

# 10-301/601: Introduction to Machine Learning Lecture 16 – Learning Theory (Infinite Case)

Matt Gormley & Henry Chai

10/23/24

# Front Matter

- Announcements
  - HW5 released 10/9, due 10/27 at 11:59 PM
  - HW6 released 10/27, due 11/2 at 11:59 PM
  - Discussion post series on Piazza about Societal Impacts of ML
    - “All (substantive) contributions from students in these Piazza posts will be automatically endorsed and count towards the Piazza extra credit portion of your grade”

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

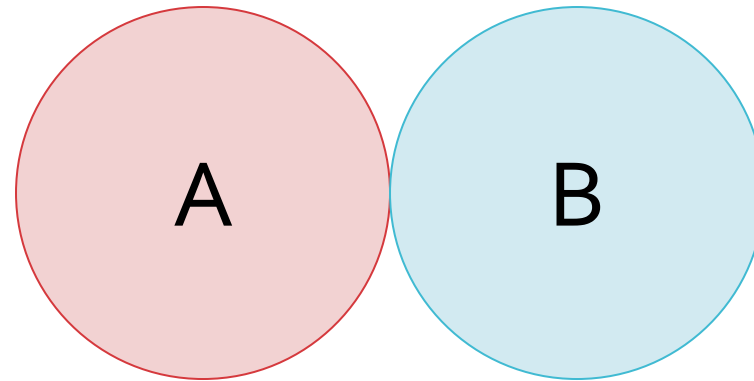
- For a finite hypothesis set  $\mathcal{H}$  and arbitrary distribution  $p^*$ , given a training data set  $S$  s.t.  $|S| = M$ , all  $h \in \mathcal{H}$  have

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

with probability at least  $1 - \delta$ .

# The Union Bound...

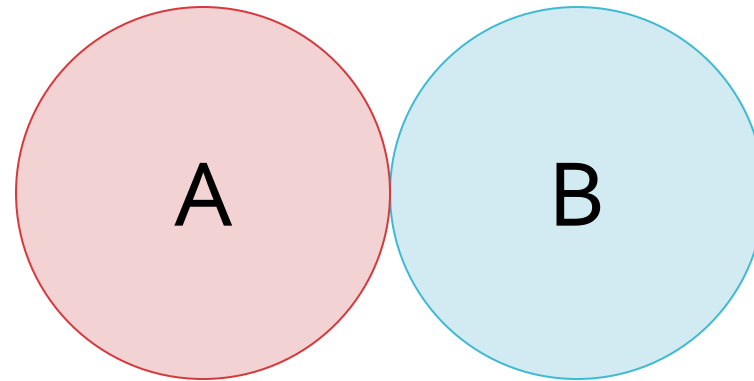
$$P\{A \cup B\} \leq P\{A\} + P\{B\}$$



# The Union Bound is Bad!

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

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

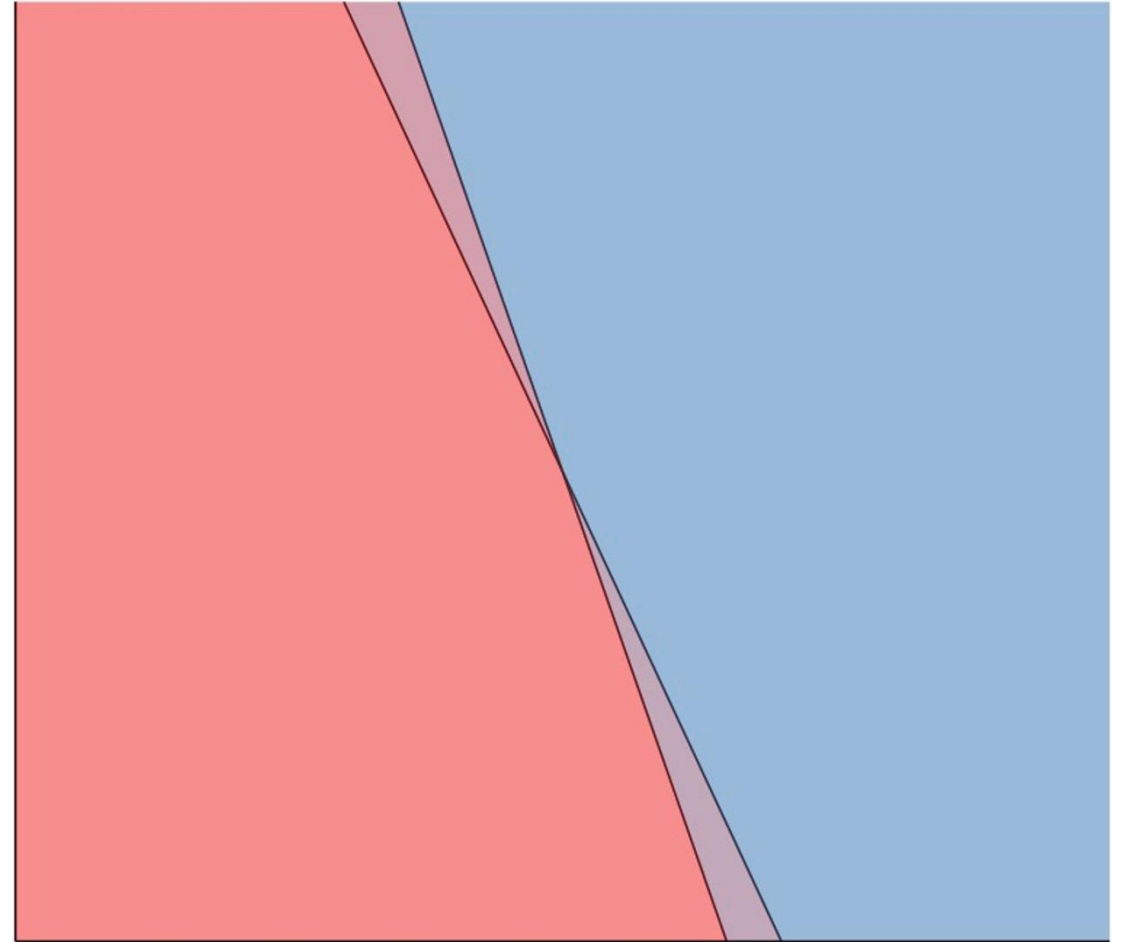


# Intuition

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!

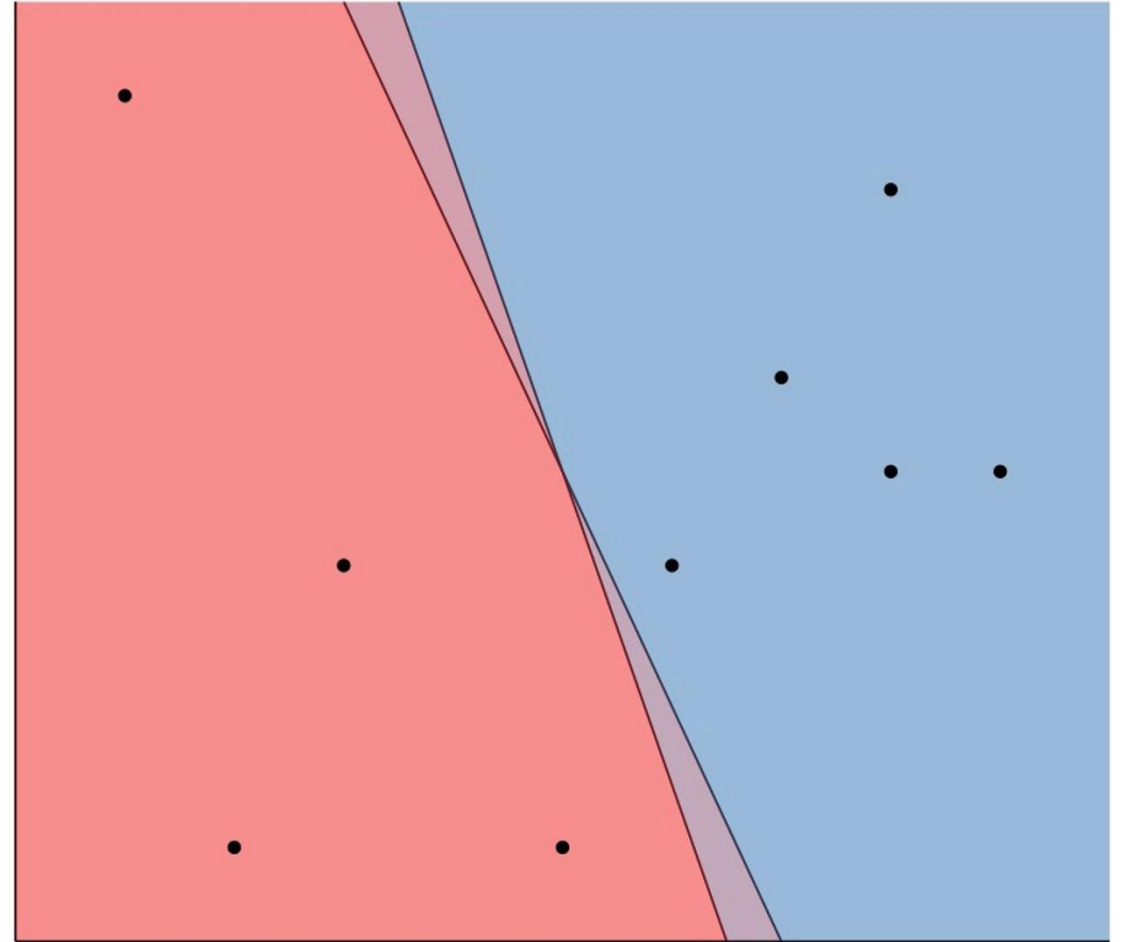


# Intuition

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!



