# 10-301/601: Introduction to Machine Learning Lecture 15 – Learning Theory (Infinite Case)

Henry Chai & Matt Gormley

10/23/23

# Front Matter

- Announcements
  - HW5 released 10/9, due 10/27 (Friday) at 11:59 PM
  - Exam 3 scheduled
    - Tuesday, December 12$^{th}$ from 5:30 PM to 8:30 PM
  - Sign up for peer tutoring! See Piazza for more details

  - Exam 1 exit poll also on Piazza

# Recall - Theorem 1: Finite, Realizable Case

- For a *finite* hypothesis set $\mathcal{H}$ such that $c^* \in \mathcal{H}$ (*realizable*) 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 $R(h) \leq \epsilon$

## Recall - Theorem 1: Finite, Realizable Case

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

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

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

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

# Statistical Learning Theory Corollary

- For a *finite* hypothesis set $\mathcal{H}$ such that $c^* \in \mathcal{H}$ (*realizable*) and arbitrary distribution $p^*$, given a training dataset $S$ where $|S| = M$, 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 $1 - \delta$.

# Recall - 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

# Statistical Learning Theory Corollary

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

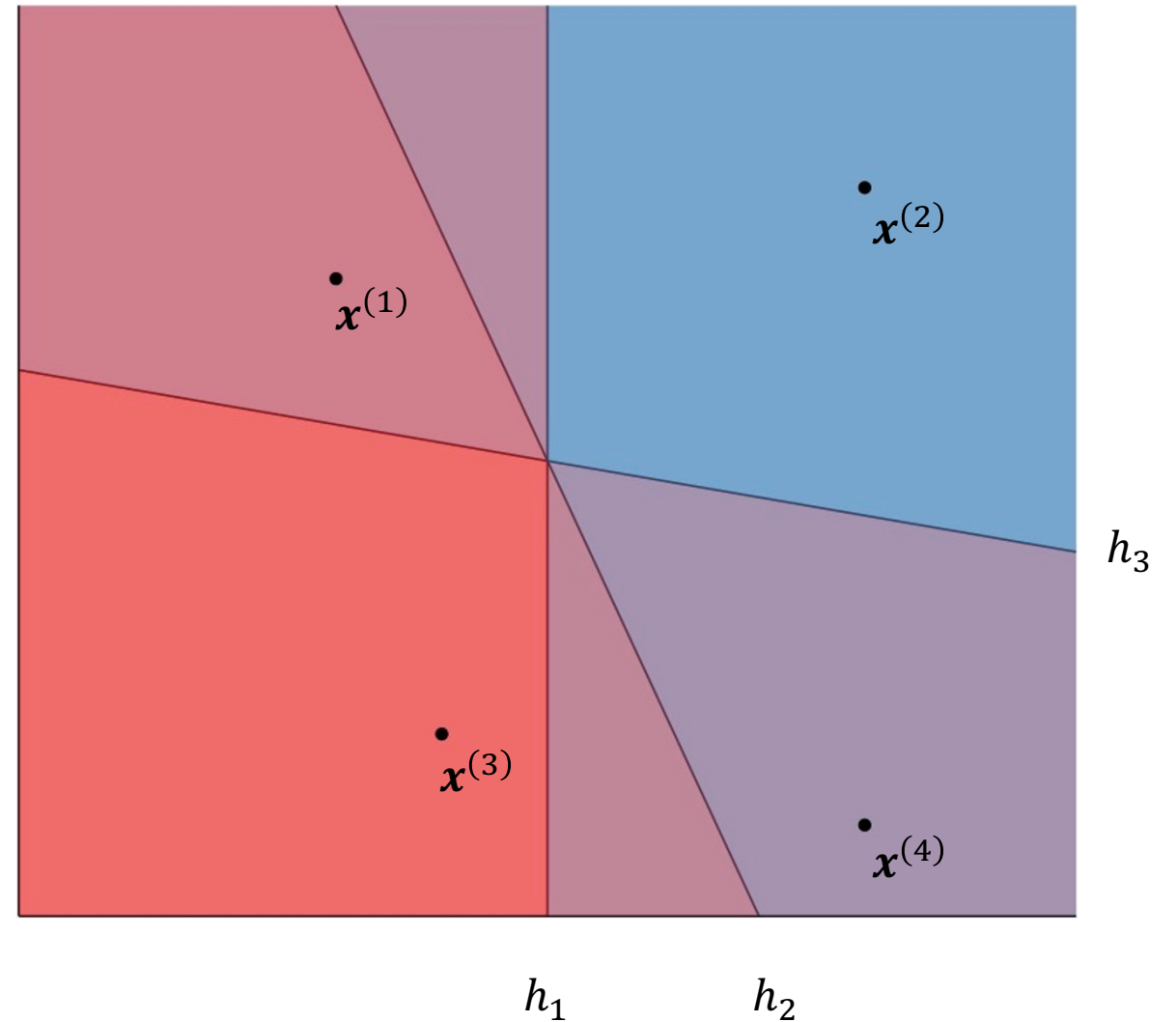$$R(h) \leq \hat{R}(h) + \sqrt{\frac{1}{2M}\left(\ln(|\mathcal{H}|) + \ln\left(\frac{2}{\delta}\right)\right)}$$

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 $S$ where $|S| = M$, 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 $1 - \delta$.

