



10-301/10-601 Introduction to Machine Learning

Machine Learning Department
School of Computer Science
Carnegie Mellon University

Course Overview

Matt Gormley & Geoff Gordon

Lecture 1

Aug. 25, 2025

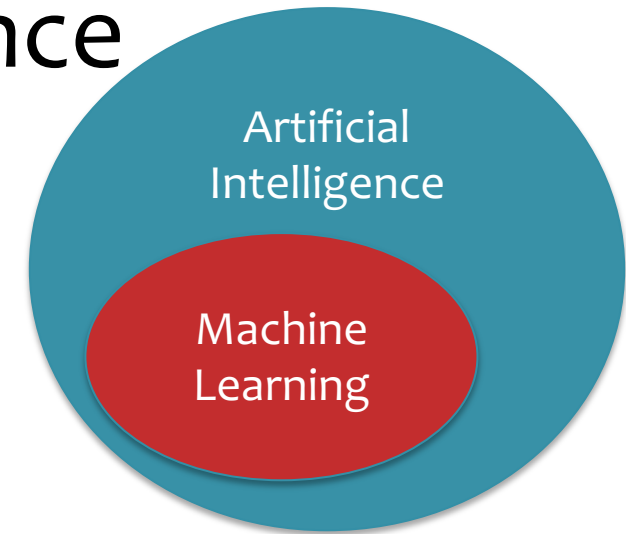
WHAT IS MACHINE LEARNING?

Artificial Intelligence

The basic goal of AI is to develop intelligent machines.

This consists of many sub-goals:

- Perception
- Reasoning
- Control / Motion / Manipulation
- Planning
- Communication
- Creativity
- Learning

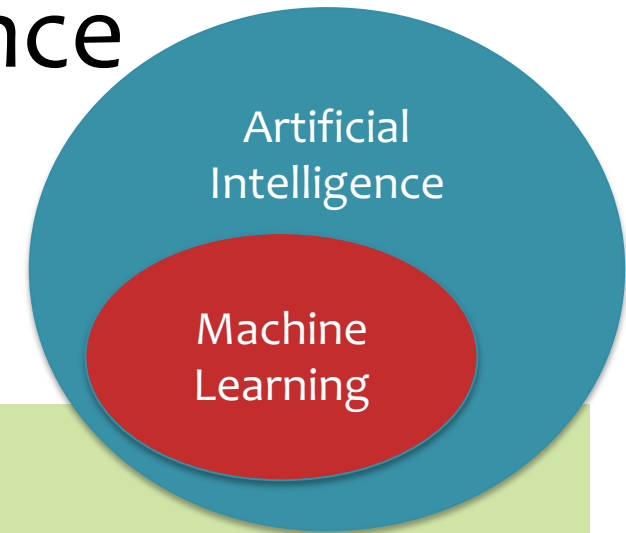


Artificial Intelligence

The basic goal of AI is to develop intelligent machines.

This consists of many sub-goals:

- Perception
- Reasoning
- Control / Motion / Manipulation
- Planning
- Communication
- Creativity
- Learning

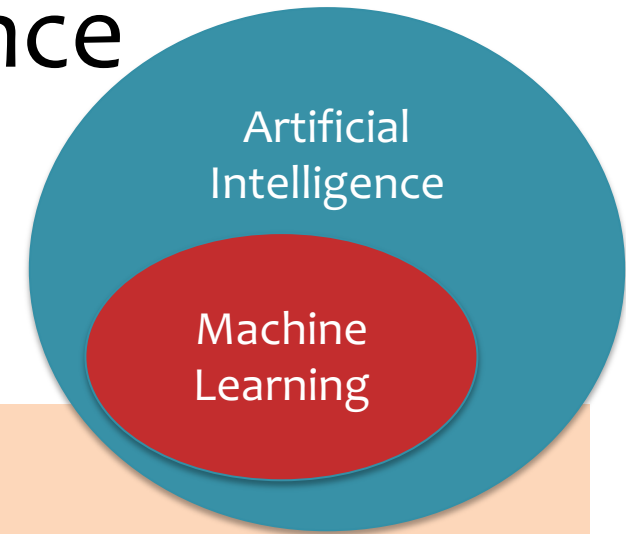


Artificial Intelligence

The basic goal of AI is to develop intelligent machines.