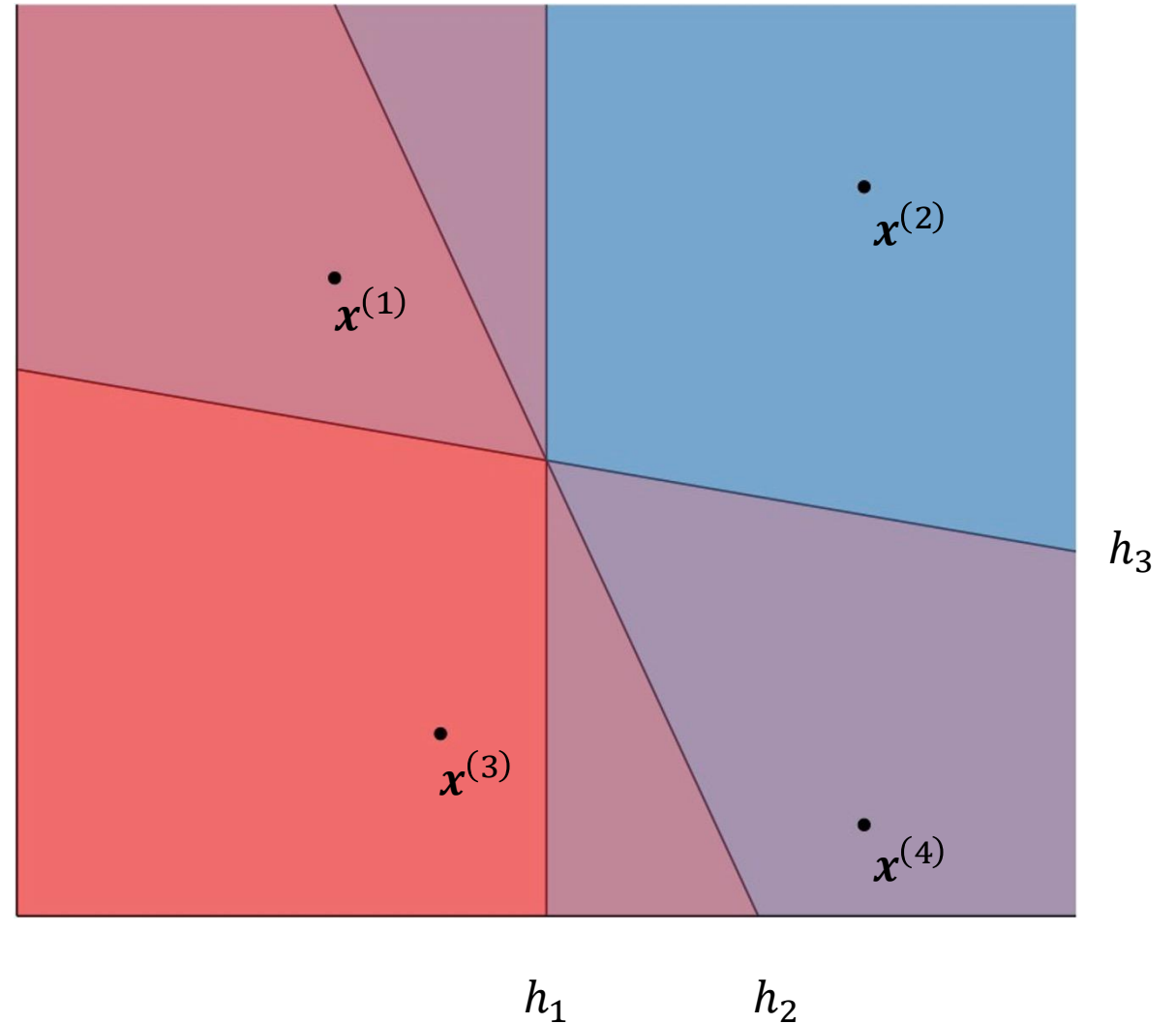
# Labellings

- 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
- 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

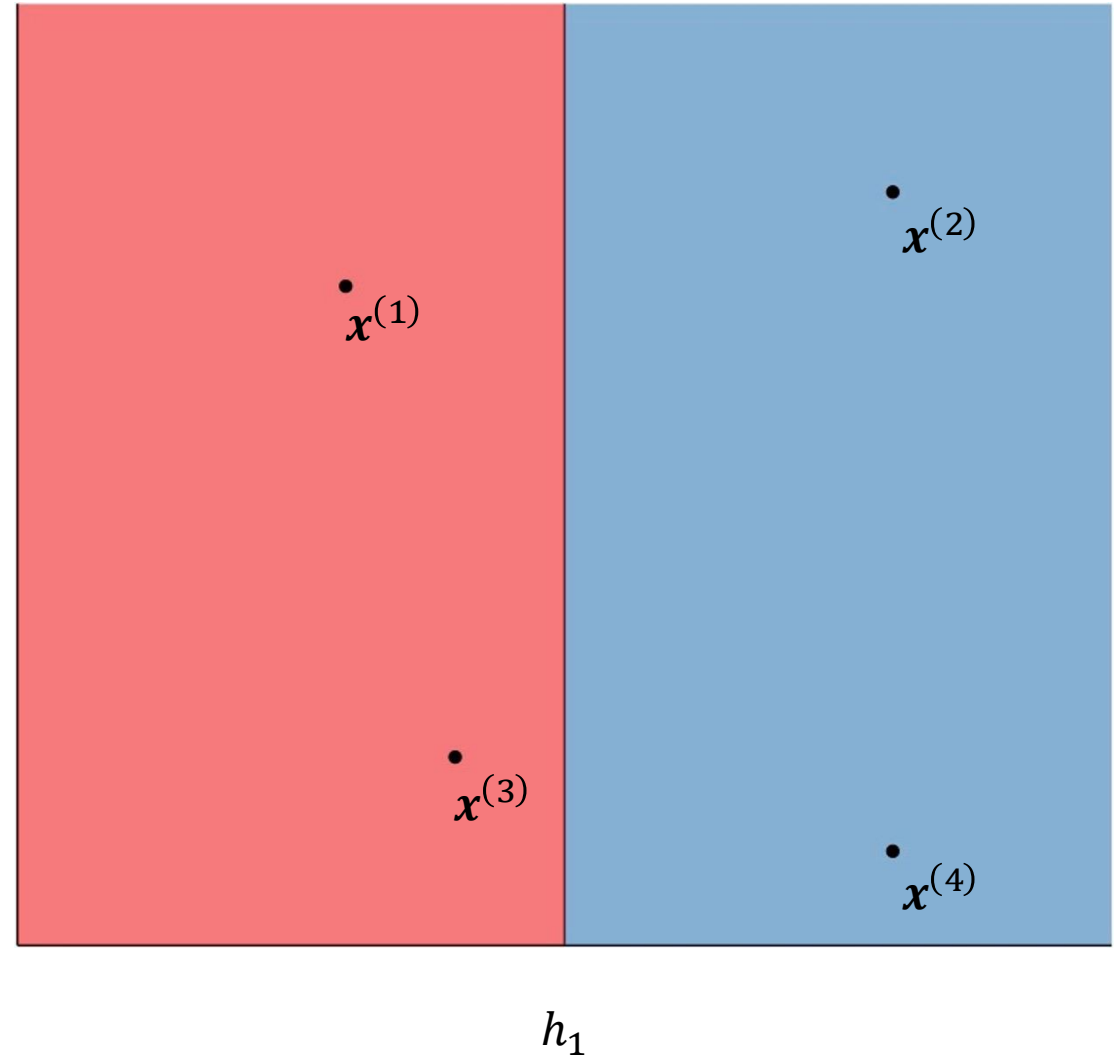
$$\mathcal{H} = \{h_1, h_2, h_3\}$$



# Example: Labellings

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

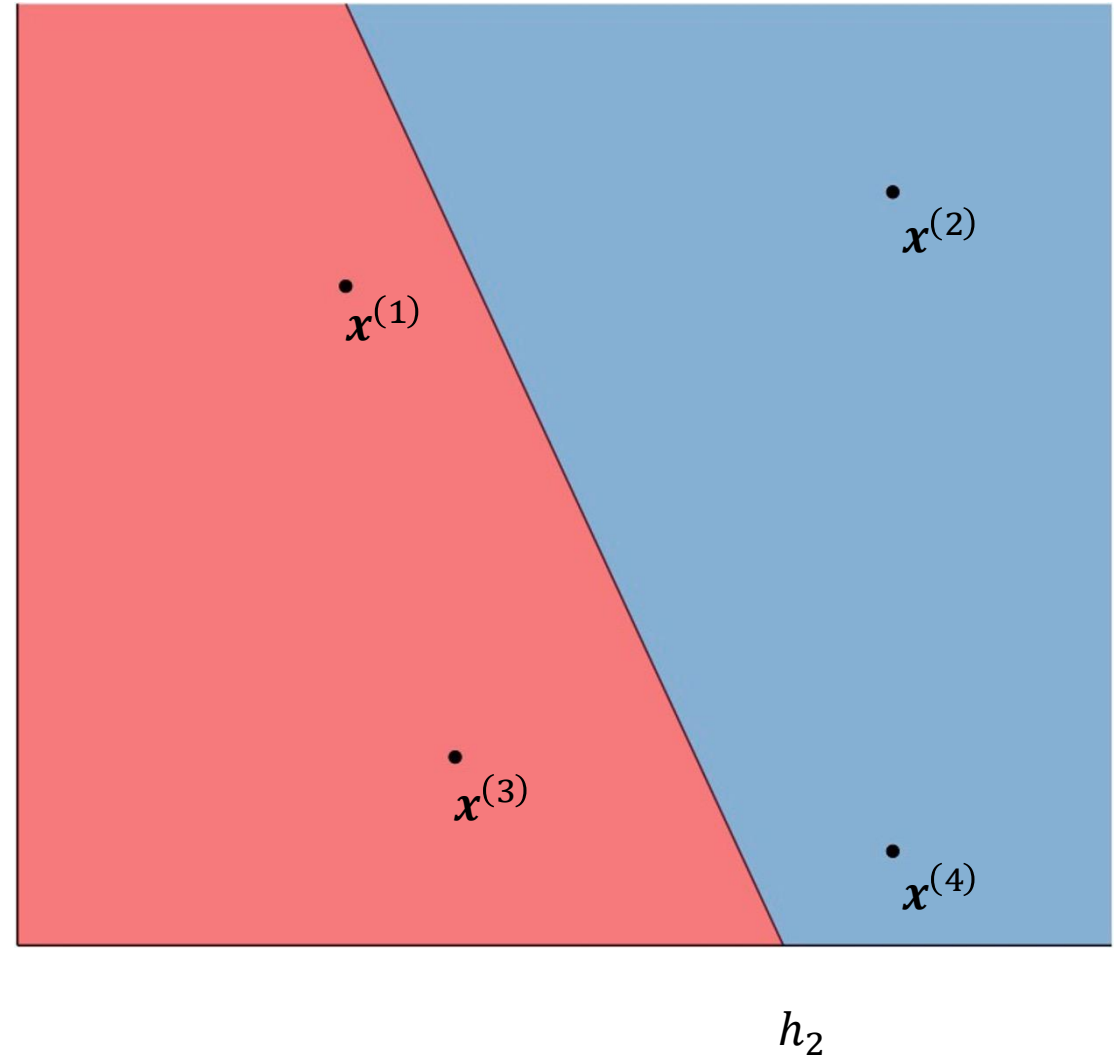
$$\begin{aligned} & (h_1(\mathbf{x}^{(1)}), h_1(\mathbf{x}^{(2)}), h_1(\mathbf{x}^{(3)}), h_1(\mathbf{x}^{(4)})) \\ &= (-1, +1, -1, +1) \end{aligned}$$