# Labellings

- Given some finite set of data points $S = \{x^{(1)}, \ldots, x^{(M)}\}$ and some hypothesis $h \in \mathcal{H}$, applying $h$ to each point in $S$ results in a **labelling**
  - $[h(x^{(1)}), \ldots, h(x^{(M)})]$ is a vector of $M$ +1's and -1's (recall: our discussion of PAC learning assumes binary classification)

- Given $S = \{x^{(1)}, \ldots, x^{(M)}\}$, 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
    $$\mathcal{H}(S) = \{[h(x^{(1)}), \ldots, h(x^{(M)})] \mid h \in \mathcal{H}\}$$

# Example: Labellings

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

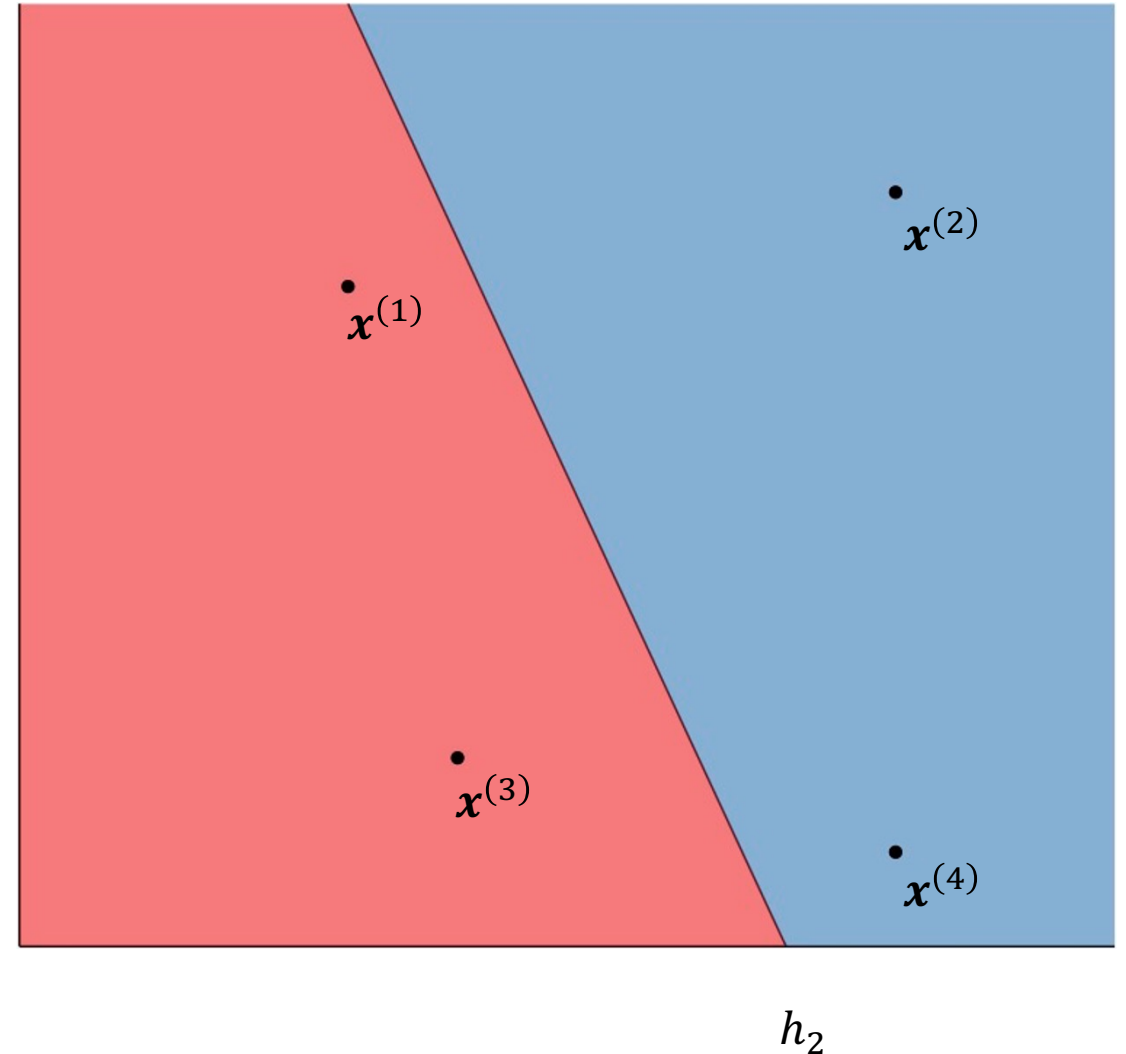# Example: Labellings

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

$\left[h_1\left(\boldsymbol{x}^{(1)}\right), h_1\left(\boldsymbol{x}^{(2)}\right), h_1\left(\boldsymbol{x}^{(3)}\right), h_1\left(\boldsymbol{x}^{(4)}\right)\right]$

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

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

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

