# PAC Learning

Matt Gormley
Lecture 14
March 5, 2018

# 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) |

## Application Areas
*Key challenges?*
NLP, Speech, Computer Vision, Robotics, Medicine, Search

## 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

# LEARNING THEORY

# Questions For Today
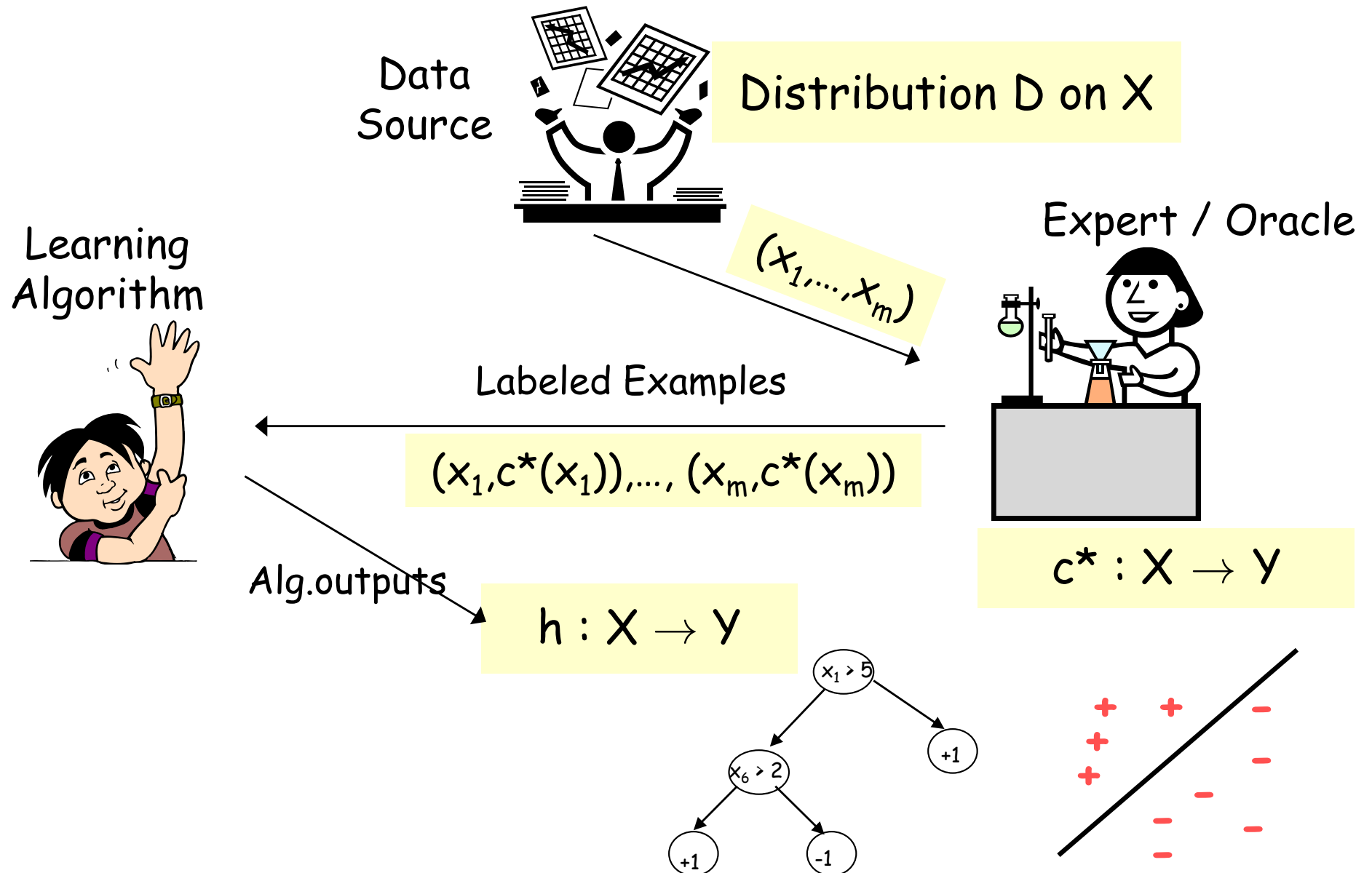
1.  Given a classifier with zero training error, what can we say about generalization error?
    (Sample Complexity, Realizable Case)

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)

# PAC/SLT models for Supervised Learning



Data Source

Distribution D on X

$(x_1,...,x_m)$

Expert / Oracle

Learning Algorithm

Labeled Examples

$(x_1,c^*(x_1)),..., (x_m,c^*(x_m))$

$c^* : X \rightarrow Y$

Alg.outputs

$h : X \rightarrow Y$

$x_1 > 5$

$+1$

$x_6 > 2$

$+1$

$-1$

# Two Types of Error

True Error (aka. **expected risk**)

$$R(h) = P_{\mathbf{x} \sim p^*(\mathbf{x})}(c^*(\mathbf{x}) \neq h(\mathbf{x}))$$

This quantity is always **unknown**

Train Error (aka. **empirical risk**)

$$\hat{R}(h) = P_{\mathbf{x} \sim \mathcal{S}}(c^*(\mathbf{x}) \neq h(\mathbf{x}))$$

$$= \frac{1}{N} \sum_{i=1}^{N} \mathbb{1}(c^*(\mathbf{x}^{(i)}) \neq h(\mathbf{x}^{(i)}))$$

$$= \frac{1}{N} \sum_{i=1}^{N} \mathbb{1}(y^{(i)} \neq h(\mathbf{x}^{(i)}))$$

We can **measure** this on the training data

where $\mathcal{S} = \{\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(N)})\}_{i=1}^{N}$ is the training data set, and $\mathbf{x} \sim \mathcal{S}$ denotes that $\mathbf{x}$ is sampled from the empirical distribution.