# Example: Labellings

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

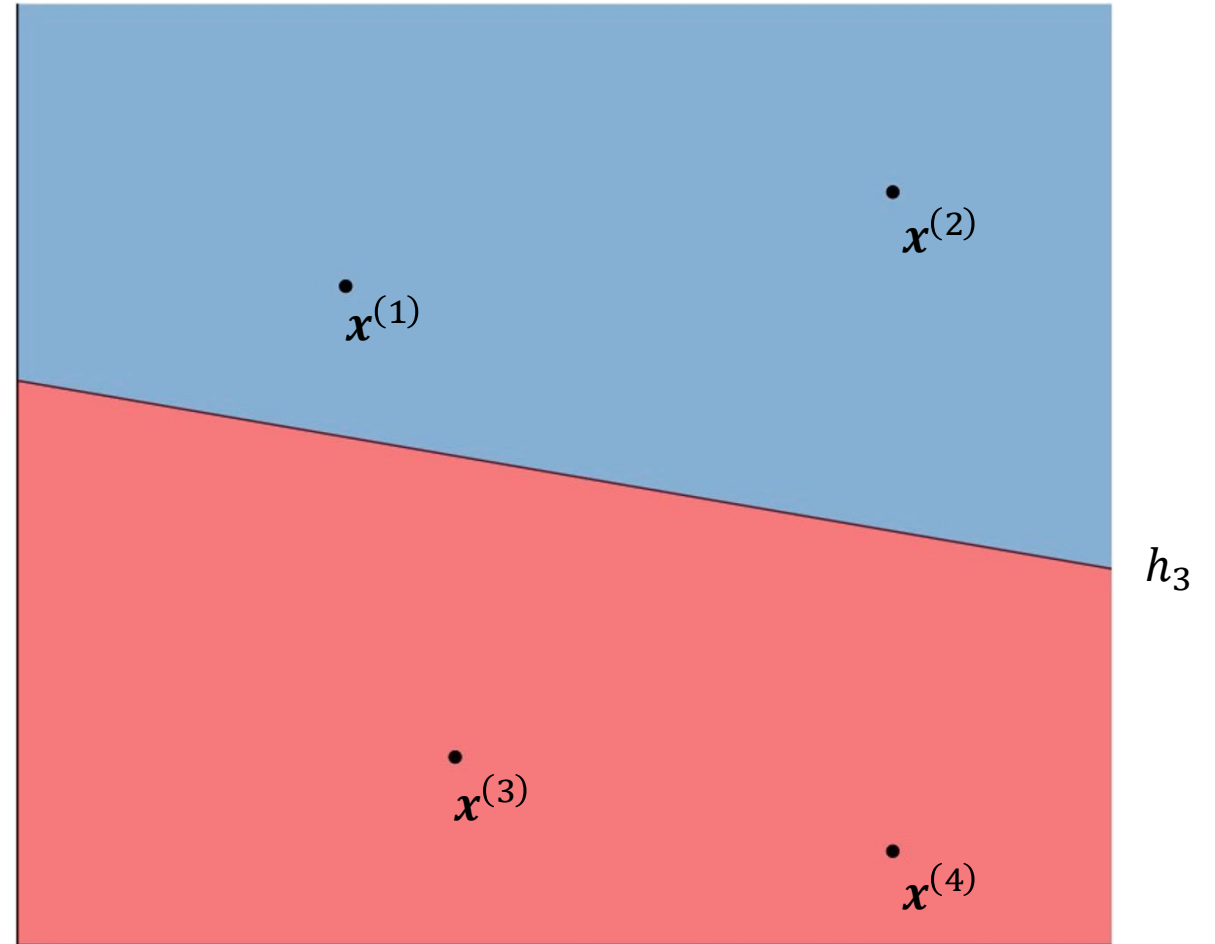
$$\begin{aligned} & \left( h_2(\mathbf{x}^{(1)}), h_2(\mathbf{x}^{(2)}), h_2(\mathbf{x}^{(3)}), h_2(\mathbf{x}^{(4)}) \right) \\ & = (-1, +1, -1, +1) \end{aligned}$$



