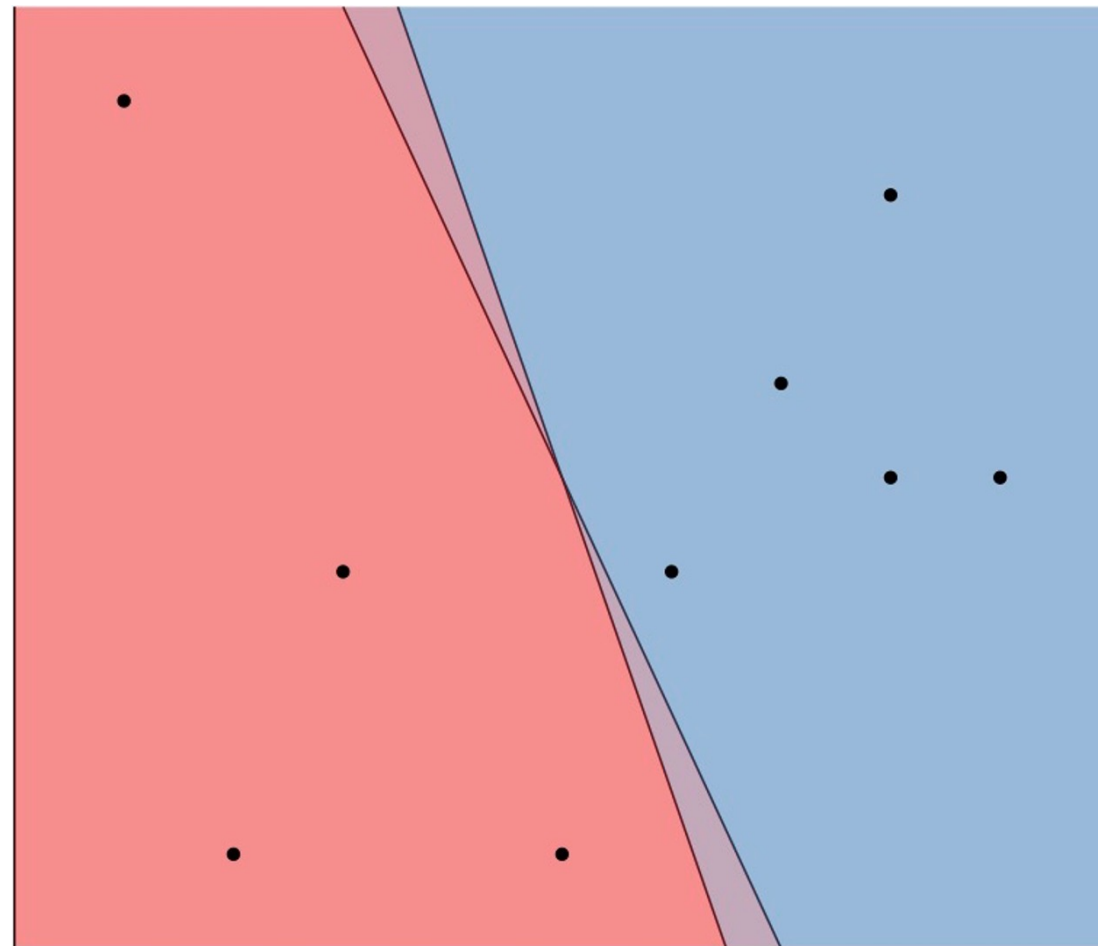$$P\{A \cup B\} = P\{A\} + P\{B\} - P\{A \cap B\}$$

A    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"
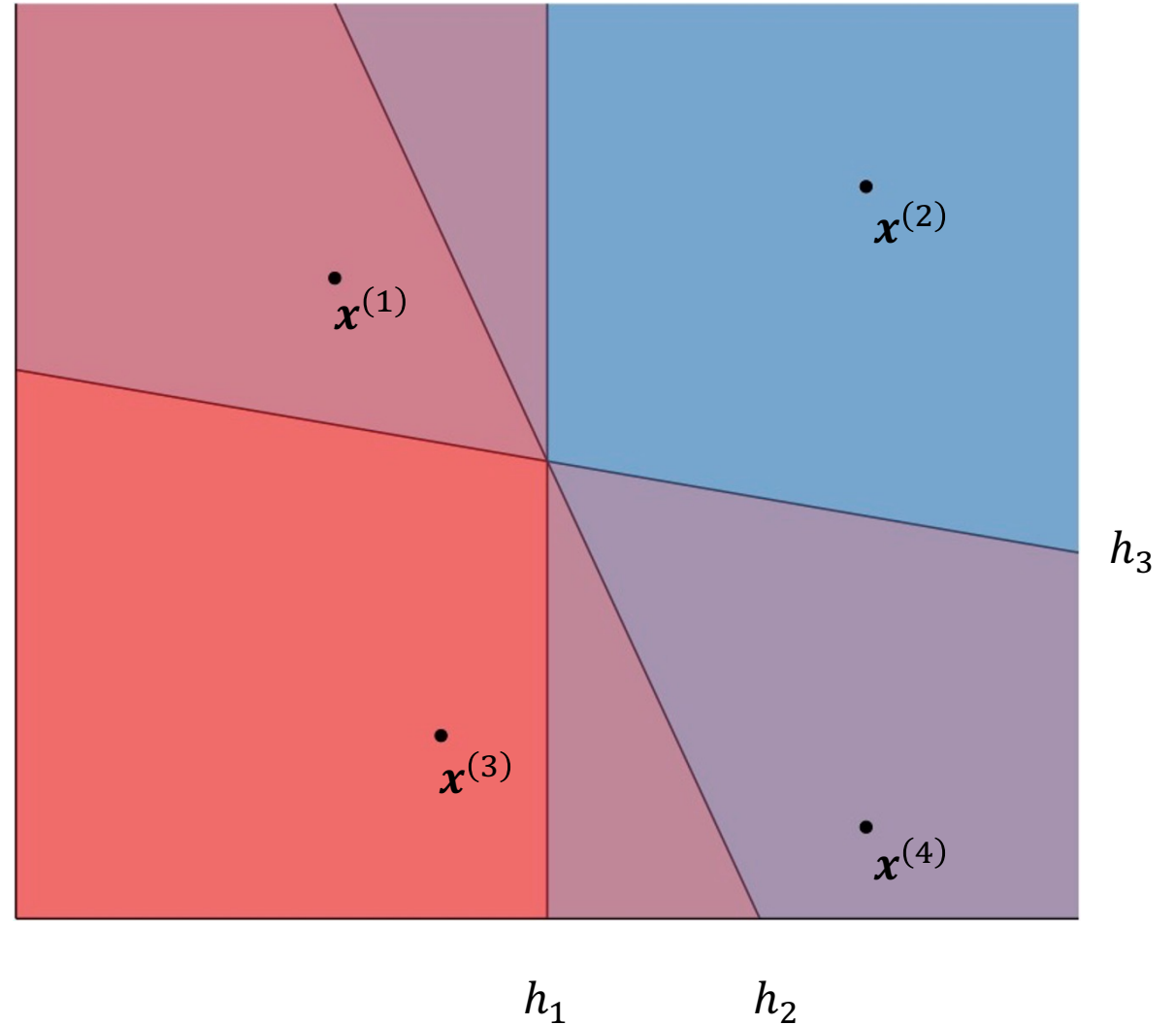
will overlap a lot!