# PAC / SLT Model

1. Generate instances from unknown distribution $p^*$

$$x^{(i)} \sim p^*(x), \forall i \qquad (1)$$

2. Oracle labels each instance with unknown function $c^*$

$$y^{(i)} = c^*(x^{(i)}), \forall i \qquad (2)$$

3. Learning algorithm chooses hypothesis $h \in \mathcal{H}$ with low(est) training error, $\hat{R}(h)$

$$\hat{h} = \underset{h}{\operatorname{argmin}} \hat{R}(h) \qquad (3)$$

4. Goal: Choose an $h$ with low generalization error $R(h)$

# Three Hypotheses of Interest

The **true function** $c^*$ is the one we are trying to learn and that labeled the training data:

$$y^{(i)} = c^*(x^{(i)}), \ \forall i \tag{1}$$

The **expected risk minimizer** has lowest true error:

$$h^* = \underset{h \in \mathcal{H}}{\arg\min} \ R(h) \tag{2}$$

The **empirical risk minimizer** has lowest training error:

$$\hat{h} = \underset{h \in \mathcal{H}}{\arg\min} \ \hat{R}(h) \tag{3}$$

# PAC LEARNING

# Probably Approximately Correct (PAC) Learning

*Whiteboard:*

- PAC Criterion

- Meaning of "Probably Approximately Correct"

- PAC Learnable

- Consistent Learner

- Sample Complexity

# Generalization and Overfitting

*Whiteboard:*

- Realizable vs. Agnostic Cases
- Finite vs. Infinite Hypothesis Spaces

# PAC Learning

The **PAC criterion** is that our learner produces a high accuracy learner with high probability:

$$P(|R(h) - \hat{R}(h)| \leq \epsilon) \geq 1 - \delta \tag{1}$$

Suppose we have a learner that produces a hypothesis $h \in \mathcal{H}$ given a sample of $N$ training examples. The algorithm is called **consistent** if for every $\epsilon$ and $\delta$, there exists a positive number of training examples $N$ such that for any distribution $p^*$, we have that:

$$P(|R(h) - \hat{R}(h)| > \epsilon) < \delta \tag{2}$$

The **sample complexity** is the minimum value of $N$ for which this statement holds. If $N$ is finite for some learning algorithm, then $\mathcal{H}$ is said to be **learnable**. If $N$ is a polynomial function of $\frac{1}{\epsilon}$ and $\frac{1}{\delta}$ for some learning algorithm, then $\mathcal{H}$ is said to be **PAC learnable**.

# SAMPLE COMPLEXITY RESULTS

# 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).
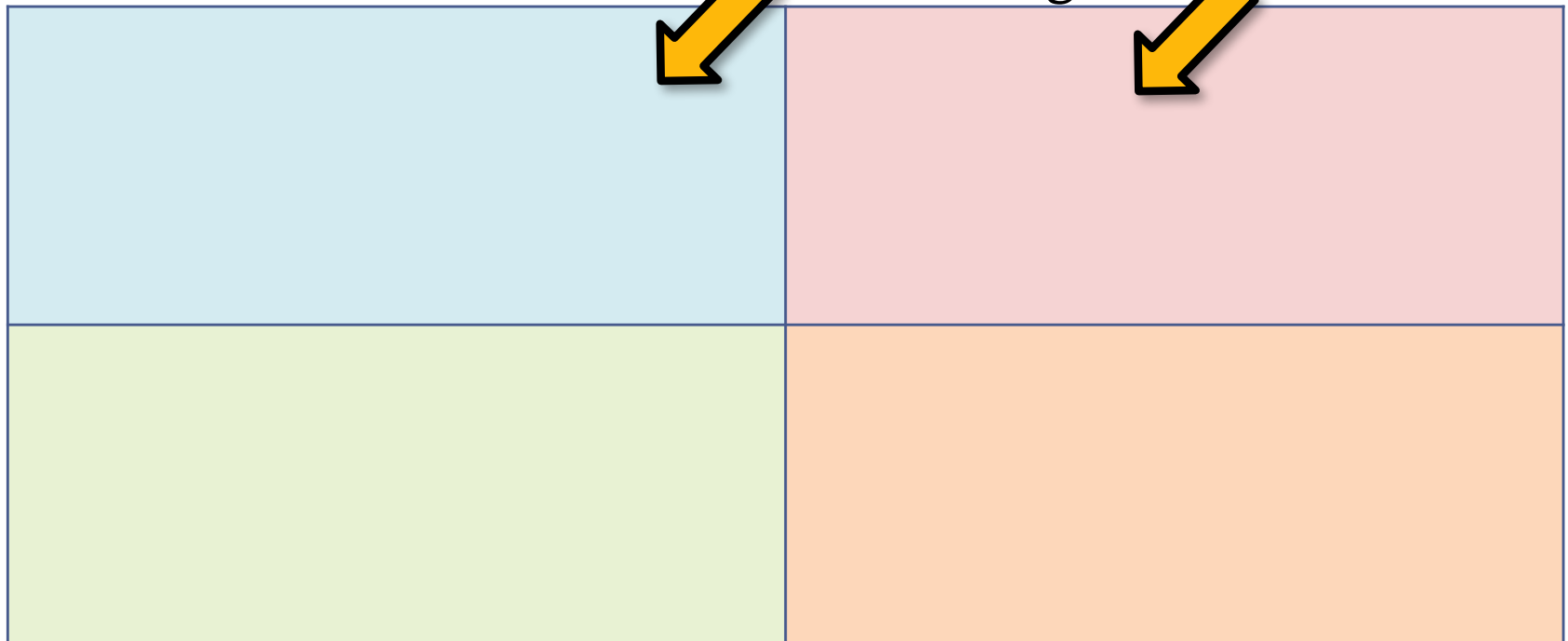
**Four Cases we care about…**