# Example: Labellings

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

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

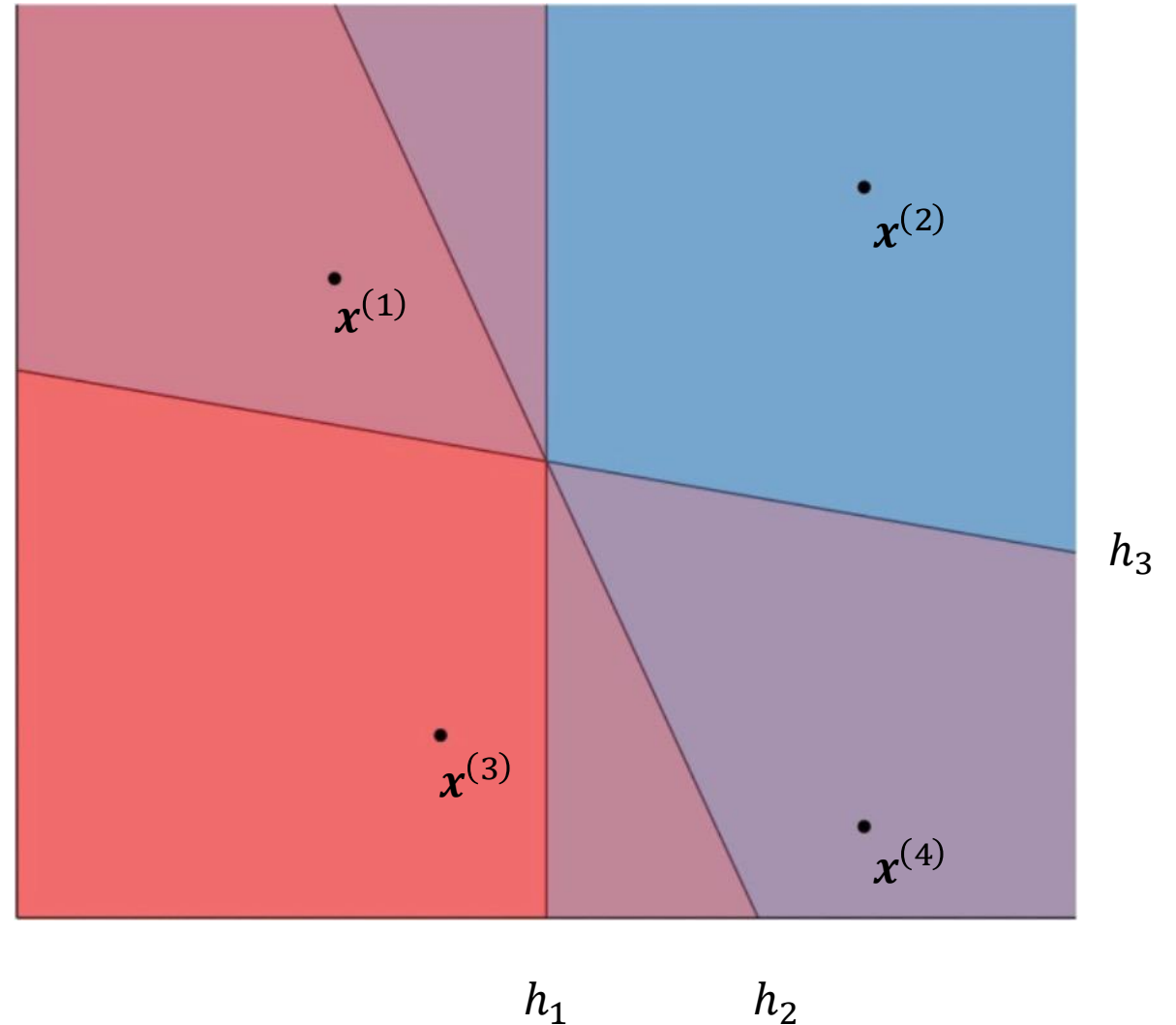


# Example: Labellings

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

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

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

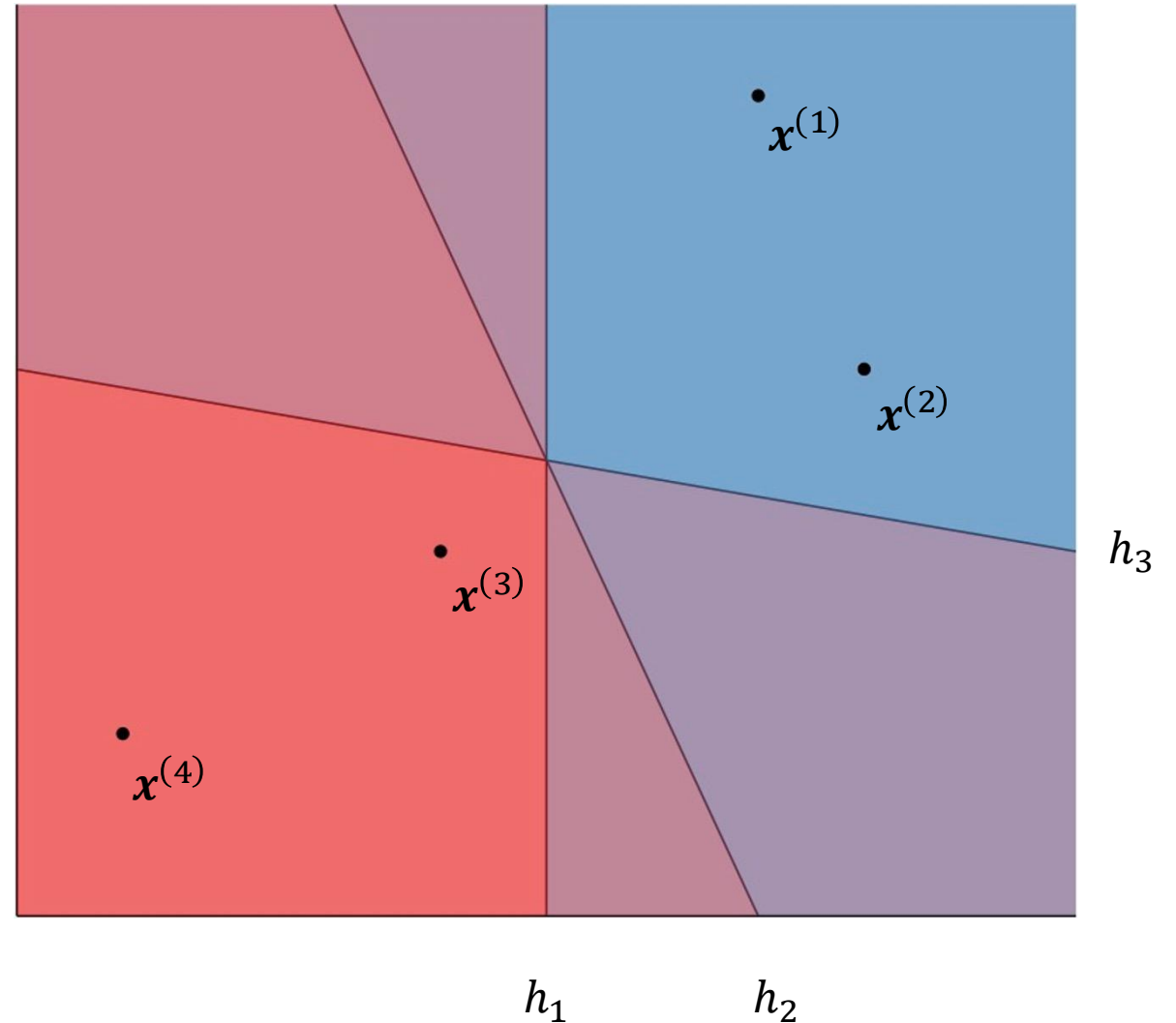


# Example: Labellings

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

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

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

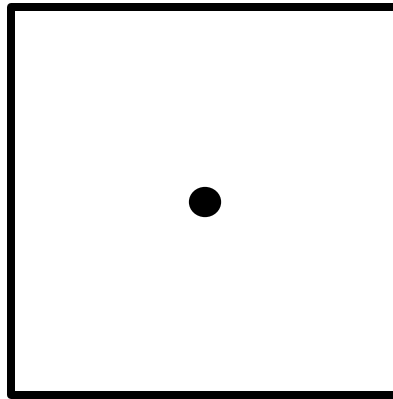


# VC-Dimension

- $\mathcal{H}(S)$  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

- $\mathbf{x} \in \mathbb{R}^2$  and  $\mathcal{H} =$  all 2-dimensional linear separators
- What is  $VC(\mathcal{H})$ ?
  - Can  $\mathcal{H}$  shatter some set of 1 point?

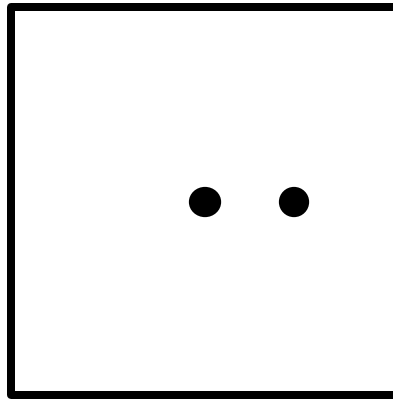


$S$



# VC-Dimension: Example

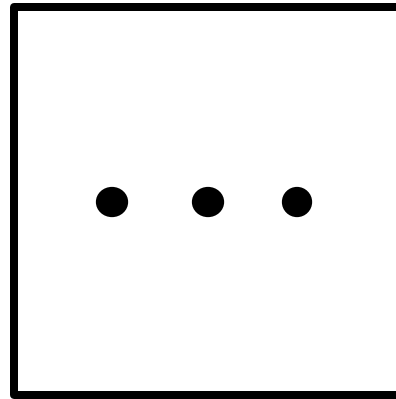
- $\mathbf{x} \in \mathbb{R}^2$  and  $\mathcal{H} =$  all 2-dimensional linear separators
- What is  $VC(\mathcal{H})$ ?
  - Can  $\mathcal{H}$  shatter some set of 1 point?
  - Can  $\mathcal{H}$  shatter some set of 2 points?



$S$

# VC-Dimension: Example

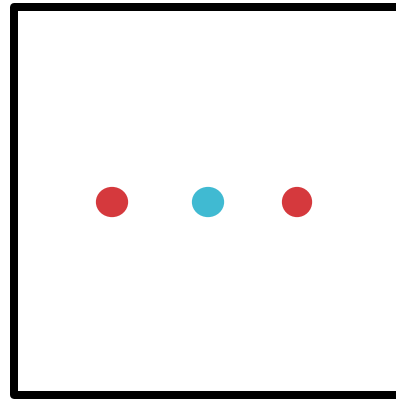
- $\mathbf{x} \in \mathbb{R}^2$  and  $\mathcal{H} =$  all 2-dimensional linear separators
- What is  $VC(\mathcal{H})$ ?
  - Can  $\mathcal{H}$  shatter some set of 1 point?
  - Can  $\mathcal{H}$  shatter some set of 2 points?
  - Can  $\mathcal{H}$  shatter some set of 3 points?



$S$

# VC-Dimension: Example

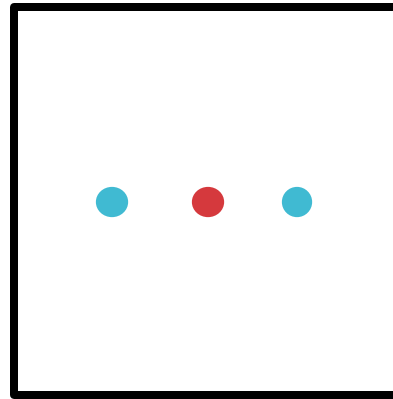
- $\mathbf{x} \in \mathbb{R}^2$  and  $\mathcal{H} =$  all 2-dimensional linear separators
- What is  $VC(\mathcal{H})$ ?
  - Can  $\mathcal{H}$  shatter some set of 1 point?
  - Can  $\mathcal{H}$  shatter some set of 2 points?
  - Can  $\mathcal{H}$  shatter some set of 3 points?



$S$

# VC-Dimension: Example

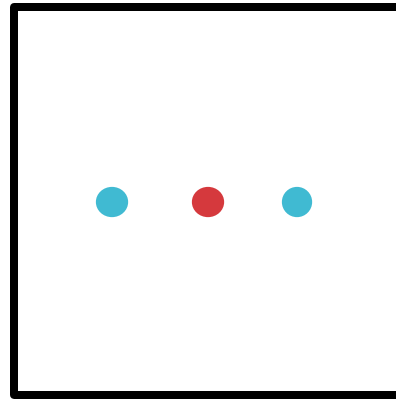
- $\mathbf{x} \in \mathbb{R}^2$  and  $\mathcal{H} =$  all 2-dimensional linear separators
- What is  $VC(\mathcal{H})$ ?
  - Can  $\mathcal{H}$  shatter some set of 1 point?
  - Can  $\mathcal{H}$  shatter some set of 2 points?
  - Can  $\mathcal{H}$  shatter some set of 3 points?



$S$

# VC-Dimension: Example

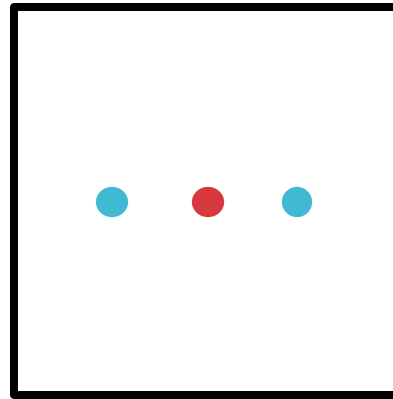
- $\mathbf{x} \in \mathbb{R}^2$  and  $\mathcal{H} =$  all 2-dimensional linear separators
- What is  $VC(\mathcal{H})$ ?
  - Can  $\mathcal{H}$  shatter some set of 1 point?
  - Can  $\mathcal{H}$  shatter some set of 2 points?
  - Can  $\mathcal{H}$  shatter **some** set of 3 points?



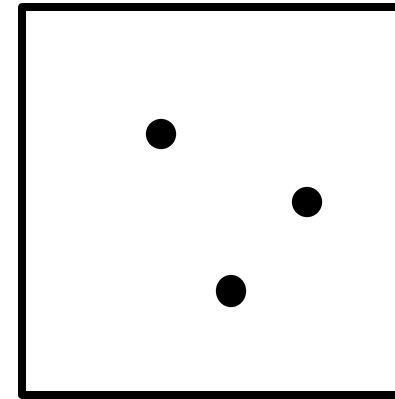
$S$

# VC-Dimension: Example

- $\mathbf{x} \in \mathbb{R}^2$  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?



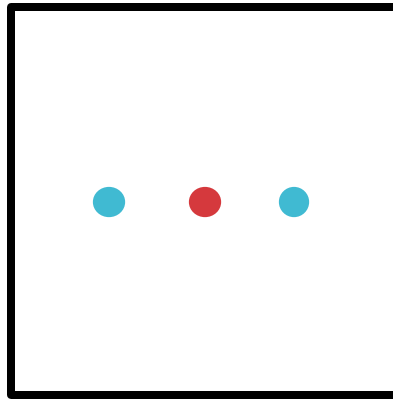
$S_1$



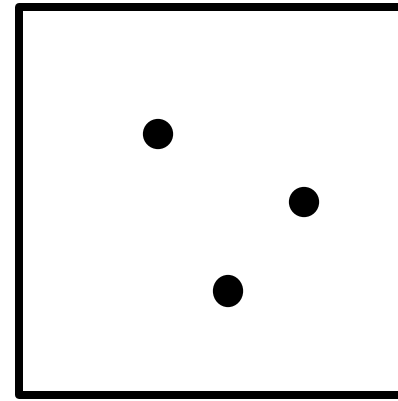
$S_2$

# VC-Dimension: Example

- $\mathbf{x} \in \mathbb{R}^2$  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?