# 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 $S = \left( \boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(M)} \right)$ 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

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

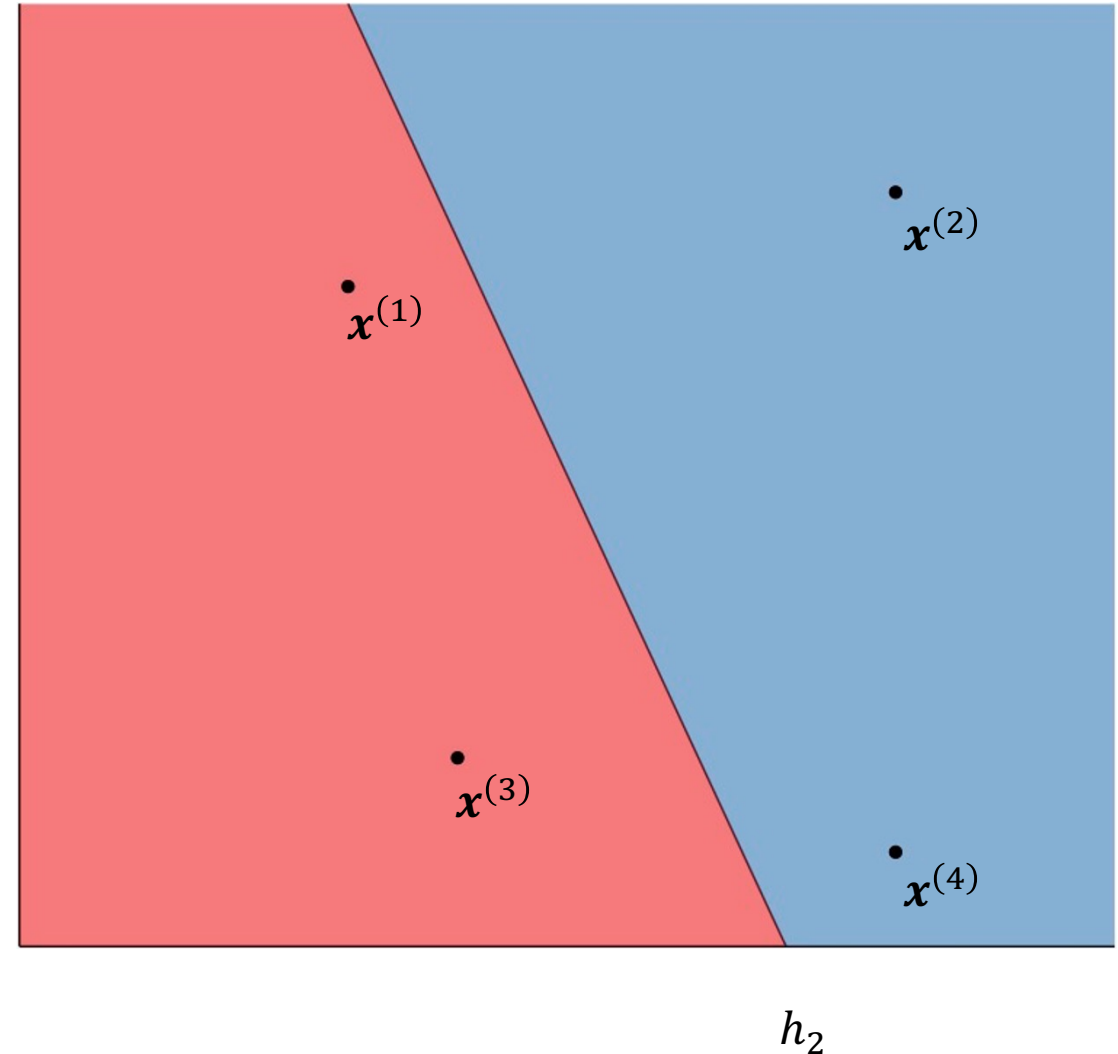# Example: Labellings

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

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

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

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

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

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

$h_1$

# Example: Labellings

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

$\left(h_2\big(\boldsymbol{x}^{(1)}\big), h_2\big(\boldsymbol{x}^{(2)}\big), h_2\big(\boldsymbol{x}^{(3)}\big), h_2\big(\boldsymbol{x}^{(4)}\big)\right)$

$= (-1, +1, -1, +1)$



$h_2$

# Example: Labellings

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

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

$$= (+1, +1, -1, -1)$$

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

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

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

$h_3$

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

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

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

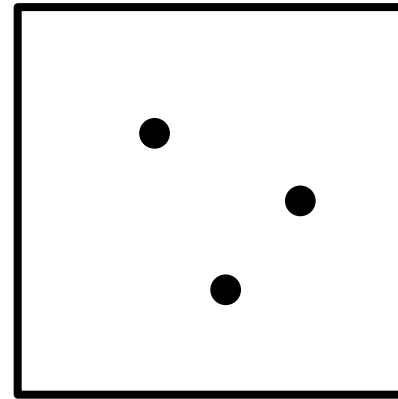- If $\exists\, S$ s.t. $|S| = M$ and $\mathcal{H}$ shatters $S$, then $g_{\mathcal{H}}(M) = 2^M$

# Growth Function: Example

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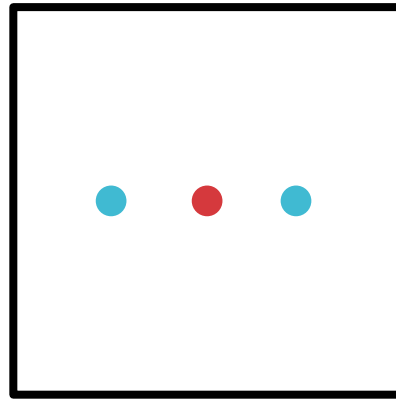
- What is $g_{\mathcal{H}}(3)$?

## Growth Function: Example

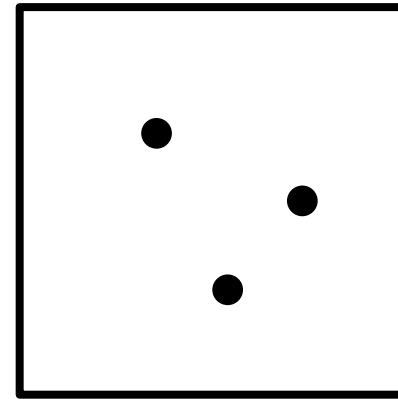- $x^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators

- What is $g_{\mathcal{H}}(3)$?

# Growth Function: Example

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

- What is $g_{\mathcal{H}}(3)$?

# Growth Function: Example

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

- What is $g_{\mathcal{H}}(3)$?



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

# Growth Function: Example

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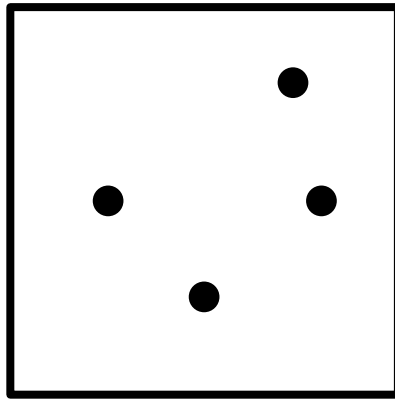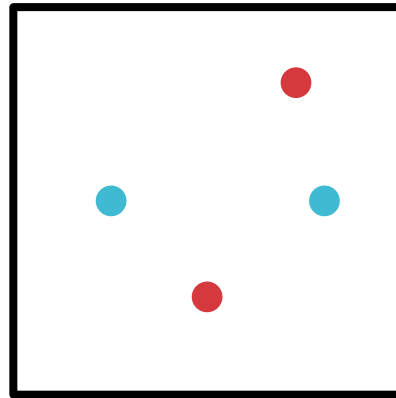
- $g_{\mathcal{H}}(3) = 8 = 2^3$

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

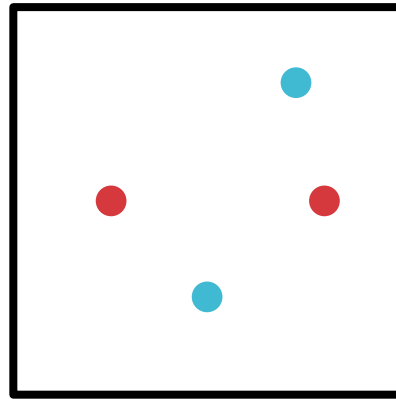## Growth Function: Example

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

- What is $g_{\mathcal{H}}(4)$?