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

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

$h_1$

# Example: Labellings

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

$\left[h_1\big(\boldsymbol{x}^{(1)}\big), h_1\big(\boldsymbol{x}^{(2)}\big), h_1\big(\boldsymbol{x}^{(3)}\big), h_1\big(\boldsymbol{x}^{(4)}\big)\right]$

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

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

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

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

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

$h_2$

# Example: Labellings

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

$\left[h_1\big(\boldsymbol{x}^{(1)}\big), h_1\big(\boldsymbol{x}^{(2)}\big), h_1\big(\boldsymbol{x}^{(3)}\big), h_1\big(\boldsymbol{x}^{(4)}\big)\right]$

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

# Example: Labellings

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

$\mathcal{H}(S)$

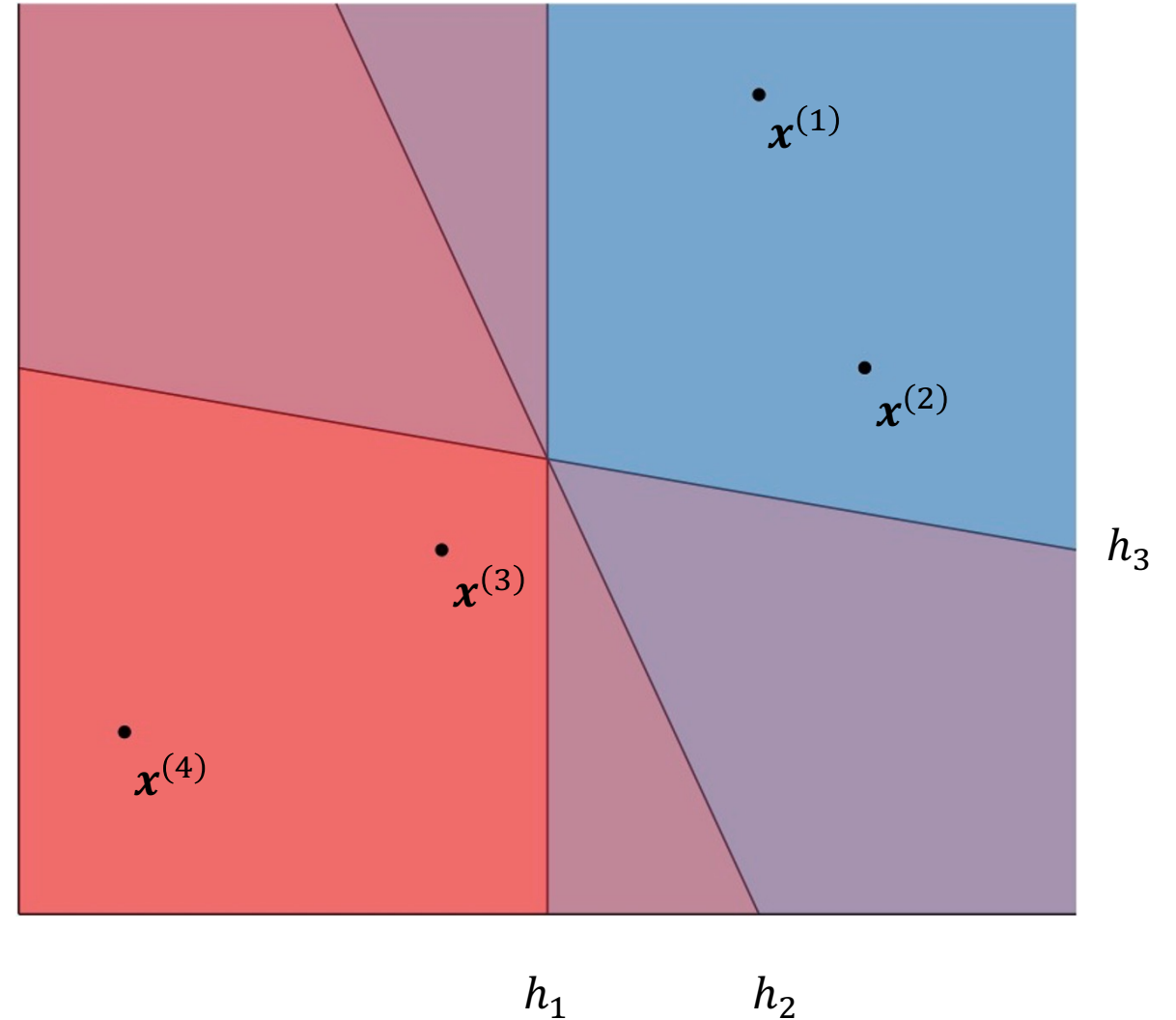$= \{[+1, +1, -1, -1], [-1, +1, -1, +1]\}$