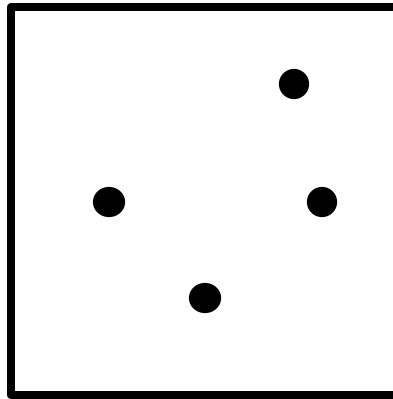
$$|\mathcal{H}(S_1)| = 6$$



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

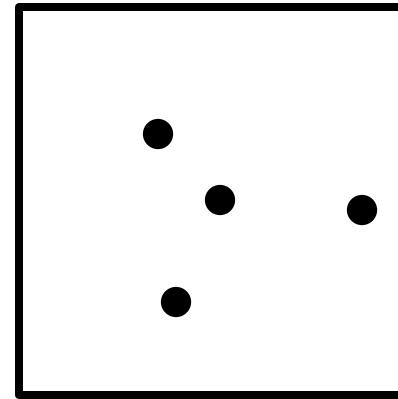
# VC-Dimension: Example

- $\mathbf{x} \in \mathbb{R}^2$  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



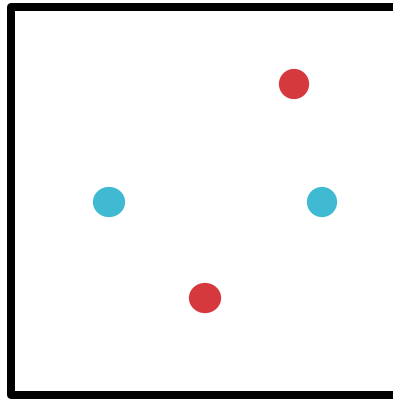
$S_2$

At least one point  
inside the convex hull



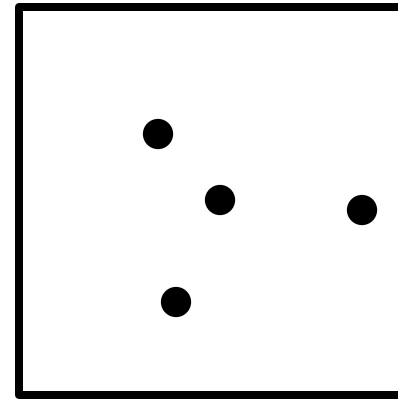
# VC-Dimension: Example

- $\mathbf{x} \in \mathbb{R}^2$  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

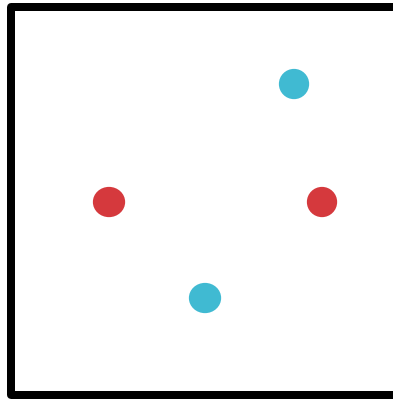


$S_2$

At least one point  
inside the convex hull

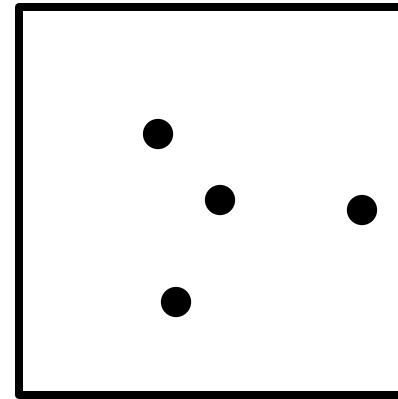
# VC-Dimension: Example

- $\mathbf{x} \in \mathbb{R}^2$  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

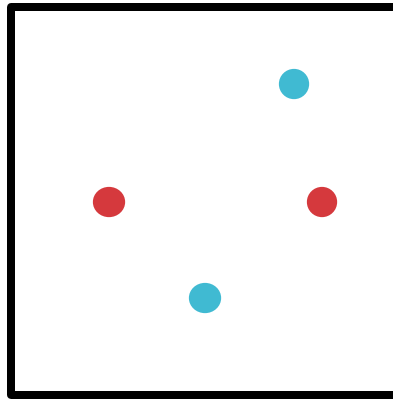


$S_2$

At least one point  
inside the convex hull

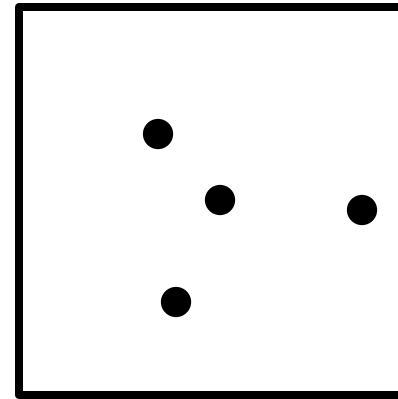
# VC-Dimension: Example

- $\mathbf{x} \in \mathbb{R}^2$  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

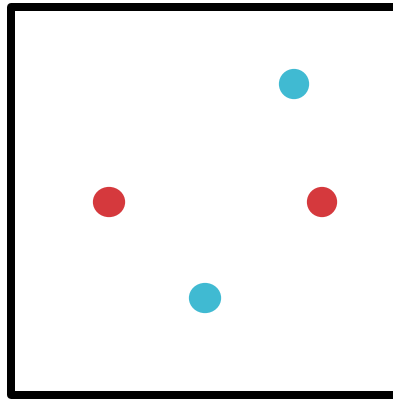


$S_2$

At least one point  
inside the convex hull

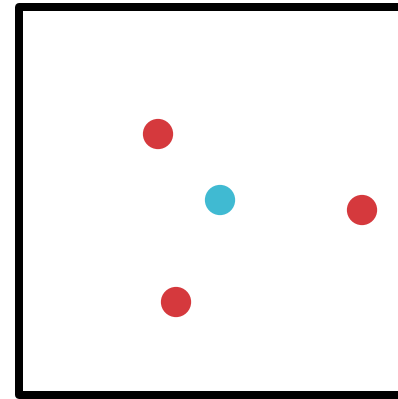
# VC-Dimension: Example

- $\mathbf{x} \in \mathbb{R}^2$  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

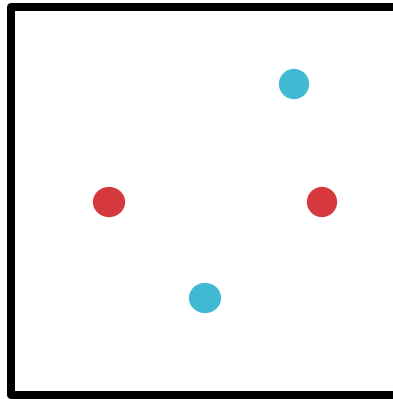


$S_2$

At least one point  
inside the convex hull

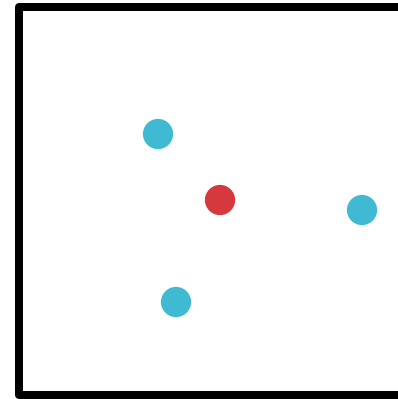
# VC-Dimension: Example

- $\mathbf{x} \in \mathbb{R}^2$  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

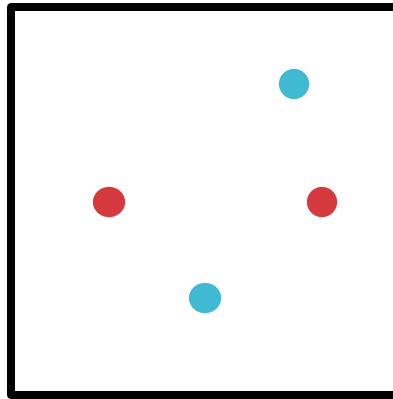


$S_2$

At least one point  
inside the convex hull

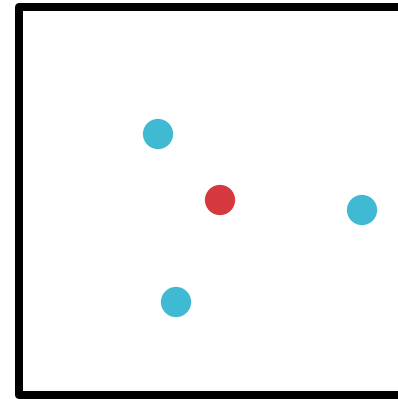
# VC-Dimension: Example

- $\mathbf{x} \in \mathbb{R}^2$  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

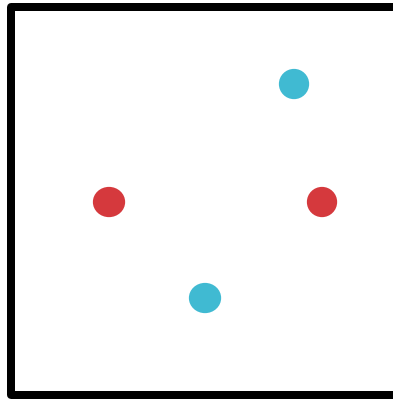


$$|\mathcal{H}(S_2)| = 14$$

At least one point  
inside the convex hull

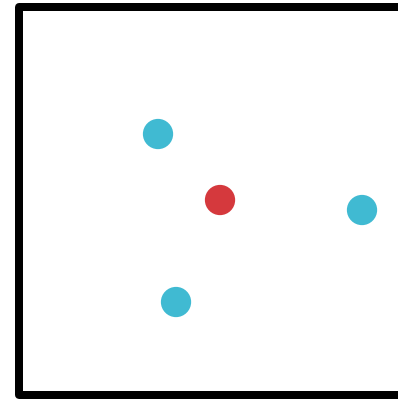
# VC-Dimension: Example

- $\mathbf{x} \in \mathbb{R}^2$  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