# Growth Function: Example

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

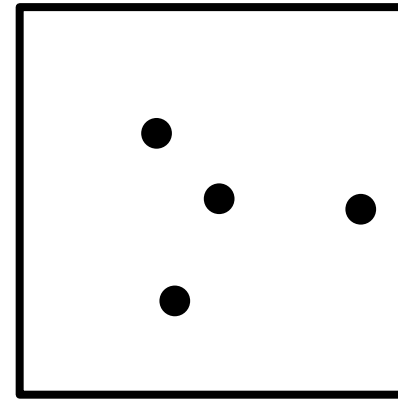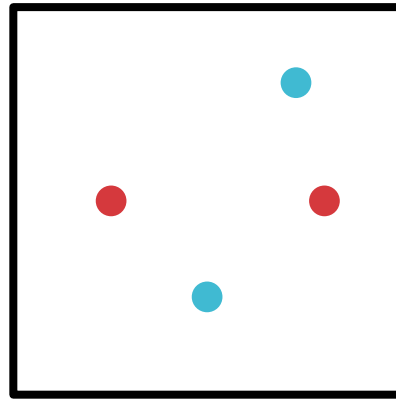- What is $g_{\mathcal{H}}(4)$?

# Growth Function: Example

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

- What is $g_{\mathcal{H}}(4)$?

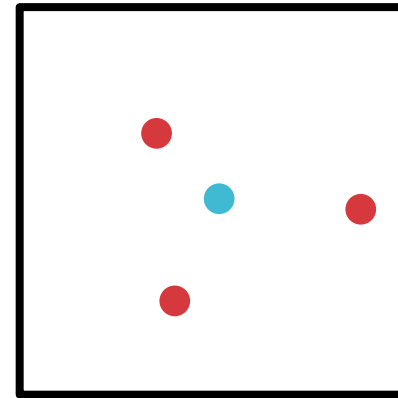$$|\mathcal{H}(S_1)| = 14$$

# Growth Function: Example

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

- What is $g_{\mathcal{H}}(4)$?



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

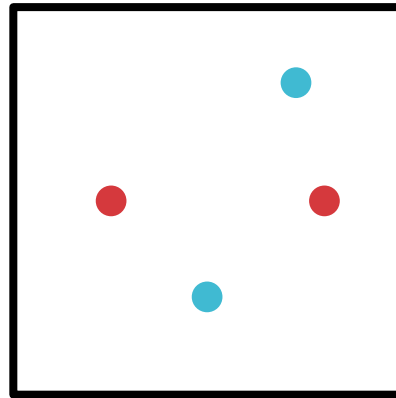# Growth Function: Example

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

- What is $g_{\mathcal{H}}(4)$?



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

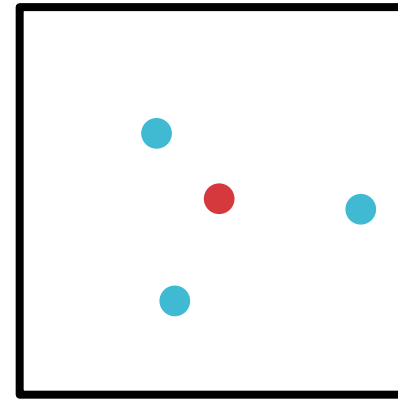# Growth Function: Example

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

- What is $g_{\mathcal{H}}(4)$?



$$|\mathcal{H}(S_1)| = 14$$
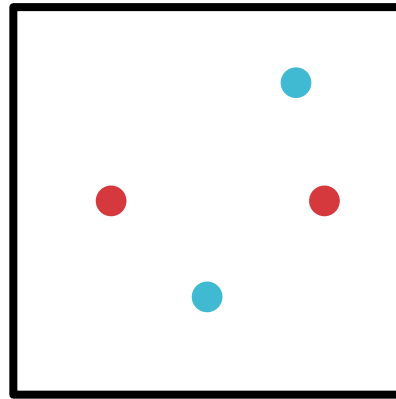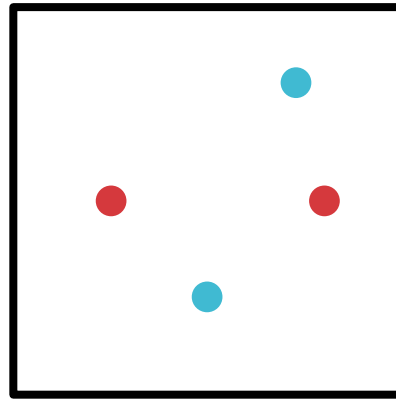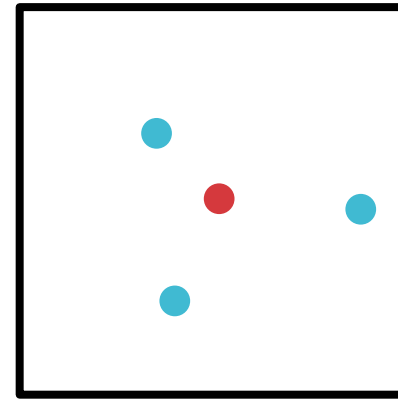
# Growth Function: Example

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

- $g_{\mathcal{H}}(4) = 14 < 2^4$

$$|\mathcal{H}(S_1)| = 14 \qquad\qquad |\mathcal{H}(S_2)| = 14$$

## 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

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

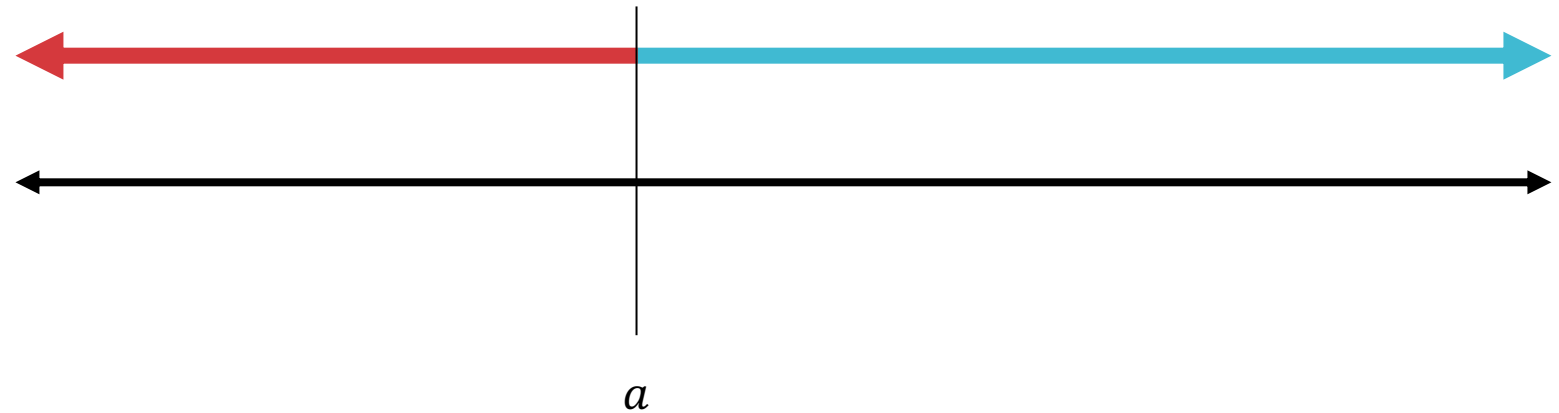then with probability at least $1 - \delta$, all $h \in \mathcal{H}$ with $R(h) \geq \epsilon$ have $\hat{R}(h) > 0$