This consists of many sub-goals:

- Perception
- Reasoning
- Control / Motion / Manipulation
- Planning
- Communication
- Creativity
- Learning

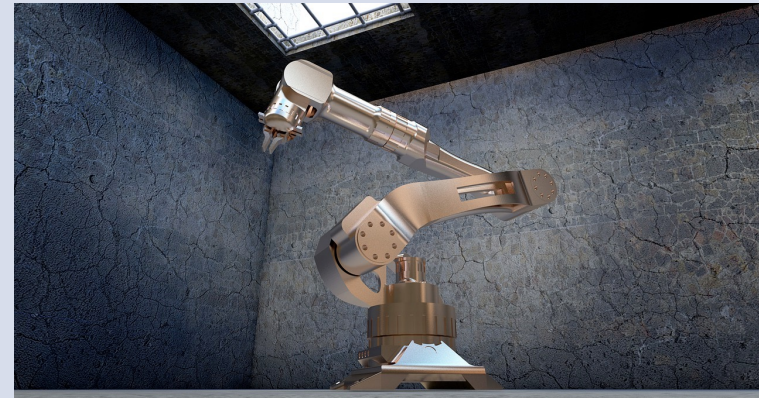
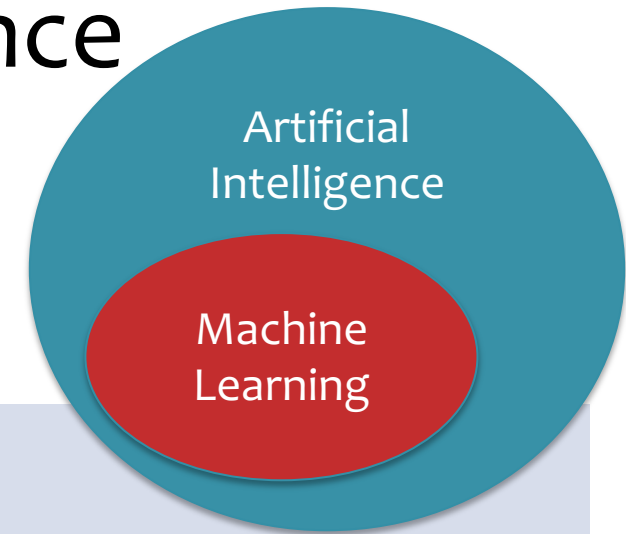


Artificial Intelligence

The basic goal of AI is to develop intelligent machines.

This consists of many sub-goals:

- Perception
- Reasoning
- Control / Motion / Manipulation
- Planning
- Communication
- Creativity
- Learning

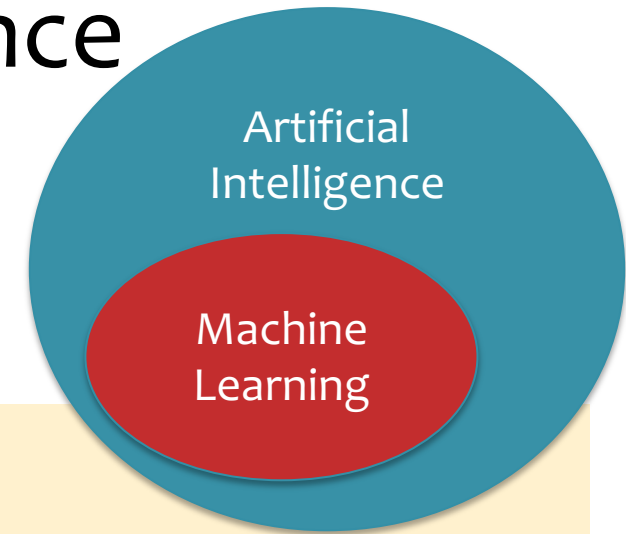


Artificial Intelligence

The basic goal of AI is to develop intelligent machines.