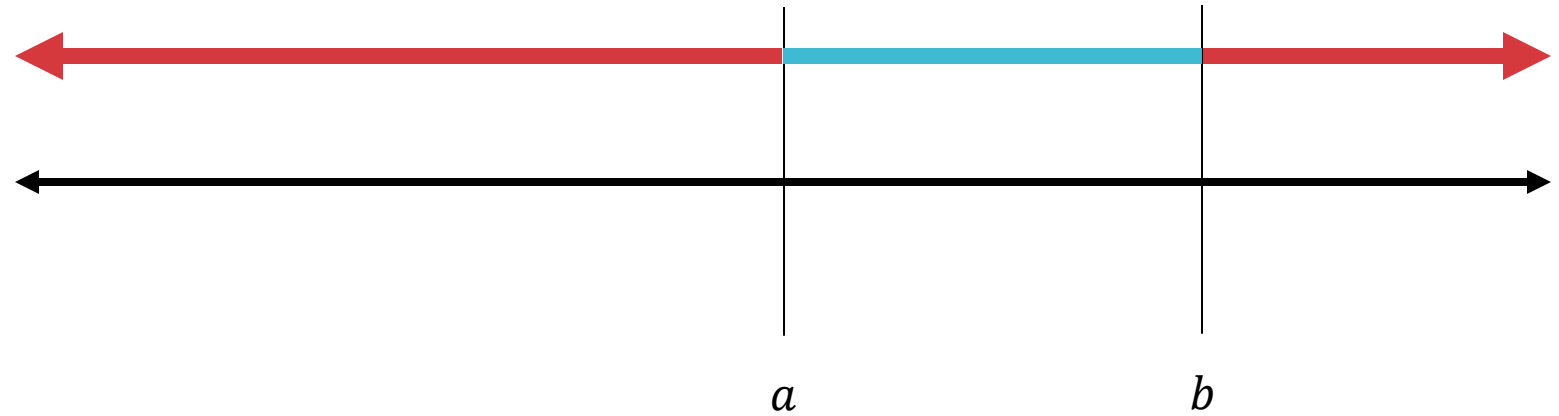
## VC-Dimension: Example

- $\mathbf{x} \in \mathbb{R}^d$  and  $\mathcal{H} =$  all  $d$ -dimensional linear separators
- $VC(\mathcal{H}) = d + 1$



# VC-Dimension: Example

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

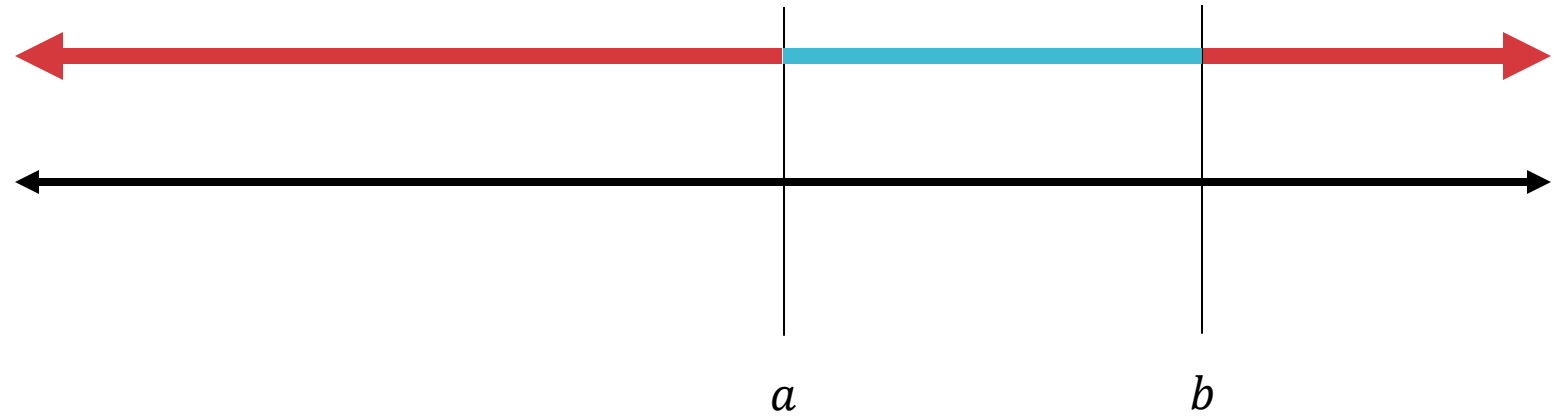


## Poll Question 1:

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

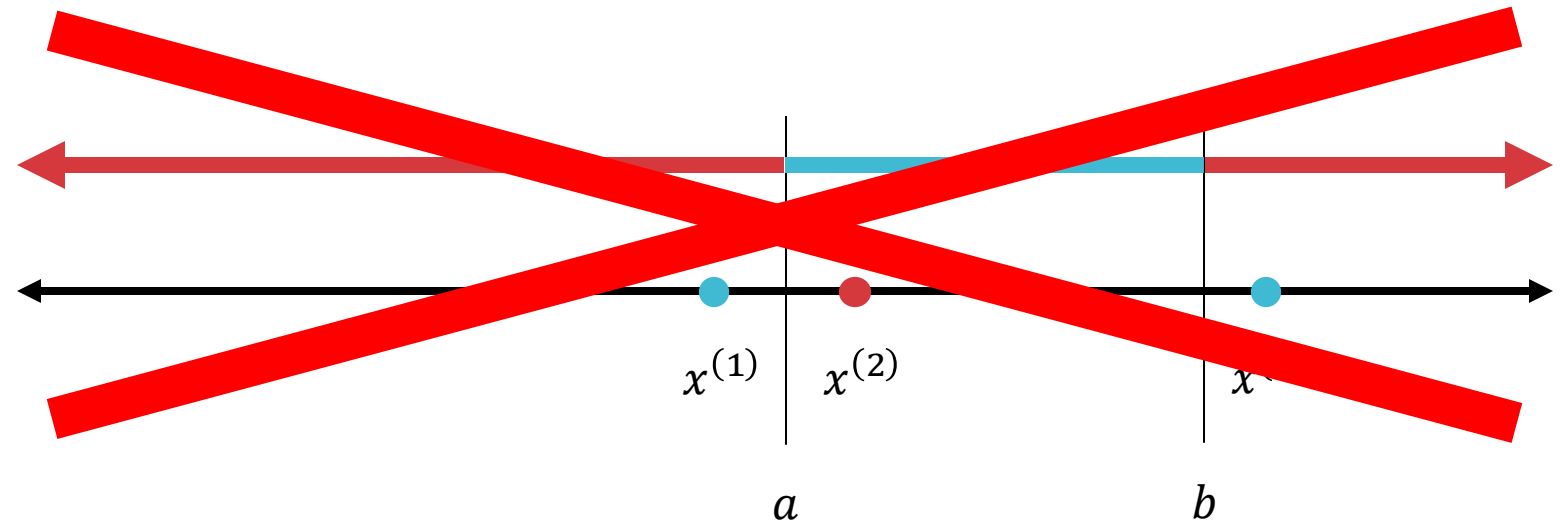
- A. 0
- B. 1
- C. 1.5 (TOXIC)
- D. 2
- E. 3

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



# VC-Dimension: Example

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



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

## 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

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

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 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

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

# 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

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

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

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

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

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

# Learning Theory Learning Objectives

You should be able to...

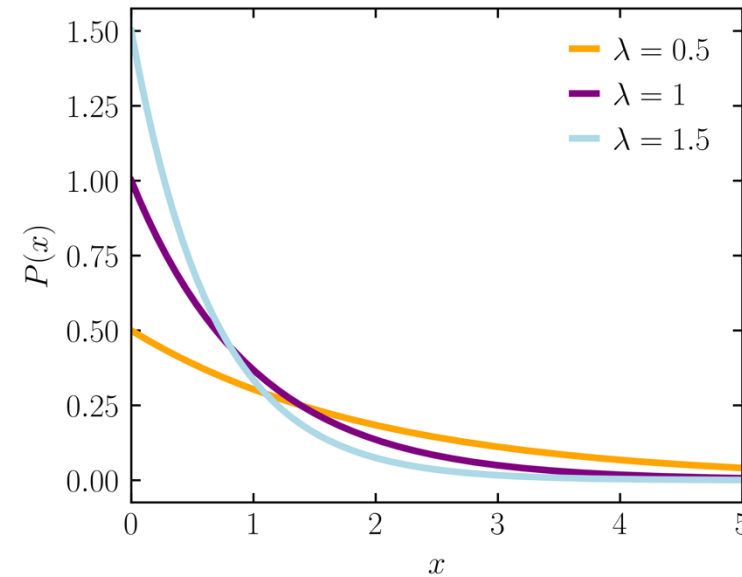
- Identify the properties of a learning setting and assumptions required to ensure low generalization error
- Distinguish true error, train error, test error
- Define PAC and explain what it means to be approximately correct and what occurs with high probability
- Apply sample complexity bounds to real-world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

# Recall: Probabilistic Learning

- Previously:
  - (Unknown) Target function,  $c^*: \mathcal{X} \rightarrow \mathcal{Y}$
  - Classifier,  $h: \mathcal{X} \rightarrow \mathcal{Y}$
  - Goal: find a classifier,  $h$ , that best approximates  $c^*$
- Now:
  - (Unknown) Target *distribution*,  $y \sim p^*(Y|\mathbf{x})$
  - Distribution,  $p(Y|\mathbf{x})$
  - Goal: find a distribution,  $p$ , that best approximates  $p^*$

# Recall: Maximum Likelihood Estimation (MLE)

- Insight: every valid probability distribution has a finite amount of probability mass as it must sum/integrate to 1
- Idea: set the parameter(s) so that the likelihood of the samples is maximized
- Intuition: assign as much of the (finite) probability mass to the observed data *at the expense of unobserved data*
- Example: the exponential distribution



# Bernoulli Distribution MLE

- A Bernoulli random variable takes value **1** with probability  $\phi$  and value **0** with probability  $1 - \phi$