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

  - $g_{\mathcal{H}}(M) = O\left(M^{d_{VC}(\mathcal{H})}\right)$ (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
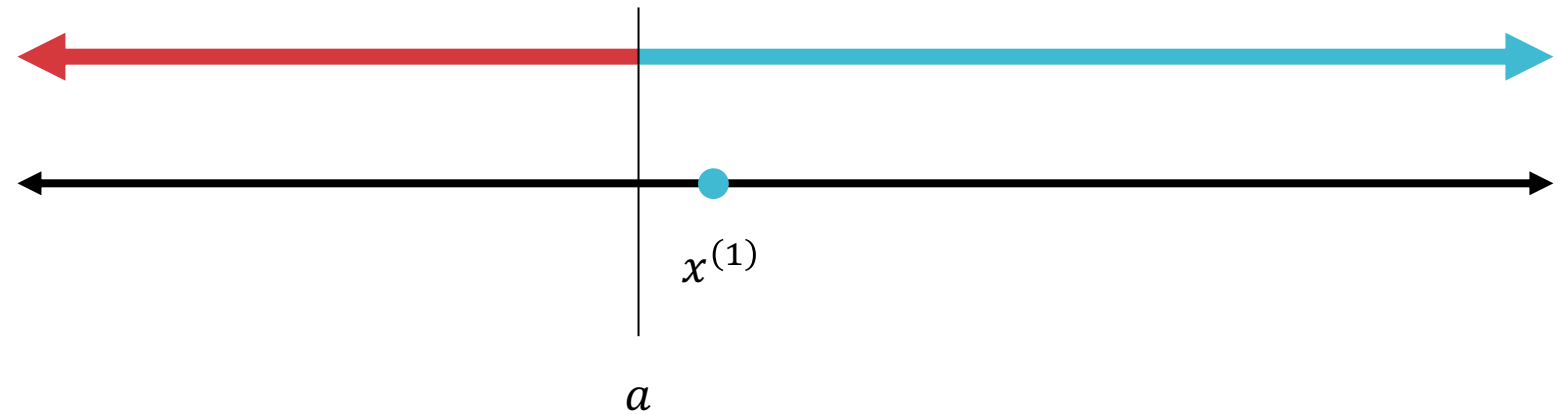
VC-Dimension: Example

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

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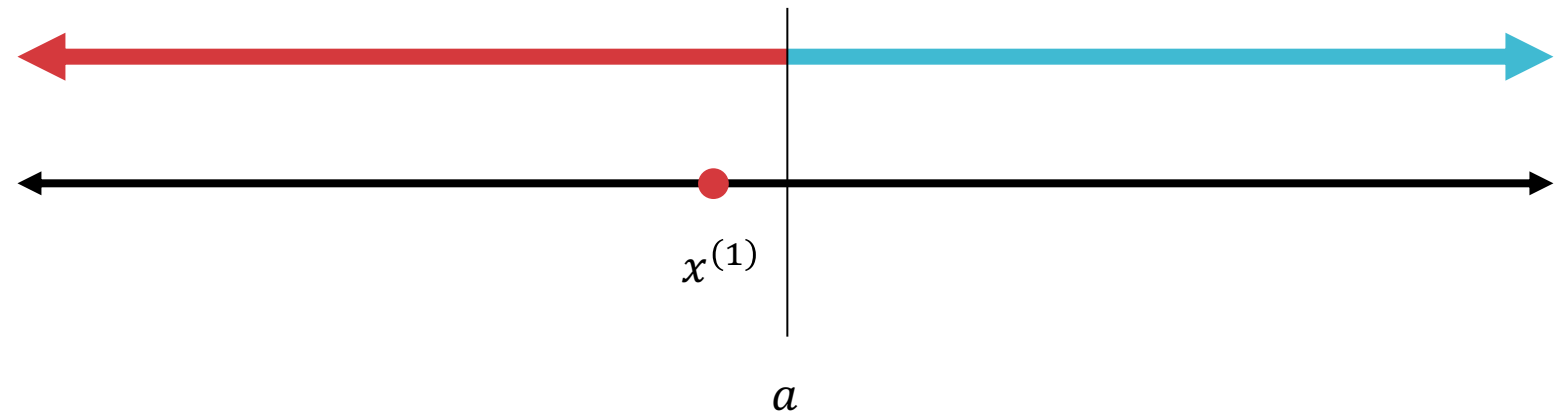
- What is $d_{VC}(\mathcal{H})$?

# VC-Dimension: Example

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



$x^{(1)}$

$a$

- 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 $h(x; a) = \text{sign}(x - a)$



$x^{(1)}$

$a$

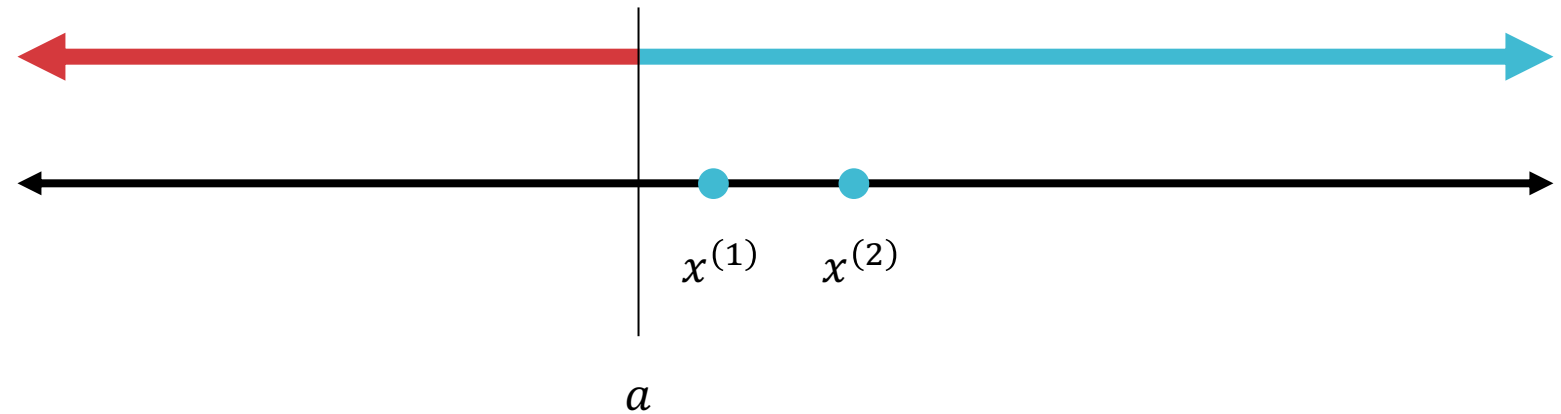- What is $d_{VC}(\mathcal{H})$?
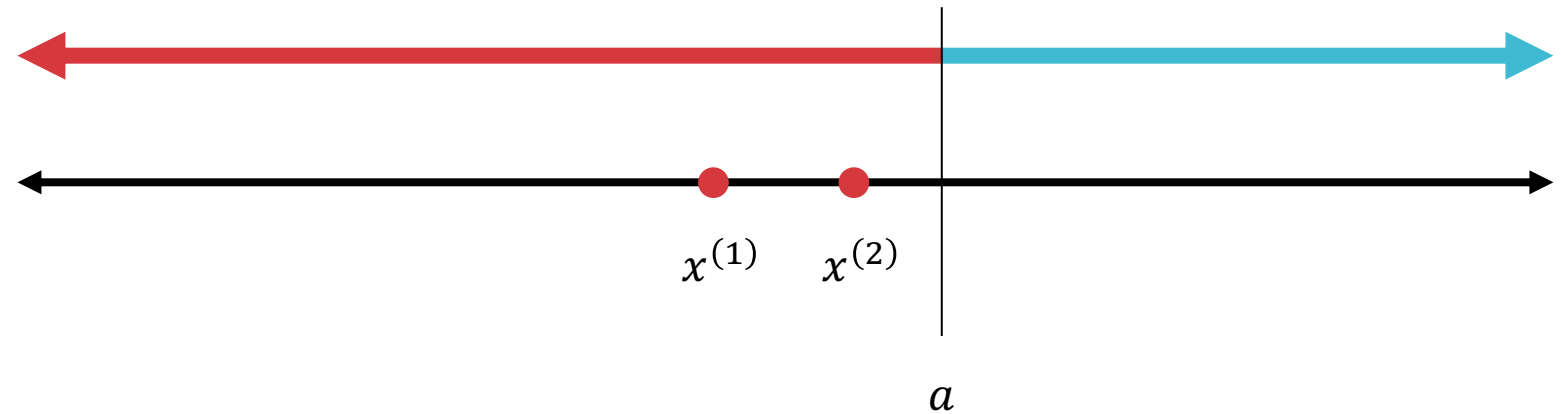
# VC-Dimension: Example

