



# 10-601 Introduction to Machine Learning

Machine Learning Department  
School of Computer Science  
Carnegie Mellon University

## PAC Learning + Oracles, Sampling, Generative vs. Discriminative

Matt Gormley  
Lecture 16  
Oct. 24, 2018

# Q&A

**Q:** Why do we shuffle the examples in SGD?

**A:** This is how we do sampling *without* replacement

1. **Theoretically** we can show sampling **without replacement** is not significantly worse than sampling with replacement (Shamir, 2016)
2. **Practically** sampling without replacement tends to work better

**Q:** What is “bias”?

**A:** That depends. The word “bias” shows up all over machine learning! Watch out...

1. The additive term in a linear model (i.e.  $b$  in  $w^T x + b$ )
2. Inductive bias is the principle by which a learning algorithm generalizes to unseen examples
3. Bias of a model in a societal sense may refer to racial, socio-economic, gender biases that exist in the predictions of your model
4. The difference between the expected predictions of your model and the ground truth (as in “bias-variance tradeoff”)

(See your TAs excellent post here:  
<https://piazza.com/class/jkmt7l4of093k5?cid=383>)

# Reminders

- **Midterm Exam**
  - Thursday Evening 6:30 – 9:00 (2.5 hours)
  - Room and seat assignments announced on Piazza
  - You may bring one 8.5 x 11 cheatsheet

# Sample Complexity Results

**Definition 0.1.** The **sample complexity** of a learning algorithm is the number of examples required to achieve arbitrarily small error (with respect to the optimal hypothesis) with high probability (i.e. close to 1).

**Four Cases we care about...**

	Realizable	Agnostic
Finite $ \mathcal{H} $	$N \geq \frac{1}{\epsilon} [\log( \mathcal{H} ) + \log(\frac{1}{\delta})]$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $R(h) \geq \epsilon$ have $\hat{R}(h) > 0$ .	$N \geq \frac{1}{2\epsilon^2} [\log( \mathcal{H} ) + \log(\frac{2}{\delta})]$ labeled examples are sufficient so that with probability $(1 - \delta)$ for all $h \in \mathcal{H}$ we have that $ R(h) - \hat{R}(h)  < \epsilon$ .
Infinite $ \mathcal{H} $	$N = O(\frac{1}{\epsilon} [\text{VC}(\mathcal{H}) \log(\frac{1}{\epsilon}) + \log(\frac{1}{\delta})])$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $R(h) \geq \epsilon$ have $\hat{R}(h) > 0$ .	$N = O(\frac{1}{\epsilon^2} [\text{VC}(\mathcal{H}) + \log(\frac{1}{\delta})])$ labeled examples are sufficient so that with probability $(1 - \delta)$ for all $h \in \mathcal{H}$ we have that $ R(h) - \hat{R}(h)  \leq \epsilon$ .

# Generalization and Inductive Bias

## Chalkboard:

- Setting: binary classification with binary feature vectors
- Instance space vs. Hypothesis space
- Counting: # of instances, # leaves in a full decision tree, # of full decision trees, # of labelings of training examples
- Algorithm: keep all full decision trees consistent with the training data and do a majority vote to classify
- Case study: training size is all, all-but-one, all-but-two, all-but-three,...

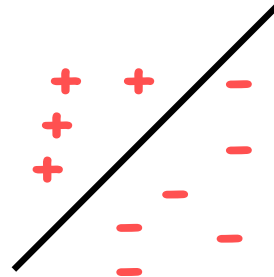
# **VC DIMENSION**



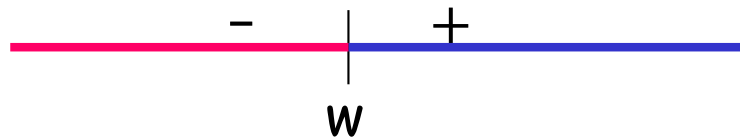
# What if $H$ is infinite?