- The pmf of the Bernoulli distribution is

$$p(x|\phi) = \phi^x(1 - \phi)^{1-x}$$

# Coin Flipping MLE

- A Bernoulli random variable takes value **1** (or heads) with probability  $\phi$  and value **0** (or tails) with probability  $1 - \phi$

- The pmf of the Bernoulli distribution is

$$p(x|\phi) = \phi^x(1 - \phi)^{1-x}$$

- Given  $N$  iid samples  $\{x^{(1)}, \dots, x^{(N)}\}$ , the log-likelihood is

$$\ell(\phi) = \sum_{i=1}^N \log p(x^{(i)}|\phi) = \sum_{i=1}^N \log \phi^{x^{(i)}}(1 - \phi)^{1-x^{(i)}}$$

$$= \sum_{i=1}^N x^{(i)} \log \phi + (1 - x^{(i)}) \log(1 - \phi)$$

$$= N_1 \log \phi + N_0 \log(1 - \phi)$$

- where  $N_1$  is the number of **1**'s in  $\{x^{(1)}, \dots, x^{(N)}\}$  and  $N_0$  is the number of **0**'s

# Coin Flipping MLE

- A Bernoulli random variable takes value **1** (or heads) with probability  $\phi$  and value **0** (or tails) with probability  $1 - \phi$

- The pmf of the Bernoulli distribution is

$$p(x|\phi) = \phi^x(1 - \phi)^{1-x}$$

- The partial derivative of the log-likelihood is

$$\frac{\partial \ell}{\partial \phi} = \frac{N_1}{\phi} - \frac{N_0}{1 - \phi}$$

- where  $N_1$  is the number of **1**'s in  $\{x^{(1)}, \dots, x^{(N)}\}$  and  $N_0$  is the number of **0**'s

# Coin Flipping MLE

- A Bernoulli random variable takes value **1** (or heads) with probability  $\phi$  and value **0** (or tails) with probability  $1 - \phi$

- The pmf of the Bernoulli distribution is

$$p(x|\phi) = \phi^x(1 - \phi)^{1-x}$$

- The partial derivative of the log-likelihood is

$$\frac{N_1}{\hat{\phi}} - \frac{N_0}{1 - \hat{\phi}} = 0 \rightarrow \frac{N_1}{\hat{\phi}} = \frac{N_0}{1 - \hat{\phi}}$$

$$\rightarrow N_1(1 - \hat{\phi}) = N_0\hat{\phi} \rightarrow N_1 = \hat{\phi}(N_0 + N_1)$$

$$\rightarrow \hat{\phi} = \frac{N_1}{N_0 + N_1}$$

- where  $N_1$  is the number of **1**'s in  $\{x^{(1)}, \dots, x^{(N)}\}$  and  $N_0$  is the number of **0**'s



## Poll Question 2:

Go to <https://justflipacoin.com/> and flip the coin 5 times. What is the MLE of your coin?

- A. 0/5
- B. 1/5
- C. 2/5
- D. 3/5
- E.  $\pi/5$  (TOXIC)
- F. 4/5
- G. 5/5

- A Bernoulli random variable takes value **1** (or heads) with probability  $\phi$  and value **0** (or tails) with probability  $1 - \phi$

- The pmf of the Bernoulli distribution is

$$p(x|\phi) = \phi^x(1 - \phi)^{1-x}$$

- The partial derivative of the log-likelihood is

$$\frac{N_1}{\hat{\phi}} - \frac{N_0}{1 - \hat{\phi}} = 0 \rightarrow \frac{N_1}{\hat{\phi}} = \frac{N_0}{1 - \hat{\phi}}$$

$$\rightarrow N_1(1 - \hat{\phi}) = N_0\hat{\phi} \rightarrow N_1 = \hat{\phi}(N_0 + N_1)$$

$$\rightarrow \hat{\phi} = \frac{N_1}{N_0 + N_1}$$

- where  $N_1$  is the number of **1**'s in  $\{x^{(1)}, \dots, x^{(N)}\}$  and  $N_0$  is the number of **0**'s

# Maximum a Posteriori (MAP) Estimation

- Insight: sometimes we have *prior* information we want to incorporate into parameter estimation
- Idea: use Bayes rule to reason about the *posterior* distribution over the parameters

- MLE finds  $\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathcal{D}|\theta)$

- MAP finds  $\hat{\theta} = \operatorname{argmax}_{\theta} p(\theta|\mathcal{D})$   
 $= \operatorname{argmax}_{\theta} p(\mathcal{D}|\theta)p(\theta)/p(\mathcal{D})$

$$= \operatorname{argmax}_{\theta} p(\mathcal{D}|\theta)p(\theta)$$

likelihood                  prior

$$= \operatorname{argmax}_{\theta} \underbrace{\log p(\mathcal{D}|\theta) + \log p(\theta)}_{\text{log-posterior}}$$

Okay, but how on earth do we pick a prior?

- Insight: sometimes we have *prior* information we want to incorporate into parameter estimation
- Idea: use Bayes rule to reason about the *posterior* distribution over the parameters

- MLE finds  $\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathcal{D}|\theta)$

- MAP finds  $\hat{\theta} = \operatorname{argmax}_{\theta} p(\theta|\mathcal{D})$   
 $= \operatorname{argmax}_{\theta} p(\mathcal{D}|\theta)p(\theta)/p(\mathcal{D})$

$$= \operatorname{argmax}_{\theta} \underbrace{\log p(\mathcal{D}|\theta) + \log p(\theta)}_{\text{log-posterior}}$$

likelihood                      prior

# Coin Flipping MAP

- A Bernoulli random variable takes value **1** (or heads) with probability  $\phi$  and value **0** (or tails) with probability  $1 - \phi$
- The pmf of the Bernoulli distribution is