- $x^{(m)} \in \mathbb{R}$ 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 $h(x; a) = \text{sign}(x - a)$



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

$a$

- What is $d_{VC}(\mathcal{H})$?
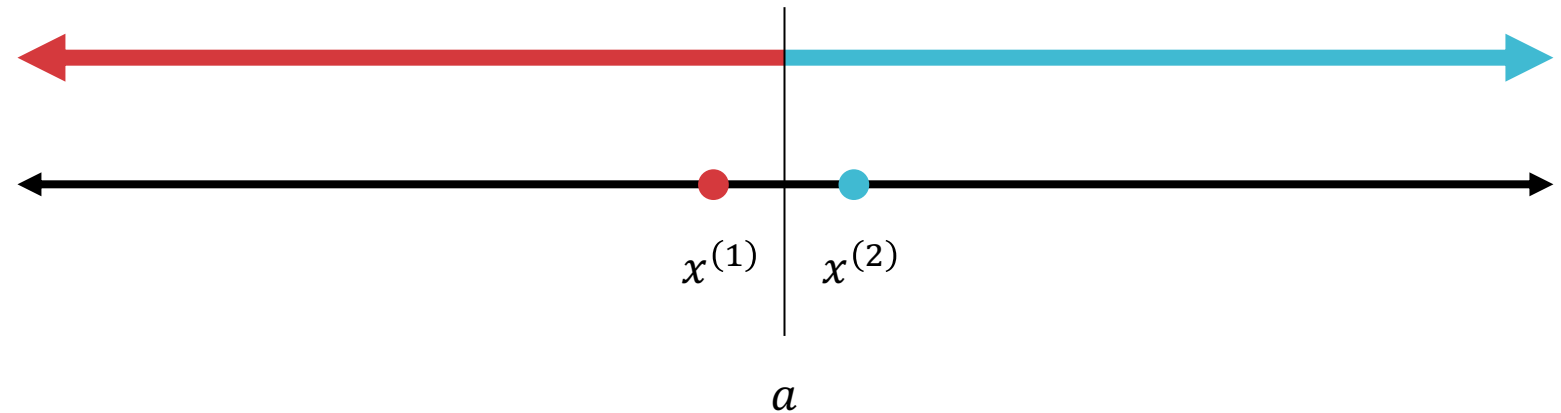
# VC-Dimension: Example

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


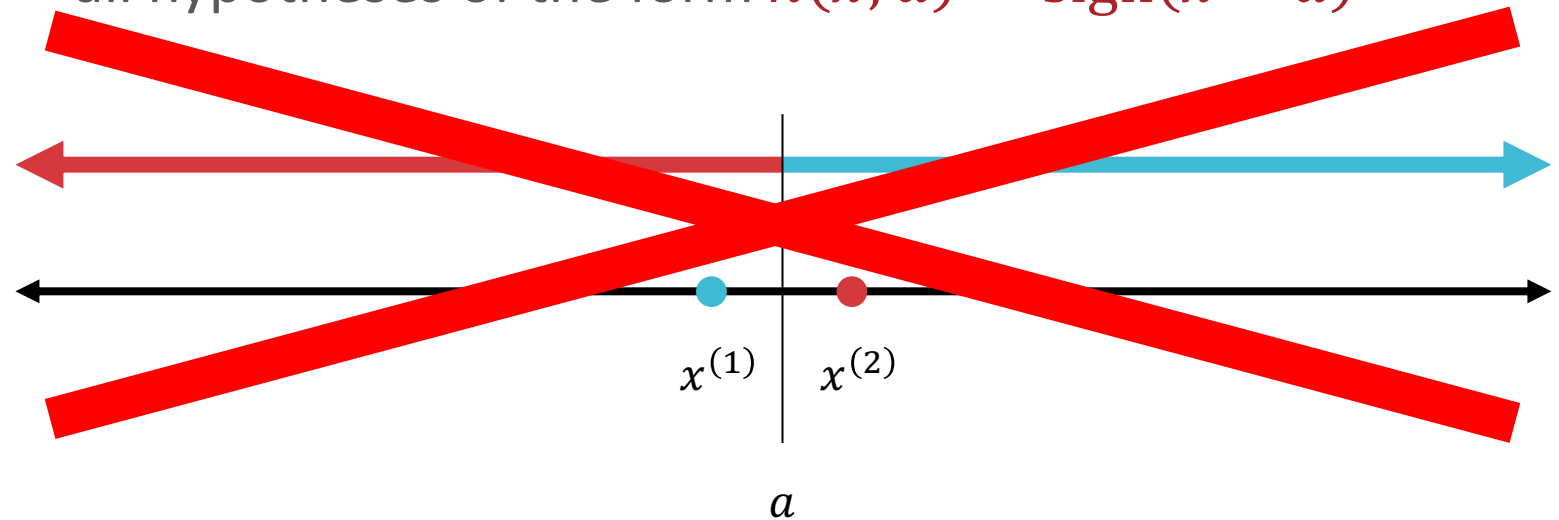
- What is $d_{VC}(\mathcal{H})$?
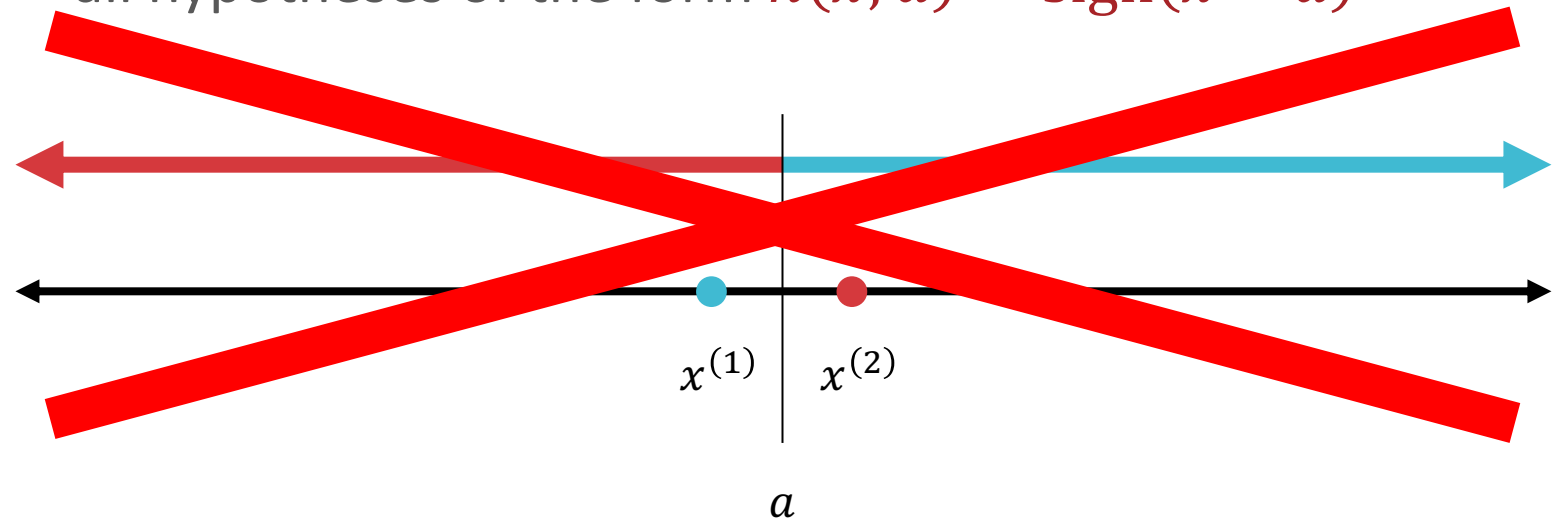
VC-Dimension: Example

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



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

$a$

- 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 $h(x; a) = \text{sign}(x - a)$