This consists of many sub-goals:

- Perception
- Reasoning
- Control / Motion / Manipulation
- Planning
- Communication
- Creativity
- Learning

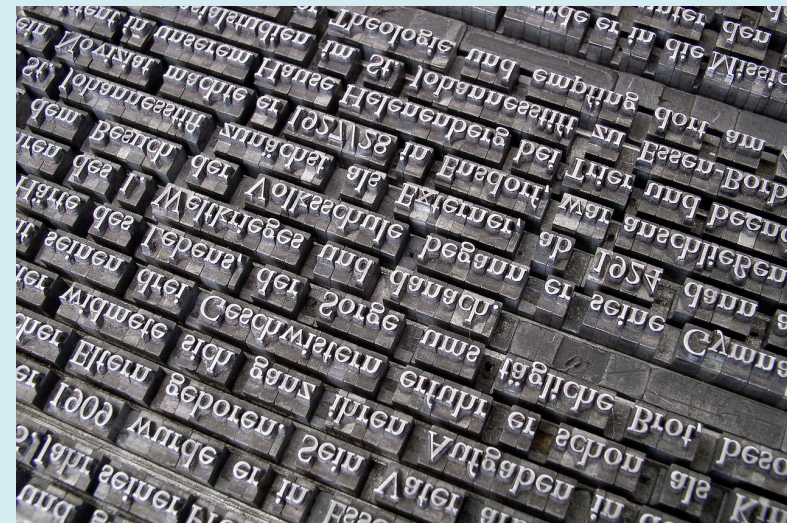
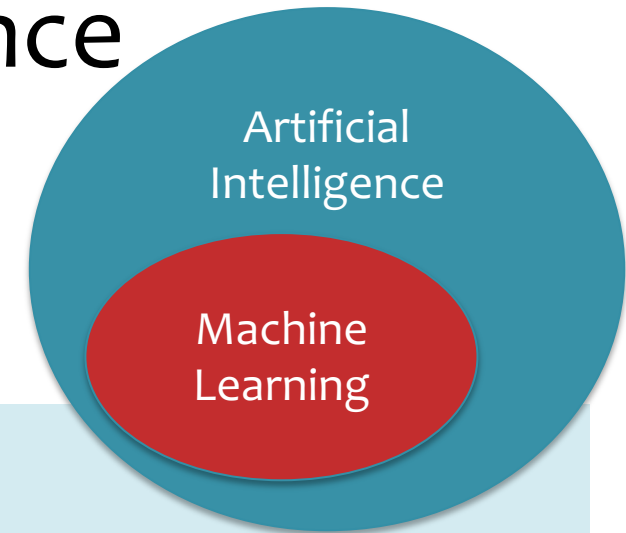


Artificial Intelligence

The basic goal of AI is to develop intelligent machines.

This consists of many sub-goals:

- Perception
- Reasoning
- Control / Motion / Manipulation
- Planning
- Communication
- Creativity
- Learning

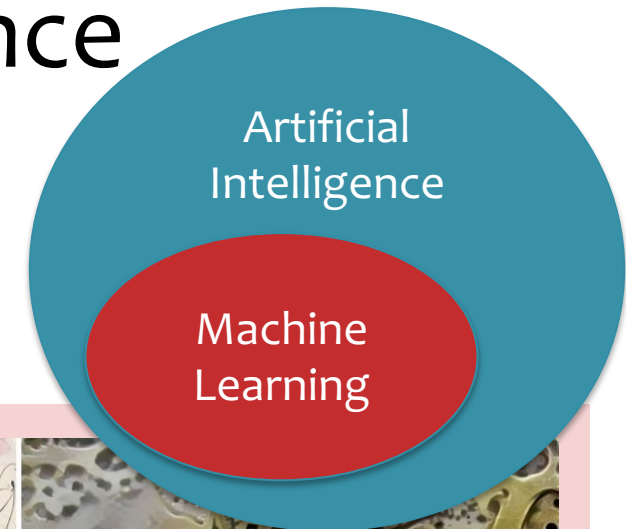


Artificial Intelligence

The basic goal of AI is to develop intelligent machines.