We'll start with the finite case…

|  | Realizable | Agnostic |
|---|---|---|
| Finite $|\mathcal{H}|$ | | |
| Infinite $|\mathcal{H}|$ | | |

# 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 $|\mathcal{H}|$ | $N \geq \frac{1}{\epsilon}\left[\log(|\mathcal{H}|) + \log(\frac{1}{\delta})\right]$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $R(h) \geq \epsilon$ have $\hat{R}(h) > 0$. | |
| Infinite $|\mathcal{H}|$ | | |

# Example: Conjunctions

*In-Class Quiz:*

Suppose H = class of conjunctions over **x** in $\{0,1\}^M$

If $M = 10$, $\varepsilon = 0.1$, $\delta = 0.01$, how many examples suffice?

|  | Realizable | Agnostic |
|---|---|---|
| Finite $|\mathcal{H}|$ | $N \geq \frac{1}{\epsilon}\left[\log(|\mathcal{H}|) + \log(\frac{1}{\delta})\right]$ labeled examples are sufficient so that with probability $(1-\delta)$ all $h \in \mathcal{H}$ with $R(h) \geq \epsilon$ have $\hat{R}(h) > 0$. |  |
| Infinite $|\mathcal{H}|$ |  |  |

# 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** $|\mathcal{H}|$ | $N \geq \frac{1}{\epsilon}\left[\log(|\mathcal{H}|) + \log(\frac{1}{\delta})\right]$ labeled examples are sufficient so that with probability $(1-\delta)$ all $h \in \mathcal{H}$ with $R(h) \geq \epsilon$ have $\hat{R}(h) > 0$. | $N \geq \frac{1}{2\epsilon^2}\left[\log(|\mathcal{H}|) + \log(\frac{2}{\delta})\right]$ 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** $|\mathcal{H}|$ | | |

1. Bound is **inversely linear in epsilon** (e.g. halving the error requires double the examples)
2. Bound is **only logarithmic in |H|** (e.g. quadrupling the hypothesis space only requires double the examples)

1. Bound is **inversely quadratic in epsilon** (e.g. halving the error requires 4x the examples)
2. Bound is **only logarithmic in |H|** (i.e. same as Realizable case)

Realizable

Agnostic

| | Realizable | Agnostic |
|---|---|---|
| Finite $|\mathcal{H}|$ | $N \geq \frac{1}{\epsilon}\left[\log(|\mathcal{H}|) + \log(\frac{1}{\delta})\right]$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $R(h) \geq \epsilon$ have $\hat{R}(h) > 0$. | $N \geq \frac{1}{2\epsilon^2}\left[\log(|\mathcal{H}|) + \log(\frac{2}{\delta})\right]$ 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 $|\mathcal{H}|$ | | |

# Generalization and Overfitting

*Whiteboard:*

- Sample Complexity Bounds (Agnostic Case)
- Corollary (Agnostic Case)
- Empirical Risk Minimization
- Structural Risk Minimization
- Motivation for Regularization

# 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 $|\mathcal{H}|$ | $N \geq \frac{1}{\epsilon}\left[\log(|\mathcal{H}|) + \log(\frac{1}{\delta})\right]$ labeled examples are sufficient so that with probability $(1-\delta)$ all $h \in \mathcal{H}$ with have $\hat{R}(h) > 0$. | $N \geq \frac{1}{2\epsilon^2}\left[\log(|\mathcal{H}|) + \log(\frac{2}{\delta})\right]$ labeled examples are sufficient so $-\delta)$ for $|R(h) -$ |
| Infinite $|\mathcal{H}|$ | | |

> We need a new definition of "complexity" for a Hypothesis space for these results (see *VC Dimension*)

# 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