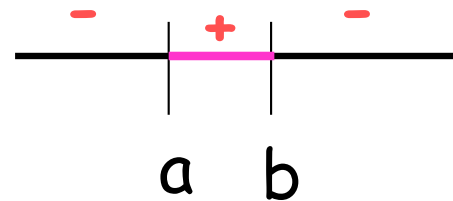
E.g., linear separators in  $\mathbb{R}^d$



E.g., thresholds on the real line



E.g., intervals on the real line



# Shattering, VC-dimension

**Definition:**

$H[S]$  - the set of splittings of dataset  $S$  using concepts from  $H$ .

$H$  shatters  $S$  if  $|H[S]| = 2^{|S|}$ .

A set of points  $S$  is shattered by  $H$  if there are hypotheses in  $H$  that split  $S$  in all of the  $2^{|S|}$  possible ways; i.e., all possible ways of classifying points in  $S$  are achievable using concepts in  $H$ .

**Definition:** VC-dimension (Vapnik-Chervonenkis dimension)

The **VC-dimension** of a hypothesis space  $H$  is the cardinality of the largest set  $S$  that can be shattered by  $H$ .

If arbitrarily large finite sets can be shattered by  $H$ , then  $\text{VCdim}(H) = \infty$



# Shattering, VC-dimension

**Definition:** VC-dimension (Vapnik-Chervonenkis dimension)

The **VC-dimension** of a hypothesis space  $H$  is the cardinality of the largest set  $S$  that can be shattered by  $H$ .

If arbitrarily large finite sets can be shattered by  $H$ , then  $VCdim(H) = \infty$

To show that VC-dimension is  $d$ :

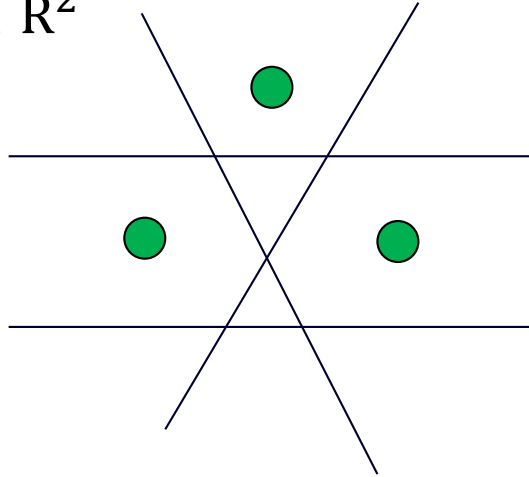
- **there exists** a set of  **$d$  points** that can be shattered
- there is **no set of  $d+1$  points** that can be shattered.