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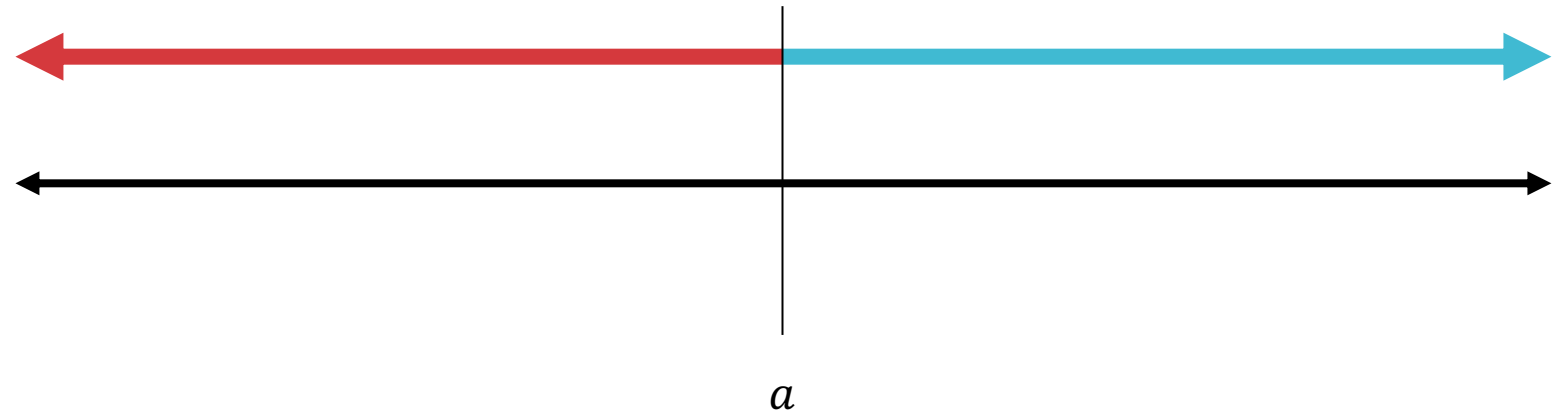
$a$

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

# VC-Dimension: Example

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



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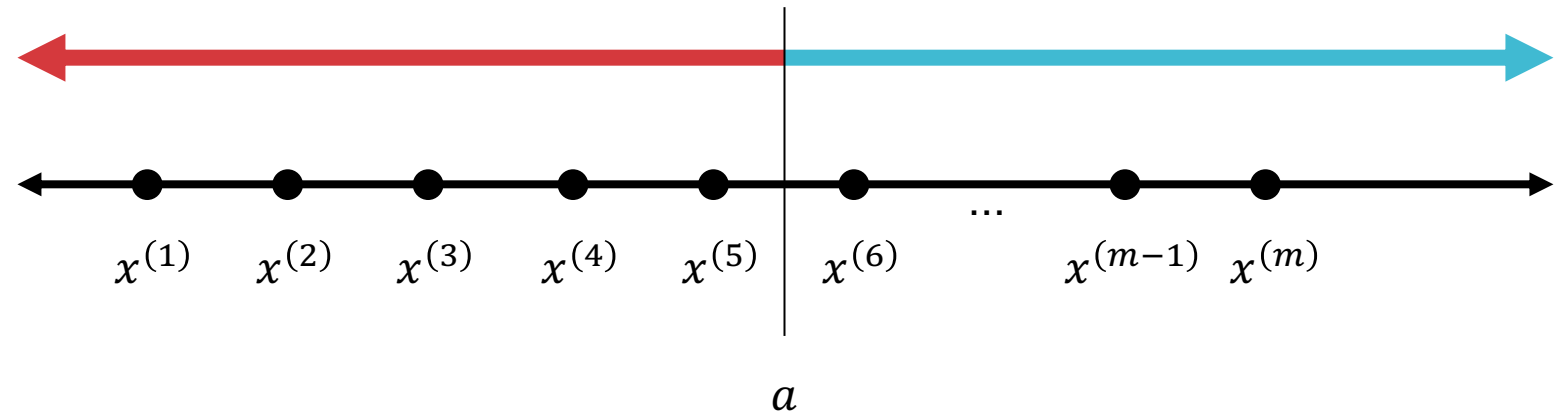
- What is $g_{\mathcal{H}}(m)$?

# VC-Dimension: Example

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



$a$

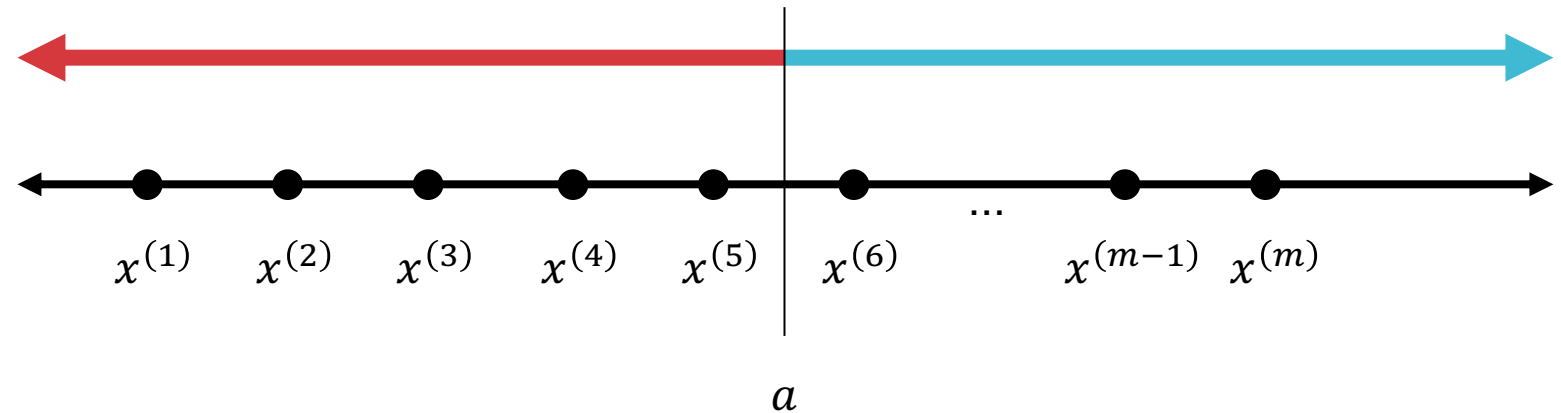- What is $g_{\mathcal{H}}(m)$?
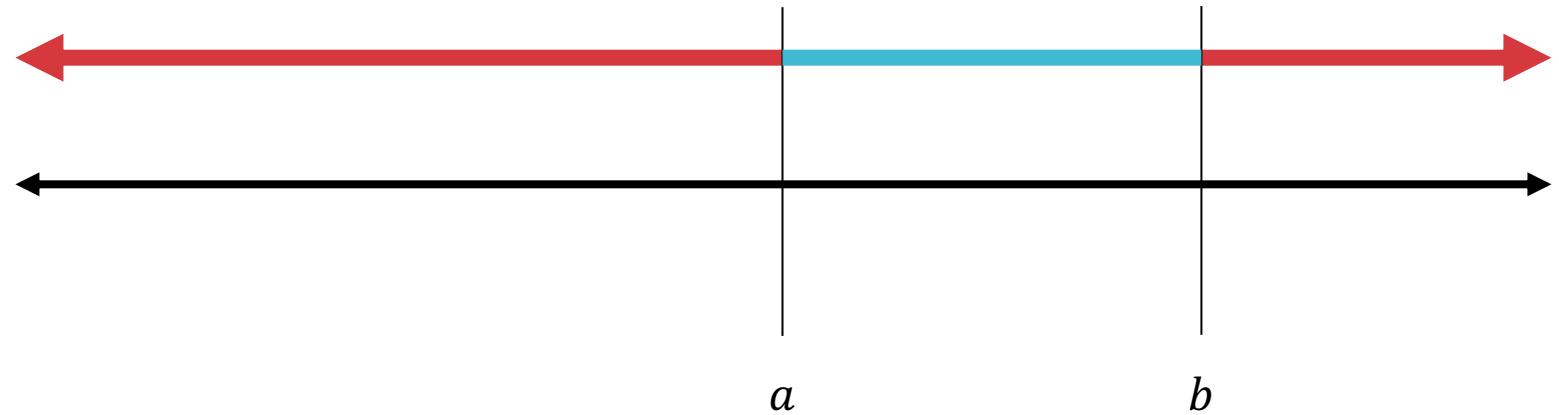
# VC-Dimension: Example