This consists of many sub-goals:

- Perception
- Reasoning
- Control / Motion / Manipulation
- Planning
- Communication
- Creativity
- Learning

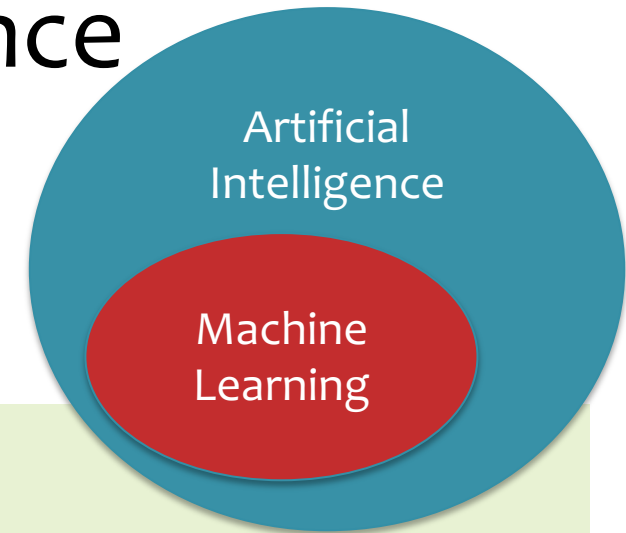


Artificial Intelligence

The basic goal of AI is to develop intelligent machines.

This consists of many sub-goals:

- Perception
- Reasoning
- Control / Motion / Manipulation
- Planning
- Communication
- Creativity
- Learning



What is Machine Learning?

The goal of this course is to provide you with a toolbox:

Machine Learning

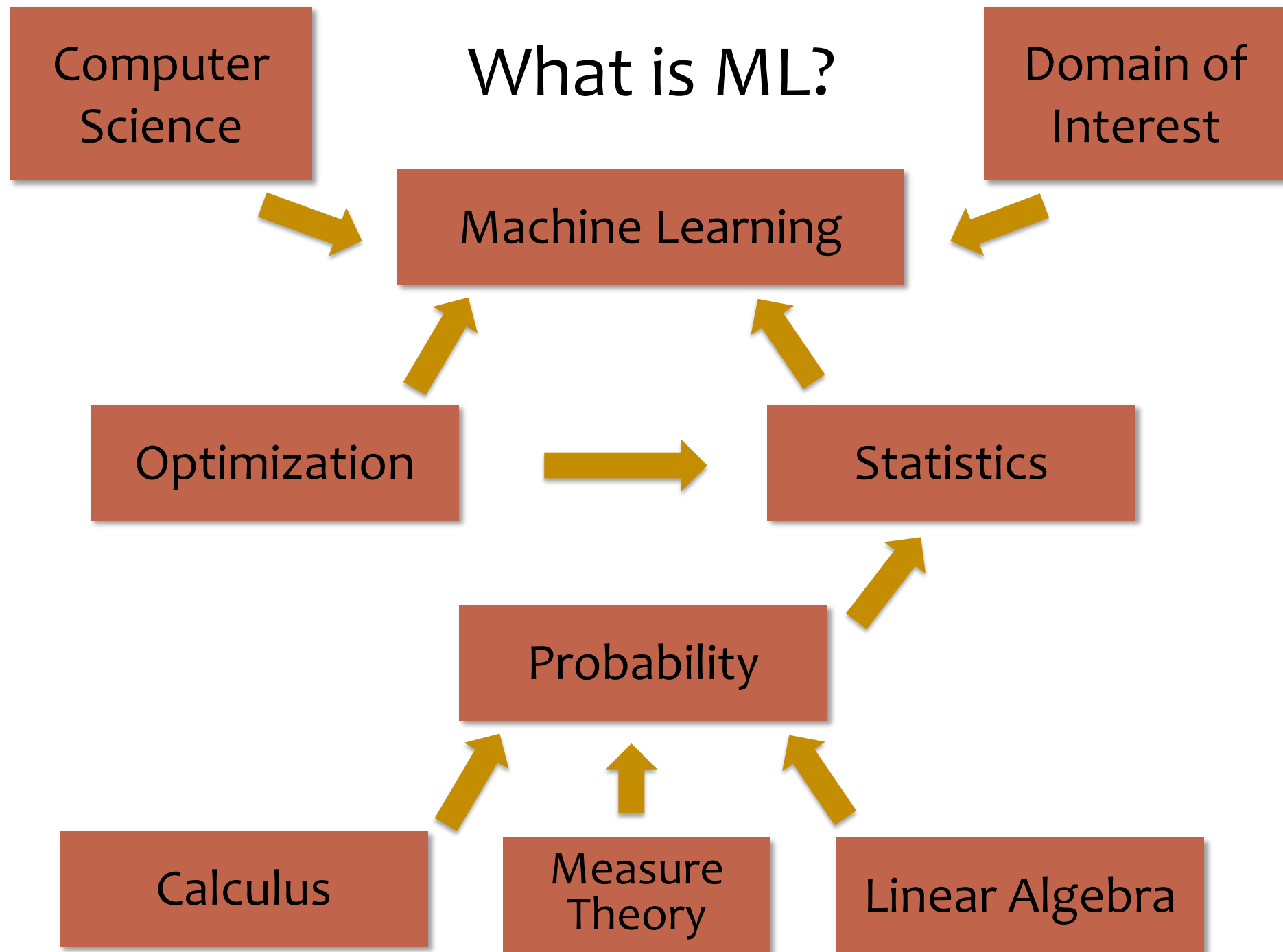
Statistics

Probability

Computer Science

Optimization





What is ML?

Speech Recognition

1. Learning to recognize spoken words

“...the SPHINX system (e.g. Lee 1989) learns speaker-specific strategies for recognizing the primitive sounds (phonemes) and words from the observed speech signal...neural network methods...hidden Markov models...”

(Mitchell, 1997)

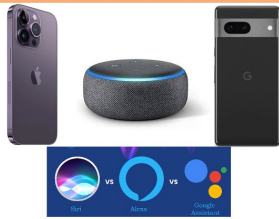


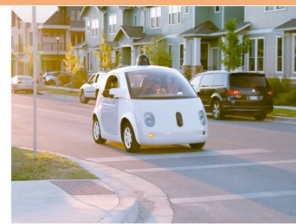
Figure from <https://botpenguin.com/alexa-vs-siri-vs-google-assistant/>

Robotics

2. Learning to drive an autonomous vehicle

“...the ALVINN system (Pomerleau 1989) has used its learned strategies to drive unassisted at 70 miles per hour for 90 miles on public highways among other cars...”

(Mitchell, 1997)



waymo.com

9

Games / Reasoning

3. Learning to beat the masters at board games

“...the world’s top computer program for backgammon, TD-GAMMON (Tesauro, 1992, 1995), learned its strategy by playing over one million practice games against itself...”

(Mitchell, 1997)



11

Computer Vision

4. Learning to recognize images

“...The recognizer is a convolution network that can be spatially replicated. From the network output, a hidden Markov model produces word scores. The entire system is globally trained to minimize word-level errors...”

(LeCun et al., 1995)

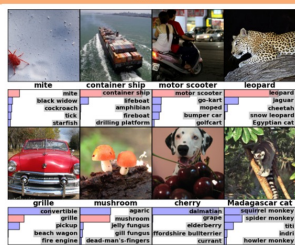


Figure from <https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf>

Learning Theory

- 5. In what cases and how well can we learn?