**Fact:** If  $H$  is finite, then  $VCdim(H) \leq \log(|H|)$ .

# Shattering, VC-dimension

E.g.,  $H$  = linear separators in  $\mathbb{R}^2$

$\text{VCdim}(H) \geq 3$

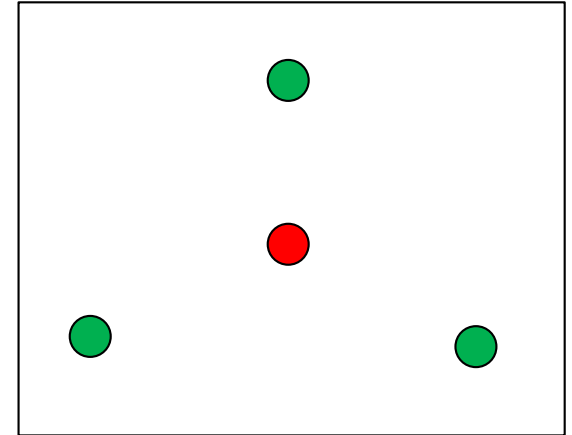


# Shattering, VC-dimension

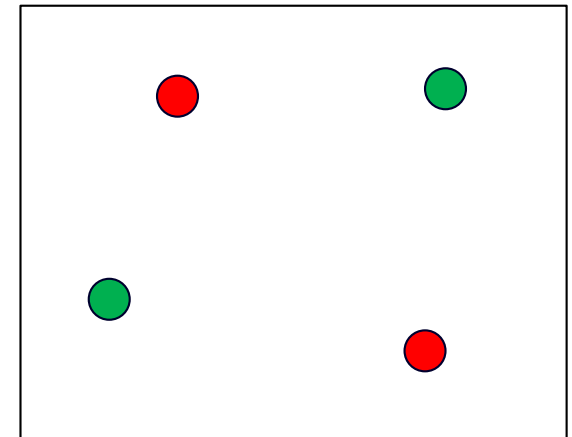
E.g.,  $H$  = linear separators in  $\mathbb{R}^2$

$\text{VCdim}(H) < 4$

Case 1: one point inside the triangle formed by the others. Cannot label inside point as positive and outside points as negative.



Case 2: all points on the boundary (convex hull). Cannot label two diagonally as positive and other two as negative.

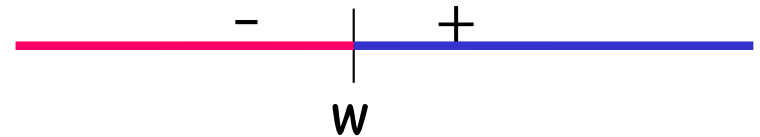


Fact:  $\text{VCdim}$  of linear separators in  $\mathbb{R}^d$  is  $d+1$

# Shattering, VC-dimension

If the VC-dimension is  $d$ , that means **there exists** a set of  $d$  points that can be shattered, but there is **no** set of  $d+1$  points that can be shattered.

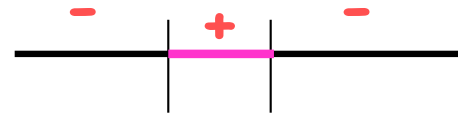
E.g.,  $H$  = Thresholds on the real line



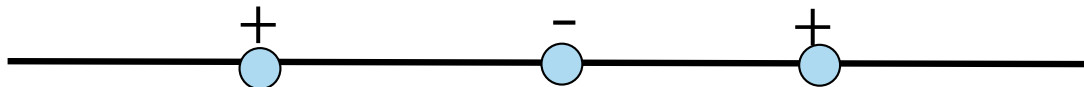
$$\text{VCdim}(H) = 1$$



E.g.,  $H$  = Intervals on the real line



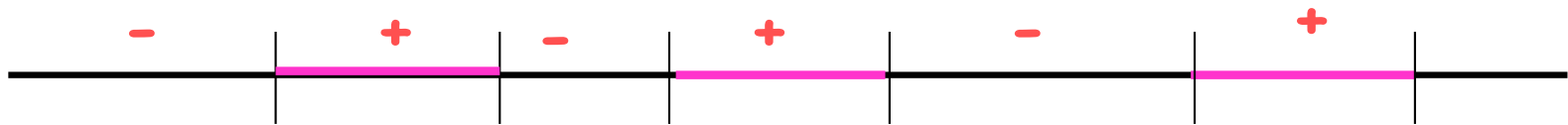
$$\text{VCdim}(H) = 2$$



# Shattering, VC-dimension

If the VC-dimension is  $d$ , that means **there exists** a set of  $d$  points that can be shattered, but there is **no** set of  $d+1$  points that can be shattered.

E.g.,  $H = \text{Union of } k \text{ intervals on the real line}$   $\text{VCdim}(H) = 2k$



$$\text{VCdim}(H) \geq 2k$$

A sample of size  $2k$  shatters  
(treat each pair of points as a  
separate case of intervals)

$$\text{VCdim}(H) < 2k + 1$$



# Sample Complexity Results

**Definition 0.1.** The **sample complexity** of a learning algorithm is the number of examples required to achieve arbitrarily small error (with respect to the optimal hypothesis) with high probability (i.e. close to 1).