- $x^{(m)} \in \mathbb{R}$ 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 $d_{VC}(\mathrm{H})$ and $g_{\mathrm{H}}(m)$ for 1-dimensional positive intervals?
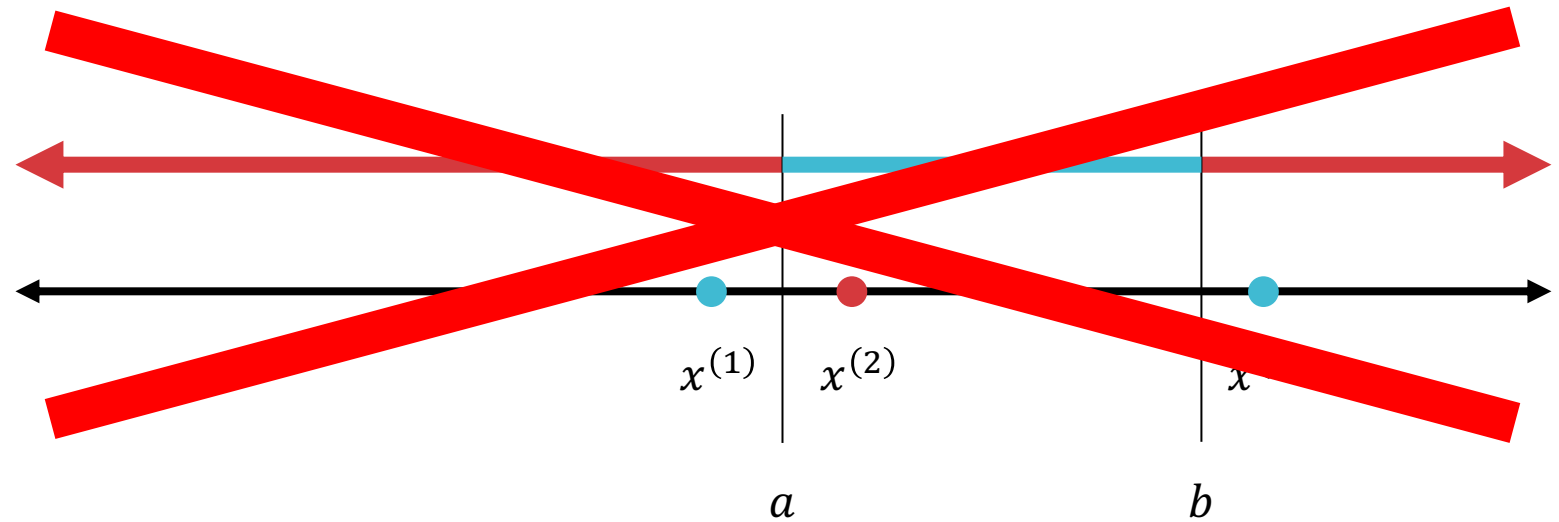
## 1 and $m + 1$

## 2 and $m + 1$

2 and $\frac{1}{2}(m^2 + m + 1)$

3 and $\frac{1}{2}(m^2 + m + 1)$

# VC-Dimension: Example

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



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

$a$ $b$

- 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



$a$            $b$

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

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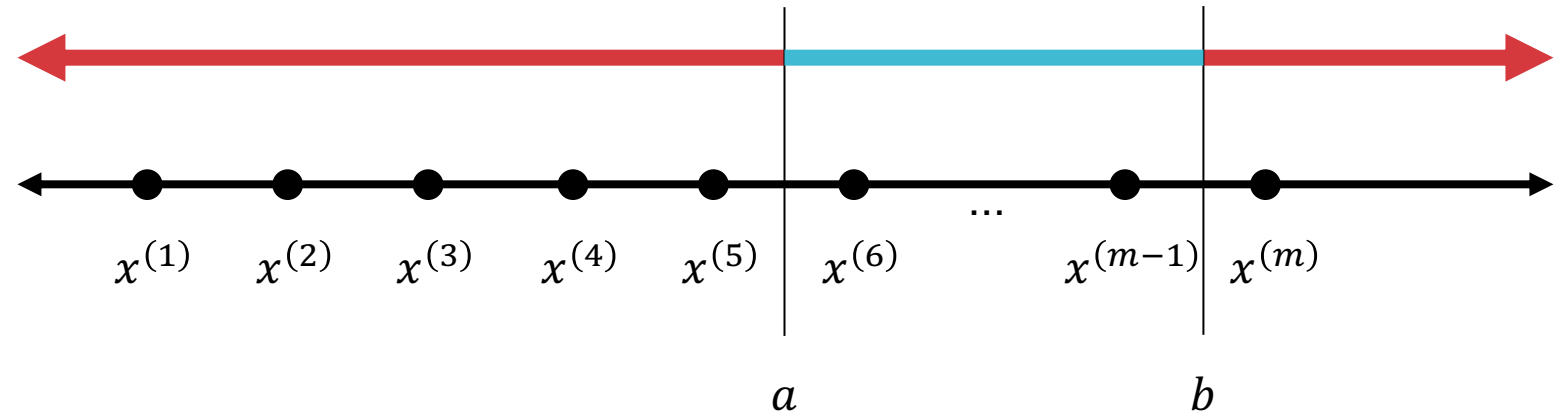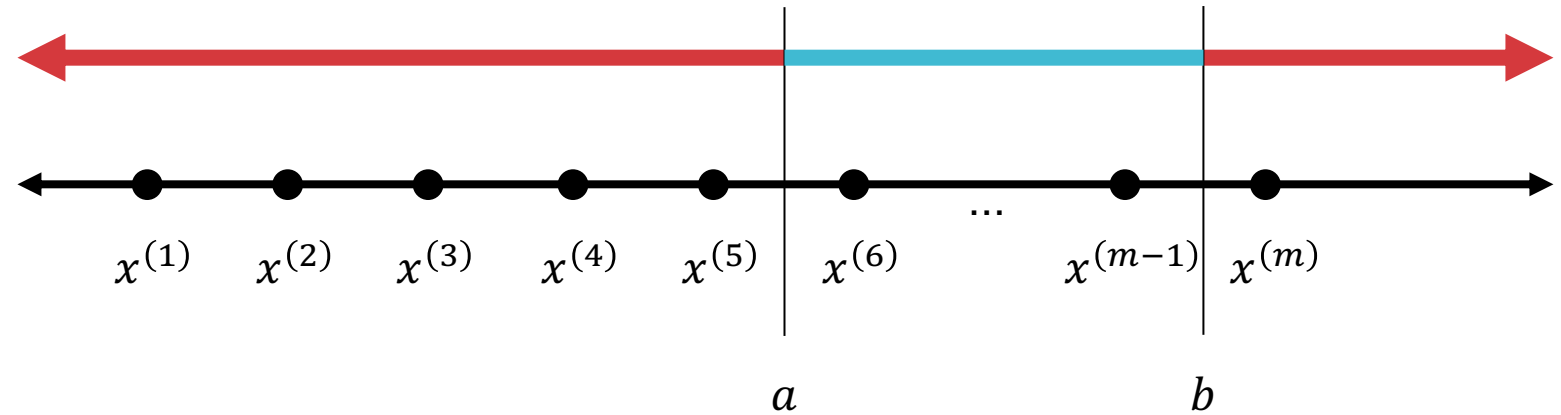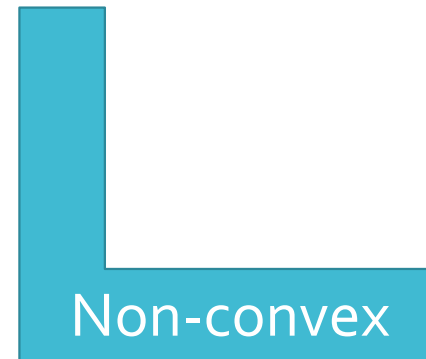
- $x^{(m)} \in \mathbb{R}$ and $\mathcal{H} = $ all 1-dimensional positive intervals



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

Convex

Convex

Non-convex

Non-convex

# 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

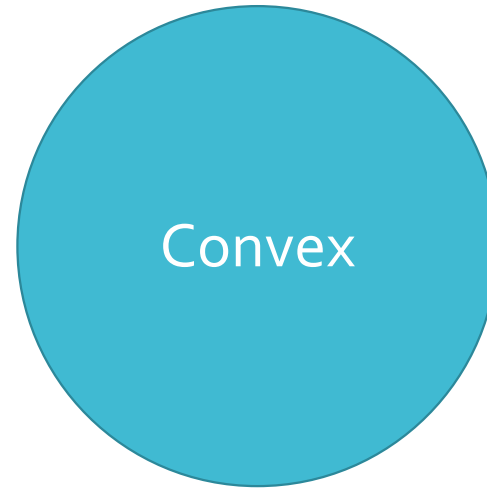- What are $d_{VC}(\mathcal{H})$ and $g_{\mathcal{H}}(M)$?