Sample Complexity Results

Definition 6.1. 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....

Finite $ N $	$N = \frac{1}{2} \log(\log(N)) + \log(\frac{1}{2})$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $k \in N$ with $R(k) \geq c$ have $R(k) > 0$.	$N = \frac{1}{2} \log(\log(N)) + \log(\frac{1}{2})$ labeled examples are sufficient so that with probability $(1 - \delta)$ for all $k \in N$ we have that $ R(k) - R(k) \leq \epsilon$.
Infinite $ N $	$N = O(\frac{1}{\epsilon} \log(N) \log(\frac{1}{\delta}) + \log(\frac{1}{\delta}))$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $k \in N$ with $R(k) \geq c$ have $R(k) > 0$.	$N = O(\frac{1}{\epsilon} \log(N) \log(\frac{1}{\delta}) + \log(\frac{1}{\delta}))$ labeled examples are sufficient so that with probability $(1 - \delta)$ for all $k \in N$ we have that $ R(k) - R(k) \leq \epsilon$.

Two Types of Error

① **True Error (aka. expected risk) (aka. Squared Loss Error)**
 $R(x) = P_{\text{avg}} \sum_{i=1}^n (x^{(i)}(n) + h(x^{(i)}))$ ← *analysis solution*

② **Test Error (aka. empirical risk)**
 $\hat{R}(x) = P_{\text{avg}} \sum_{i=1}^n (x^{(i)}(n) \neq h(x^{(i)}))$
 $= \frac{1}{N} \sum_{i=1}^N \mathbb{I}(x^{(i)}(n) \neq h(x^{(i)}))$
 $= \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y^{(i)} \neq h(x^{(i)}))$ ← *known, complete*

$S = \{x^{(1)}, \dots, x^{(N)}\}$

PAC Learning

Q: Can we bound $R(h)$ in terms of $\hat{R}(h)$?

A: Yes!

PAC states $R(h)$ is Probably Approximately Correct

PAC theorem yields hypothesis h , which is
 $R(h) \approx 0$
 with high probability $P(R(h) \approx 0) \approx 1$


Def: PAC Criterion

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

1. How many examples do we need to learn?
2. How do we quantify our ability to generalize to unseen data?
3. Which algorithms are better suited to specific learning settings?

Speech Recognition

1. Learning to recognize spoken words

THEN	NOW
<p>“...the SPHINX system (e.g. Lee 1989) learns speaker-specific strategies for recognizing the primitive sounds (phonemes) and words from the observed speech signal...neural network methods...hidden Markov models...”</p> <p>(Mitchell, 1997)</p>	 <p>The 'NOW' column features images of modern devices and AI assistants. At the top are an iPhone, an Amazon Echo, and a Google Pixel. Below them is a graphic comparing three AI assistants: Siri (represented by a colorful wave icon), Alexa (represented by a blue circular icon), and Google Assistant (represented by four colored dots). The graphic includes 'vs' text between the icons and labels for each assistant at the bottom.</p>

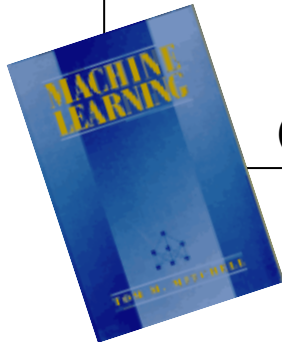
Robotics

2. Learning to drive an autonomous vehicle

THEN

“...the ALVINN system (Pomerleau 1989) has used its learned strategies to drive unassisted at 70 miles per hour for 90 miles on public highways among other cars...”