**Four Cases we care about...**

	Realizable	Agnostic
Finite $ \mathcal{H} $	$N \geq \frac{1}{\epsilon} [\log( \mathcal{H} ) + \log(\frac{1}{\delta})]$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $R(h) \geq \epsilon$ have $\hat{R}(h) > 0$ .	$N \geq \frac{1}{2\epsilon^2} [\log( \mathcal{H} ) + \log(\frac{2}{\delta})]$ labeled examples are sufficient so that with probability $(1 - \delta)$ for all $h \in \mathcal{H}$ we have that $ R(h) - \hat{R}(h)  < \epsilon$ .
Infinite $ \mathcal{H} $	$N = O(\frac{1}{\epsilon} [\text{VC}(\mathcal{H}) \log(\frac{1}{\epsilon}) + \log(\frac{1}{\delta})])$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $R(h) \geq \epsilon$ have $\hat{R}(h) > 0$ .	$N = O(\frac{1}{\epsilon^2} [\text{VC}(\mathcal{H}) + \log(\frac{1}{\delta})])$ labeled examples are sufficient so that with probability $(1 - \delta)$ for all $h \in \mathcal{H}$ we have that $ R(h) - \hat{R}(h)  \leq \epsilon$ .

# SLT-style Corollaries

**Corollary 3 (Realizable, Infinite  $|\mathcal{H}|$ ).** For some  $\delta > 0$ , with probability at least  $(1 - \delta)$ , for any hypothesis  $h$  in  $\mathcal{H}$  consistent with the data (i.e. with  $\hat{R}(h) = 0$ ),

$$R(h) \leq O \left( \frac{1}{N} \left[ \text{VC}(\mathcal{H}) \ln \left( \frac{N}{\text{VC}(\mathcal{H})} \right) + \ln \left( \frac{1}{\delta} \right) \right] \right) \quad (1)$$

**Corollary 4 (Agnostic, Infinite  $|\mathcal{H}|$ ).** For some  $\delta > 0$ , with probability at least  $(1 - \delta)$ , for all hypotheses  $h$  in  $\mathcal{H}$ ,

$$R(h) \leq \hat{R}(h) + O \left( \sqrt{\frac{1}{N} \left[ \text{VC}(\mathcal{H}) + \ln \left( \frac{1}{\delta} \right) \right]} \right) \quad (2)$$

# Generalization and Overfitting

*Whiteboard:*

- Empirical Risk Minimization
- Structural Risk Minimization
- Motivation for Regularization



# Questions For Today

1. Given a classifier with zero training error, what can we say about generalization error?  
(Sample Complexity, Realizable Case)
2. Given a classifier with low training error, what can we say about generalization error?  
(Sample Complexity, Agnostic Case)
3. Is there a theoretical justification for regularization to avoid overfitting?  
(Structural Risk Minimization)

# Learning Theory 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 learning examples
- Distinguish between a large sample and a finite sample analysis
- Theoretically motivate regularization

The Big Picture

# **CLASSIFICATION AND REGRESSION**

# Classification and Regression: The Big Picture

## *Whiteboard*

- **Decision Rules / Models** (probabilistic generative, probabilistic discriminative, perceptron, SVM, regression)
- **Objective Functions** (likelihood, conditional likelihood, hinge loss, mean squared error)
- **Regularization** (L1, L2, priors for MAP)
- **Update Rules** (SGD, perceptron)
- **Nonlinear Features** (preprocessing, kernel trick)

# ML Big Picture

## Learning Paradigms:

*What data is available and when? What form of prediction?*

- supervised learning
- unsupervised learning
- semi-supervised learning
- reinforcement learning
- active learning
- imitation learning
- domain adaptation
- online learning
- density estimation
- recommender systems
- feature learning
- manifold learning
- dimensionality reduction
- ensemble learning
- distant supervision
- hyperparameter optimization

## Theoretical Foundations:

*What principles guide learning?*

- ☐ probabilistic
- ☐ information theoretic
- ☐ evolutionary search
- ☐ ML as optimization