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

# 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 $VC(\mathcal{H}) = \infty$

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

## VC-Dimension: Example

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

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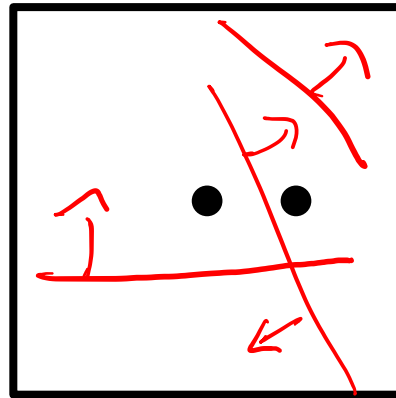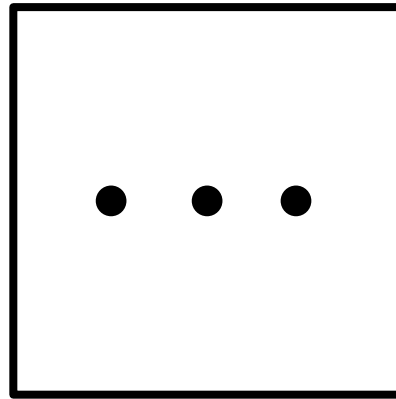  - Can $\mathcal{H}$ shatter some set of 1 point?

$S$

- $x \in \mathbb{R}^2$ 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?
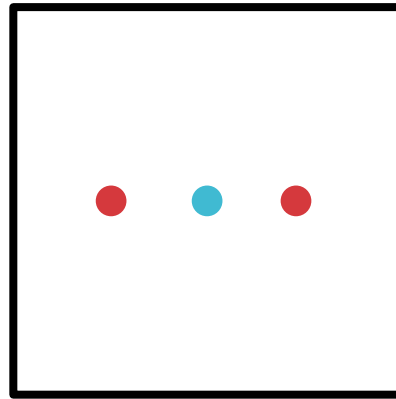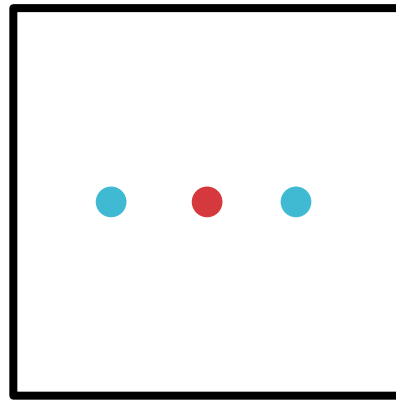
# VC-Dimension: Example



$S$

# VC-Dimension: Example

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

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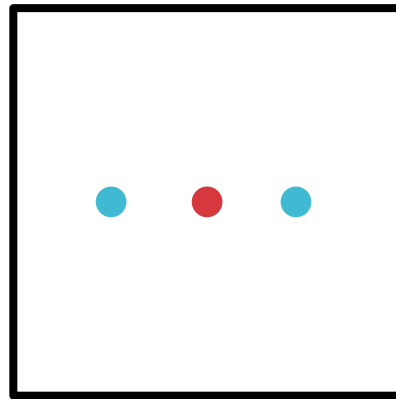  - Can $\mathcal{H}$ shatter some set of 3 points?

$S$

# VC-Dimension: Example

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

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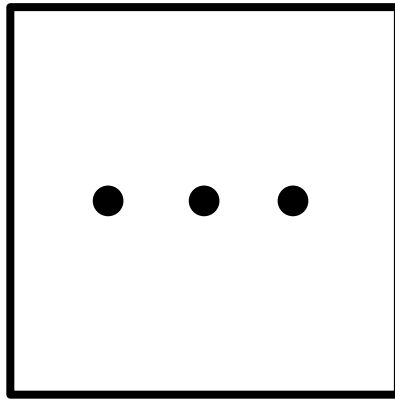  - Can $\mathcal{H}$ shatter some set of 3 points?