(Mitchell, 1997)



NOW



waymo.com

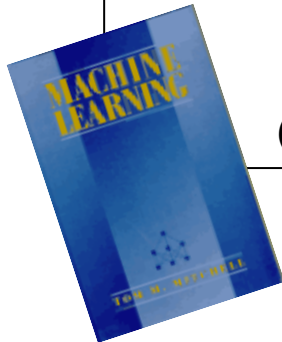
Games / Reasoning

3. Learning to beat the masters at board games

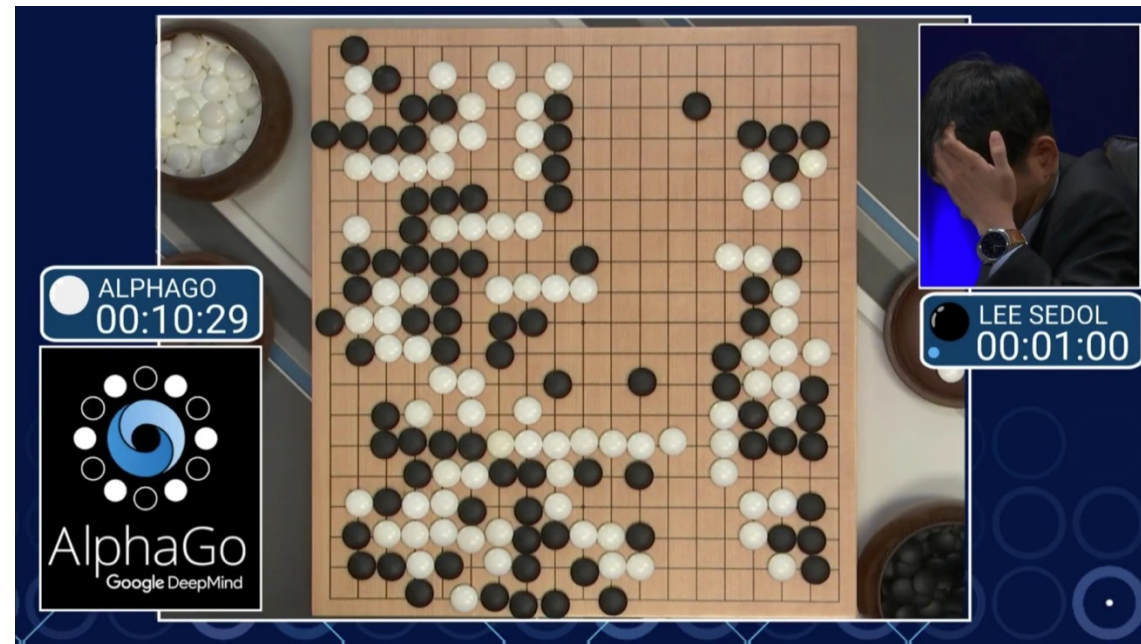
THEN

“...the world’s top computer program for backgammon, TD-GAMMON (Tesauro, 1992, 1995), learned its strategy by playing over one million practice games against itself...”

(Mitchell, 1997)



NOW

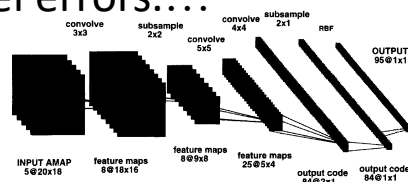


Computer Vision

4. Learning to recognize images

THEN

“...The recognizer is a convolution network that can be spatially replicated. From the network output, a hidden Markov model produces word scores. The entire system is globally trained to minimize word-level errors....”



(LeCun et al., 1995)

NOW



Learning Theory

• 5. In what cases and how well can we learn?

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$.

Two Types of Error

① True Error (aka. expected risk) (aka. Generalization Error)

$$R(h) = \mathbb{P}_{x \sim p^*(x)}(c^*(x) \neq h(x)) \quad \leftarrow \text{always unknown.}$$

② Train Error (aka. empirical risk)

$$\begin{aligned} \hat{R}(h) &= \mathbb{P}_{x \sim S}(c^*(x) \neq h(x)) \quad \leftarrow S = \{x^{(1)}, \dots, x^{(N)}\} \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{I}(c^*(x^{(i)}) \neq h(x^{(i)})) \quad \leftarrow \text{known, computable} \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y^{(i)} \neq h(x^{(i)})) \end{aligned}$$

PAC Learning

Q: Can we bound $R(h)$ in terms of $\hat{R}(h)$?