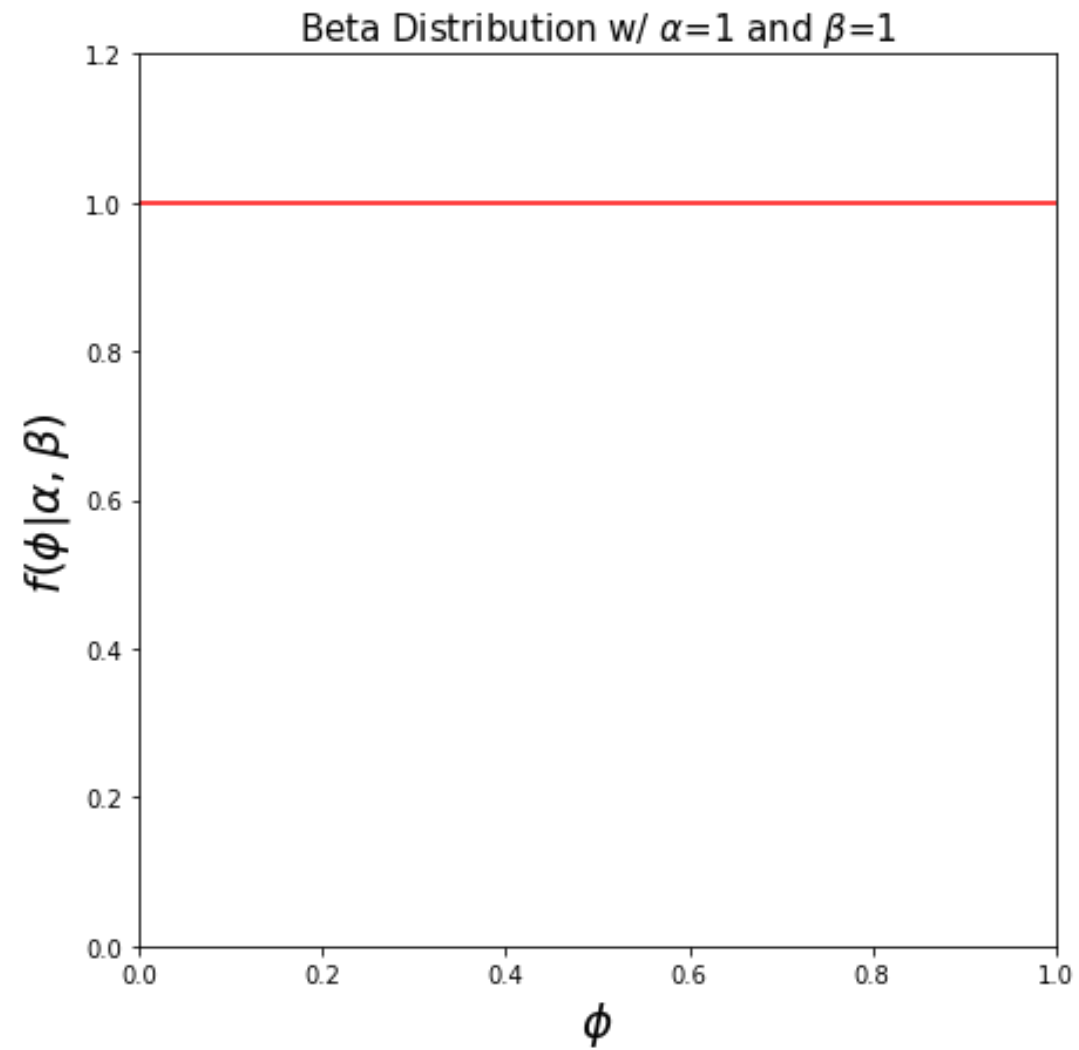
$$p(x|\phi) = \phi^x (1 - \phi)^{1-x}$$

- Assume a Beta prior over the parameter  $\phi$ , which has pdf

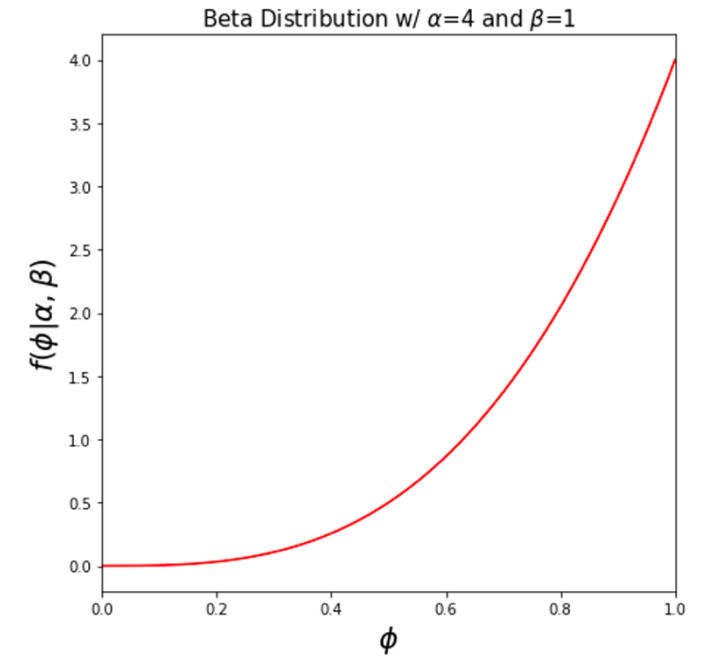
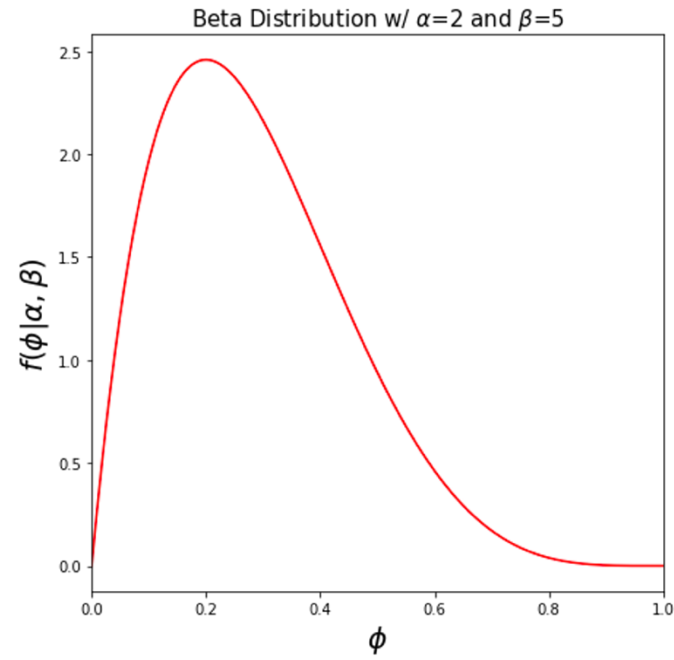
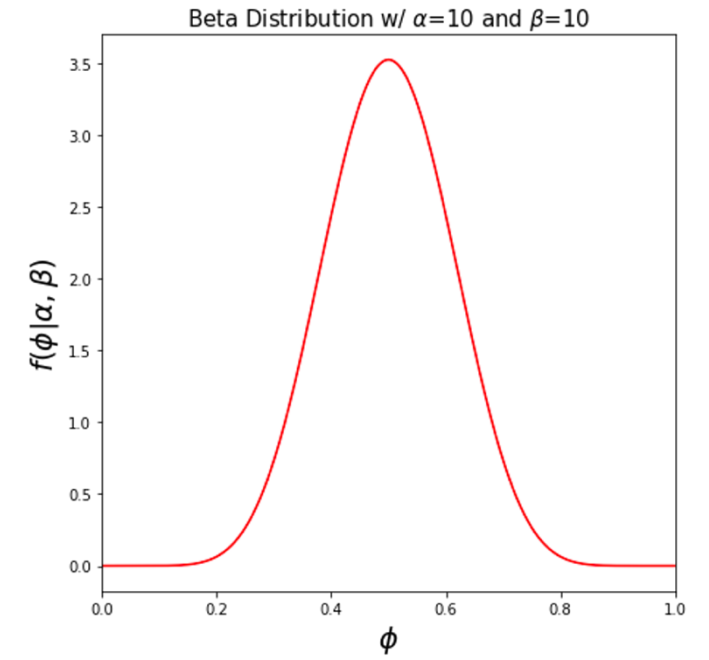
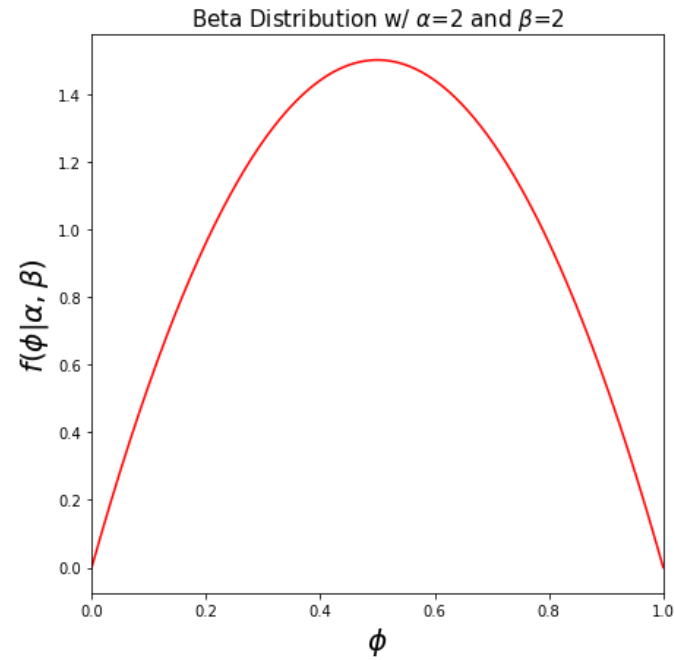
$$f(\phi|\alpha, \beta) = \frac{\phi^{\alpha-1} (1 - \phi)^{\beta-1}}{B(\alpha, \beta)}$$

where  $B(\alpha, \beta) = \int_0^1 \phi^{\alpha-1} (1 - \phi)^{\beta-1} d\phi$  is a normalizing constant to ensure the distribution integrates to **1**

# Beta Distribution



# Beta Distribution



Why use this strange looking Beta prior?

The Beta distribution is the *conjugate prior* for the Bernoulli distribution!

- A Bernoulli random variable takes value **1** (or heads) with probability  $\phi$  and value **0** (or tails) with probability  $1 - \phi$
- The pmf of the Bernoulli distribution is

$$p(x|\phi) = \phi^x (1 - \phi)^{1-x}$$

- Assume a Beta prior over the parameter  $\phi$ , which has pdf

$$f(\phi|\alpha, \beta) = \frac{\phi^{\alpha-1} (1 - \phi)^{\beta-1}}{B(\alpha, \beta)}$$

where  $B(\alpha, \beta) = \int_0^1 \phi^{\alpha-1} (1 - \phi)^{\beta-1} d\phi$  is a normalizing constant to ensure the distribution integrates to **1**

# Coin Flipping MAP

- Given  $N$  iid samples  $\{x^{(1)}, \dots, x^{(N)}\}$ , the log-posterior is

$$\begin{aligned}\ell(\phi) &= \log f(\phi|\alpha, \beta) + \sum_{n=1}^N \log p(x^{(n)}|\phi) \\ &= \log \frac{\phi^{\alpha-1} (1-\phi)^{\beta-1}}{B(\alpha, \beta)} + \sum_{n=1}^N \log \phi^{x^{(n)}} (1-\phi)^{1-x^{(n)}} \\ &= (\alpha-1) \log \phi + (\beta-1) \log(1-\phi) - \log B(\alpha, \beta) \\ &\quad + \sum_{n=1}^N x^{(n)} \log \phi + (1-x^{(n)}) \log(1-\phi) \\ &= (\alpha-1 + N_1) \log \phi + (\beta-1 + N_0) \log(1-\phi) \\ &\quad - \log B(\alpha, \beta)\end{aligned}$$



# Coin Flipping MAP

- Given  $N$  iid samples  $\{x^{(1)}, \dots, x^{(N)}\}$ , the partial derivative of the log-posterior is

$$\frac{\partial \ell}{\partial \phi} = \frac{(\alpha - 1 + N_1)}{\phi} - \frac{(\beta - 1 + N_0)}{1 - \phi}$$

$\vdots$

$$\rightarrow \hat{\phi}_{MAP} = \frac{(\alpha - 1 + N_1)}{(\beta - 1 + N_0) + (\alpha - 1 + N_1)}$$

- $\alpha - 1$  is a “pseudocount” of the number of **1**’s (or heads) you’ve “observed”
- $\beta - 1$  is a “pseudocount” of the number of **0**’s (or tails) you’ve “observed”

# Coin Flipping MAP: Example

- Suppose  $\mathcal{D}$  consists of ten 1's or heads ( $N_1 = 10$ ) and two 0's or tails ( $N_0 = 2$ ):

$$\phi_{MLE} = \frac{10}{10 + 2} = \frac{10}{12}$$

- Using a Beta prior with  $\alpha = 2$  and  $\beta = 5$ , then

$$\phi_{MAP} = \frac{(2 - 1 + 10)}{(2 - 1 + 10) + (5 - 1 + 2)} = \frac{11}{17} < \frac{10}{12}$$

# Coin Flipping MAP: Example

- Suppose  $\mathcal{D}$  consists of ten 1's or heads ( $N_1 = 10$ ) and two 0's or tails ( $N_0 = 2$ ):

$$\phi_{MLE} = \frac{10}{10 + 2} = \frac{10}{12}$$

- Using a Beta prior with  $\alpha = 101$  and  $\beta = 101$ , then

$$\phi_{MAP} = \frac{(101 - 1 + 10)}{(101 - 1 + 10) + (101 - 1 + 2)} = \frac{110}{212} \approx \frac{1}{2}$$

# Coin Flipping MAP: Example

- Suppose  $\mathcal{D}$  consists of ten **1**'s or heads ( $N_1 = 10$ ) and two **0**'s or tails ( $N_0 = 2$ ):

$$\phi_{MLE} = \frac{10}{10 + 2} = \frac{10}{12}$$

- Using a Beta prior with  $\alpha = 1$  and  $\beta = 1$ , then

$$\phi_{MAP} = \frac{(1 - 1 + 10)}{(1 - 1 + 10) + (1 - 1 + 2)} = \frac{10}{12} = \phi_{MLE}$$

# MLE/MAP Learning Objectives

You should be able to...

- Recall probability basics, including but not limited to: discrete and continuous random variables, probability mass functions, probability density functions, events vs. random variables, expectation and variance, joint probability distributions, marginal probabilities, conditional probabilities, independence, conditional independence
- State the principle of maximum likelihood estimation and explain what it tries to accomplish
- State the principle of maximum a posteriori estimation and explain why we use it
- Derive the MLE or MAP parameters of a simple model in closed form