$S$

# VC-Dimension: Example

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

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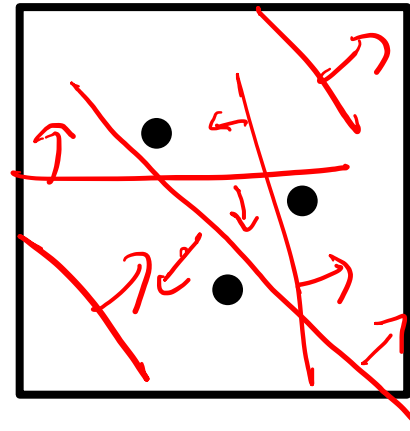  - Can $\mathcal{H}$ shatter some set of 3 points?

$S$

## VC-Dimension: Example

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

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

$S$

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

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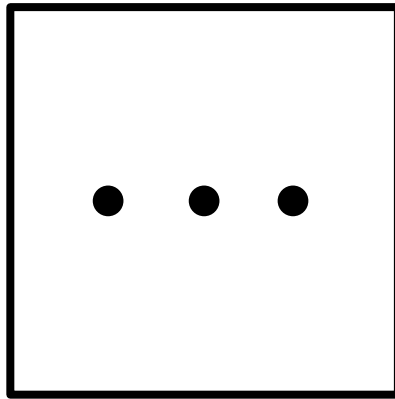  - Can $\mathcal{H}$ shatter some set of 3 points?

# VC-Dimension: Example

$S_1$

$S_2$

# VC-Dimension: Example

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

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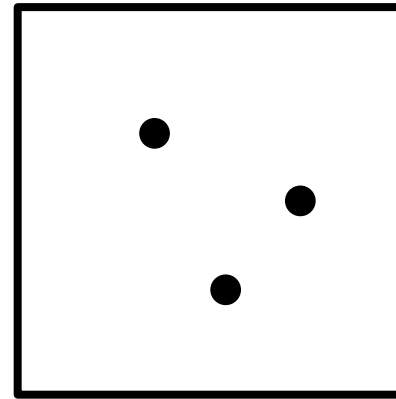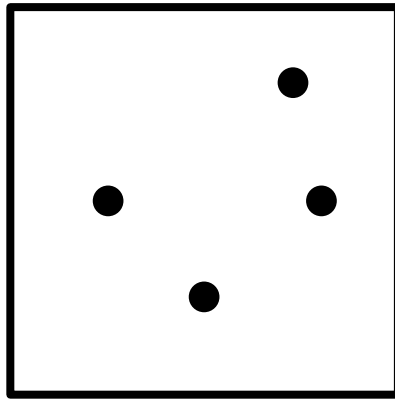  - Can $\mathcal{H}$ shatter some set of 3 points?