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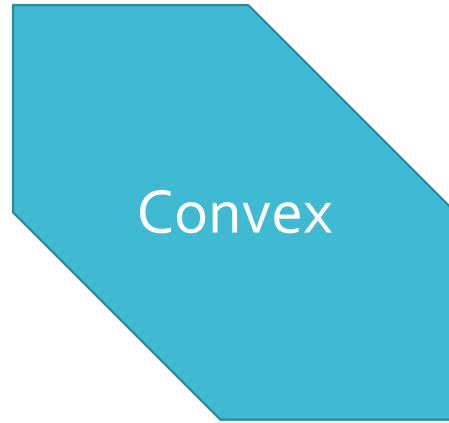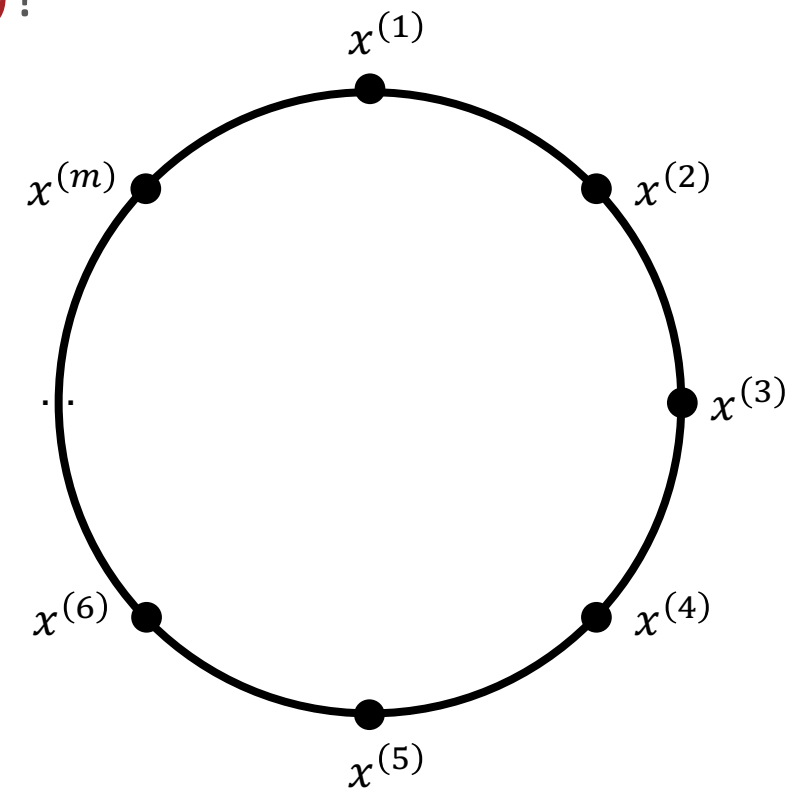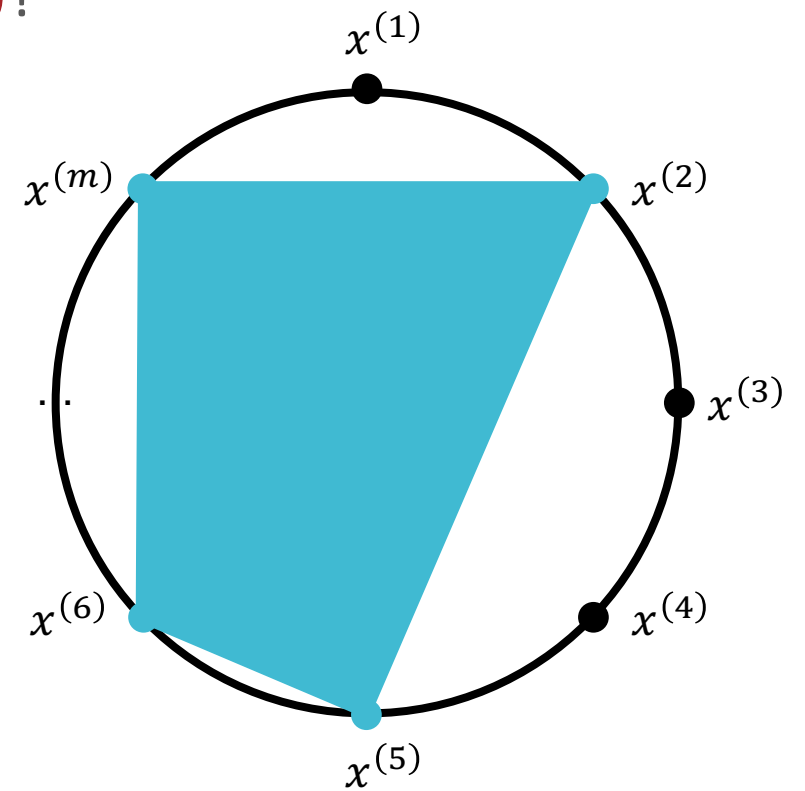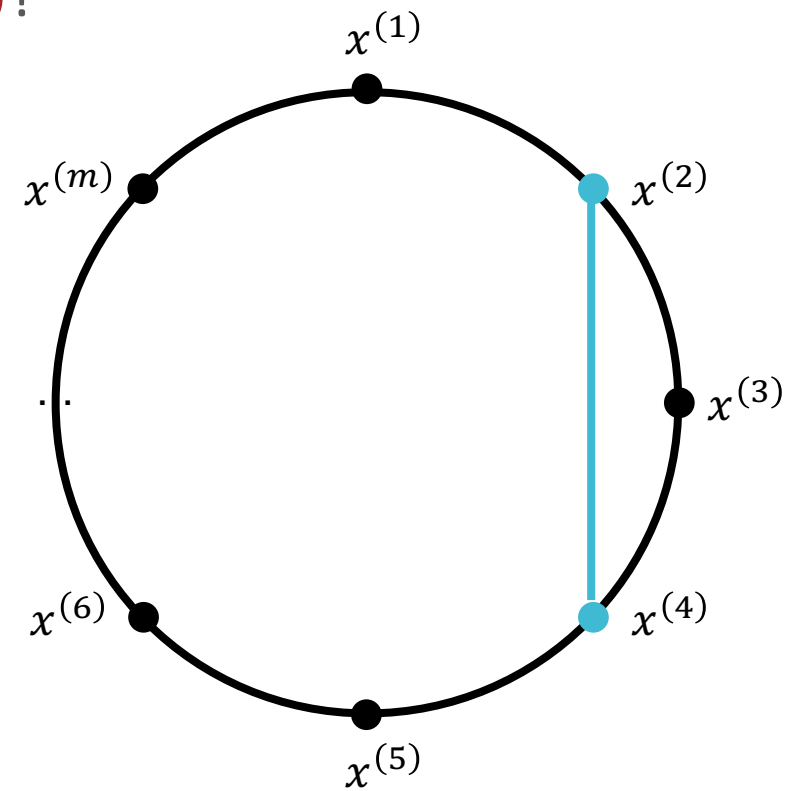
- $x^{(m)} \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional positive convex sets

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

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

# 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

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

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

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