## Problem Formulation:

*What is the structure of our output prediction?*

boolean	Binary Classification
categorical	Multiclass Classification
ordinal	Ordinal Classification
real	Regression
ordering	Ranking
multiple discrete	Structured Prediction
multiple continuous	(e.g. dynamical systems)
both discrete & cont.	(e.g. mixed graphical models)

## Facets of Building ML Systems:

*How to build systems that are robust, efficient, adaptive, effective?*

1. Data prep
2. Model selection
3. Training (optimization / search)
4. Hyperparameter tuning on validation data
5. (Blind) Assessment on test data

## Big Ideas in ML:

*Which are the ideas driving development of the field?*

- inductive bias
- generalization / overfitting
- bias-variance decomposition
- generative vs. discriminative
- deep nets, graphical models
- PAC learning
- distant rewards

## Application Areas

*Key challenges?*

NLP, Speech, Computer Vision, Robotics, Medicine, Search

# **PROBABILISTIC LEARNING**

# Probabilistic Learning

## Function Approximation

Previously, we assumed that our output was generated using a **deterministic target function**:

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

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

Our goal was to learn a hypothesis  $h(\mathbf{x})$  that best approximates  $c^*(\mathbf{x})$

## Probabilistic Learning

Today, we assume that our output is **sampled** from a conditional **probability distribution**:

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

$$y^{(i)} \sim p^*(\cdot | \mathbf{x}^{(i)})$$

Our goal is to learn a probability distribution  $p(y|\mathbf{x})$  that best approximates  $p^*(y|\mathbf{x})$



# Robotic Farming

	Deterministic	Probabilistic
Classification (binary output)	Is this a picture of a wheat kernel?	Is this plant drought resistant?
Regression (continuous output)	How many wheat kernels are in this picture?	What will the yield of this plant be?





# Oracles and Sampling

## *Whiteboard*

- Sampling from common probability distributions
  - Bernoulli
  - Categorical
  - Uniform
  - Gaussian
- Pretending to be an Oracle (Regression)
  - Case 1: Deterministic outputs
  - Case 2: Probabilistic outputs
- Probabilistic Interpretation of Linear Regression
  - Adding Gaussian noise to linear function
  - Sampling from the noise model
- Pretending to be an Oracle (Classification)
  - Case 1: Deterministic labels
  - Case 2: Probabilistic outputs (Logistic Regression)
  - Case 3: Probabilistic outputs (Gaussian Naïve Bayes)

# In-Class Exercise

1. With your neighbor, **write a function** which returns **samples from a Categorical**
  - Assume access to the `rand()` function
  - Function signature should be:  
`categorical_sample(theta)`  
where `theta` is the array of parameters
  - Make your implementation as **efficient** as possible!
2. What is the **expected runtime** of your function?

# Generative vs. Discriminative

## *Whiteboard*

- Generative vs. Discriminative Models
  - Chain rule of probability
  - Maximum (Conditional) Likelihood Estimation for Discriminative models
  - Maximum Likelihood Estimation for Generative models

# Categorical Distribution

## *Whiteboard*

- Categorical distribution details
  - Independent and Identically Distributed (i.i.d.)
  - Example: Dice Rolls

# Takeaways

- One view of what ML is trying to accomplish is **function approximation**
- The principle of **maximum likelihood estimation** provides an alternate view of learning
- **Synthetic data** can help **debug** ML algorithms
- Probability distributions can be used to **model** real data that occurs in the world  
(don't worry we'll make our distributions more interesting soon!)

# Learning Objectives

## **Oracles, Sampling, Generative vs. Discriminative**

*You should be able to...*

1. Sample from common probability distributions
2. Write a generative story for a generative or discriminative classification or regression model
3. Pretend to be a data generating oracle
4. Provide a probabilistic interpretation of linear regression
5. Use the chain rule of probability to contrast generative vs. discriminative modeling
6. Define maximum likelihood estimation (MLE) and maximum conditional likelihood estimation (MCLE)