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

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

# VC-Dimension: Example

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

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

$S_1$

All points on the convex hull

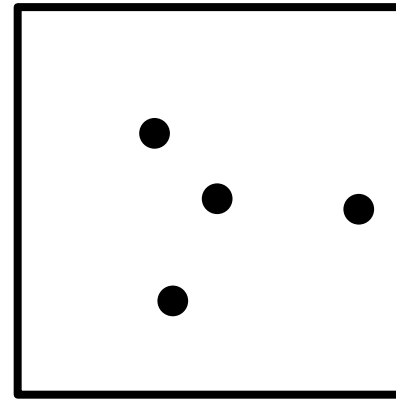$S_2$

At least one point inside the convex hull

## VC-Dimension: Example

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

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

$S_1$

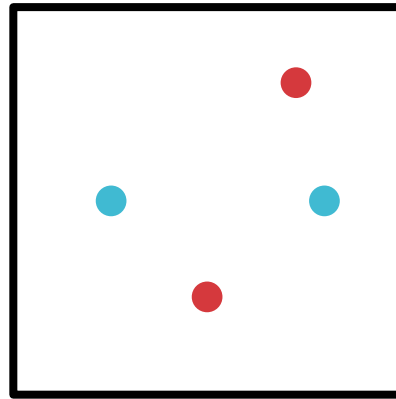All points on the convex hull

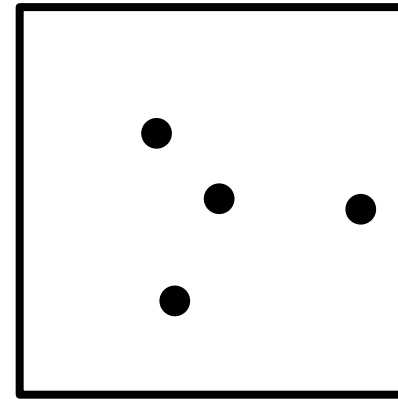$S_2$

At least one point inside the convex hull

## VC-Dimension: Example

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

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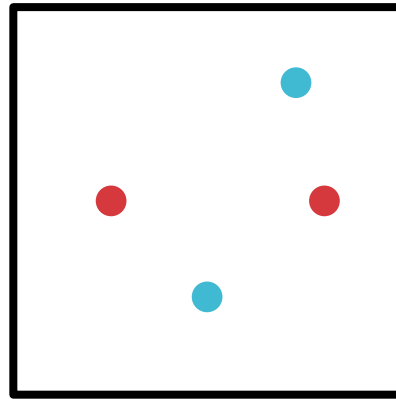
$S_1$

All points on the convex hull

$S_2$

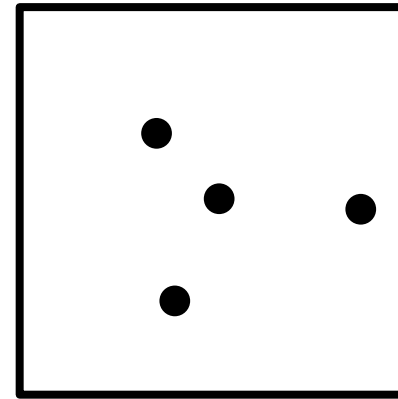At least one point inside the convex hull

## VC-Dimension: Example

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

- What is $VC(\mathcal{H})$?
  - Can $\mathcal{H}$ shatter some set of 1 point?
  - Can $\mathcal{H}$ shatter some set of 2 points?
  - Can $\mathcal{H}$ shatter some set of 3 points?
  - Can $\mathcal{H}$ shatter some set of 4 points?
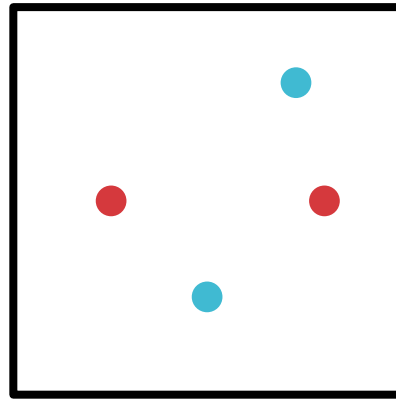
$|\mathcal{H}(S_1)| = 14$

All points on the convex hull

$S_2$

At least one point inside the convex hull

# VC-Dimension: Example

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

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



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

All points on the convex hull
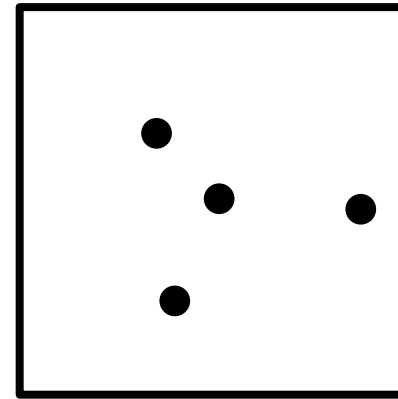
$S_2$

At least one point inside the convex hull

## VC-Dimension: Example

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

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

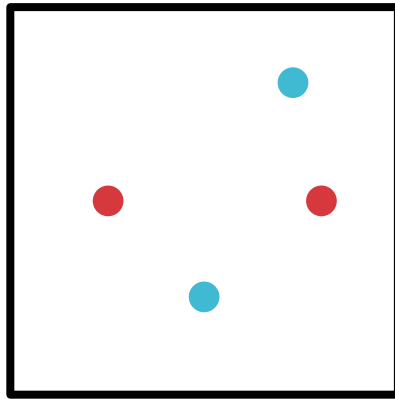$|\mathcal{H}(S_1)| = 14$

All points on the convex hull

$S_2$

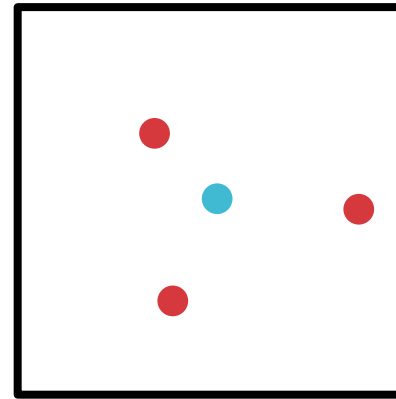At least one point inside the convex hull

# VC-Dimension: Example

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

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

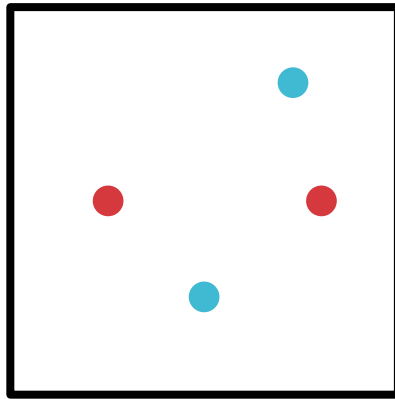$|\mathcal{H}(S_1)| = 14$

All points on the convex hull

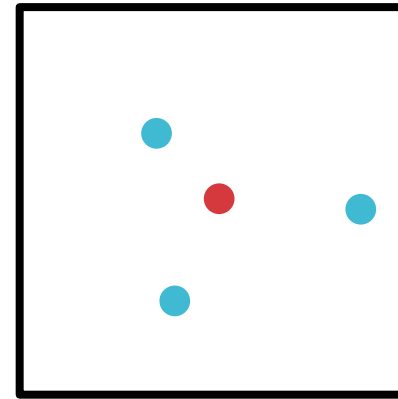$|\mathcal{H}(S_2)| = 14$

At least one point inside the convex hull

# VC-Dimension: Example

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

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

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



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

All points on the convex hull

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

At least one point inside the convex hull

# VC-Dimension: Example

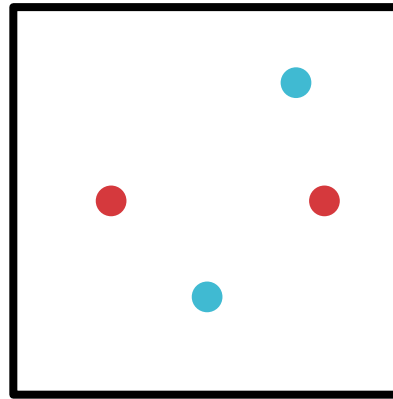- $x \in \mathbb{R}^d$ and $\mathcal{H} =$ all $d$-dimensional linear separators

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

# VC-Dimension: Example

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



$a$

# VC-Dimension: Example

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

$x^{(1)}$

$a$

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



VC-Dimension: Example

# VC-Dimension: Example

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

VC-Dimension: Example

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

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

$a$

# VC-Dimension: Example

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

# VC-Dimension: Example

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



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

$a$

- $VC(\mathcal{H}) = 1$
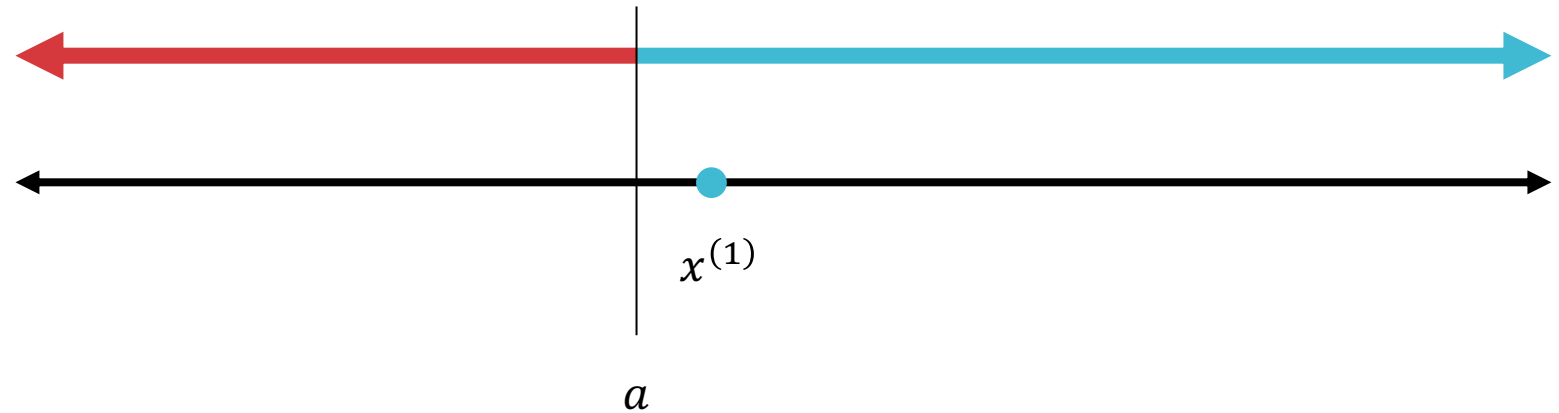
**VC-Dimension: Example**

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

$a$

$b$

## Poll Question 1:

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

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

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



$a$       $b$

VC-Dimension: Example

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



$$x^{(1)} \quad x^{(2)} \qquad x$$
$$a \qquad\qquad b$$

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

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

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

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

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 $\mathcal{H}$ such that $c^* \in \mathcal{H}$ and arbitrary distribution $p^*$, given a training dataset $S$ where $|S| = M$, 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 arbitrary distribution $p^*$, if the number of labelled training data points satisfies

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

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

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

## Statistical Learning Theory Corollary 4

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

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

with probability at least $1 - \delta$.

# Approximation Generalization Tradeoff

How well does $h$ generalize?

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

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

# Approximation Generalization Tradeoff

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

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

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

# Can we use this corollary to guide model selection?

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

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

with probability at least $1 - \delta$.

Learning Theory and Model Selection

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

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

$R(h)$

$\hat{R}(h)$ (training error)

$VC(\mathcal{H})$

error

# Learning Theory and Model Selection



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

$R(h)$ (true error)

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

$\hat{R}(h)$ (training error)

error

Best tradeoff

$VC(\mathcal{H})$

- How can we find this "best tradeoff" for linear separators?

- Use a regularizer! By (effectively) reducing the number of features our model considers, we reduce its VC-dimension.

# 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
- Theoretically motivate regularization

# Recall: Probabilistic Learning

- Previously:

  - (Unknown) Target function, $c^*: \mathcal{X} \to \mathcal{Y}$

  - Classifier, $h : \mathcal{X} \to \mathcal{Y}$

  - Goal: find a classifier, $h$, that best approximates $c^*$

- Now:

  - (Unknown) Target *distribution*, $y \sim p^*(Y|\boldsymbol{x})$

  - Distribution, $p(Y|\boldsymbol{x})$

  - Goal: find a distribution, $p$, that best approximates $p^*$

## Recall: Maximum Likelihood Estimation (MLE)

- Given independent, identically distributed observations (iid) $\mathcal{D} = \left\{ x^{(i)} \right\}_{i=1}^{N}$ from a parametrized probability distribution, MLE sets the parameters by maximizing the likelihood of the data:

$$\theta^{MLE} = \underset{\theta}{\operatorname{argmax}}\, p(\mathcal{D} \mid \theta) = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^{N} p(x^{(i)} \mid \theta)$$

- Intuition: assign as much of the (finite) probability mass to the observed data *at the expense of unobserved data*

# Recall: Maximum Likelihood Estimation (MLE)

- Given independent, identically distributed observations (iid)
  
  $\mathcal{D} = \left\{ x^{(i)} \right\}_{i=1}^{N}$ from a parametrized probability distribution, MLE sets the parameters by maximizing the *log*-likelihood of the data:

$$\theta^{MLE} = \operatorname*{argmax}_{\theta} \log p(\mathcal{D} \mid \theta) = \operatorname*{argmax}_{\theta} \sum_{i=1}^{N} \log p\left( x^{(i)} \mid \theta \right)$$

- Intuition: assign as much of the (finite) probability mass to the observed data *at the expense of unobserved data*

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

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

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

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

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

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

$$\text{where } N_i = \# \text{ of } i\text{'s in } D$$

$$\log(a^b) = b \log a$$

## 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

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

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

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

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

$$\Rightarrow \hat{\phi} = N_1 / N_0 + N_1$$

Poll Question 2:

After flipping your 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

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

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

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

- where $N_1$ is the number of $1$'s in $\left\{x^{(1)}, \ldots, x^{(N)}\right\}$ 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

$$\text{MLE} \quad \text{finds} \quad \hat{\theta} = \underset{\theta}{\text{argmax}} \; P(D|\theta)$$

$$\text{MAP} \quad \text{finds} \quad \hat{\theta} = \underset{\theta}{\text{argmax}} \; P(\theta|D)$$

$$= \underset{\theta}{\text{argmax}} \; \frac{P(D|\theta)P(\theta)}{P(D)}$$

$$= \underset{\theta}{\text{argmax}} \; \underset{\text{likelihood}}{P(D|\theta)} \; \underset{\text{prior}}{P(\theta)}$$

# Maximum a Posteriori (MAP) Estimation

1. Specify the *generative story*, i.e., the data generating distribution, including a *prior distribution*

2. Maximize the log-posterior of $\mathcal{D} = \{x^{(1)}, \ldots, x^{(N)}\}$

$$\ell_{MAP}(\theta) = \log p(\theta) + \sum_{i=1}^{N} \log p\left(x^{(i)} | \theta\right)$$

3. Solve in *closed form*: take partial derivatives, set to 0 and solve