A: Yes!

PAC stands for Probably Approximately Correct

PAC learner yields hypothesis h , which is approximately correct $R(h) \approx 0$ with high probability $\Pr(R(h) \approx 0) \approx 1$

Def: PAC Criterion

$$\Pr(\forall h, |R(h) - \hat{R}(h)| \leq \epsilon) \geq 1 - \delta$$

1. How many examples do we need to learn?
2. How do we quantify our ability to generalize to unseen data?
3. Which algorithms are better suited to specific learning settings?

What is ML?

Speech Recognition

1. Learning to recognize spoken words

“...the SPHINX system (e.g., Lee 1989) learns speaker-specific strategies for recognizing the primitive sounds (phonemes) and words from the observed speech signal...neural network methods...hidden Markov models...”

(Mitchell, 1997)



Figure from <https://botpenguin.com/alexa-vs-siri-vs-google-assistant>

Robotics

2. Learning to drive an autonomous vehicle

“...the ALVINN system (Pomerleau 1989) has used its learned strategies to drive unassisted at 70 miles per hour for 90 miles on public highways among other cars...”

(Mitchell, 1997)



waymo.com

Games / Reasoning

3. Learning to beat the masters at board games

"...the world's top computer program for backgammon, TD-GAMMON (Tesauro, 1992, 1995), learned its strategy by playing over one million practice games against itself..."

(Mitchell, 199



Computer Vision

4. Learning to recognize images

"...The recognizer is a convolution network that can be spatially replicated. From the network output, a hidden Markov model produces word scores. The entire system is globally trained to minimize word-level errors..."

(LeCun et al., 1995)



Figure from <https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf>

Learning Theory

- 5. In what cases and how well can we learn?

Sample Complexity Results

Definition 6.1. 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....

Finite $ R $	Infinite $ R $
$N = \{a^i b^j \mid i+j \leq n\}$ is a finite set of strings. Since N is finite, $ N < \aleph_0$. N is not recursive. If it were, then R would be recursive, which contradicts the assumption that R is not recursive.	$N = \{a^i b^j \mid i+j \leq n\}$ is a finite set of strings. Since N is finite, $ N < \aleph_0$. N is not recursive. If it were, then R would be recursive, which contradicts the assumption that R is not recursive.

Two Types of Error

① Test Error (aka expected risk) (aka Generalization Error)
 $R(w) = P_{S \cup \mathcal{D}}(h^w(x) \neq h(x))$ ← always unknown

② Train Error (aka empirical risk)
 $\hat{R}(w) = P_{S \cup \mathcal{D}}(h^w(x) \neq h(w))$ ← $S = \{x^{(1)}, \dots, x^{(n)}\}$
 $= \frac{1}{N} \sum_{i=1}^N \mathbb{1}(h^w(x^{(i)}) \neq h(w^{(i)}))$
 $= \frac{1}{N} \sum_{i=1}^N \mathbb{1}(f) \neq h(w^{(i)})$ ← known, can be calculated

PAC Learning

Q: Can we bound $R(n)$ in terms of $\hat{R}(n)$?

A: Yes!

PAC states R is Probably Approximately Correct

PAC theorem yields hypothesis h , which is approximately correct $R(h) \approx 0$ with high probability $P(R(h) \approx 0) \approx 1$

Def: PAC Criterion

$\frac{1}{n} \sum_{i=1}^n \ell(w_i, \hat{y}_i) - \inf_{w \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(w_i, y_i) > 1 - \delta$

1. How many examples do we need to learn?
2. How do we quantify our ability to generalize to unseen data?
3. Which algorithms are better suited to specific learning settings?

What is Machine Learning?

The goal of this course is to provide you with a toolbox:

Machine Learning

Statistics

Probability

Computer Science

Optimization

To solve all the problems above and more



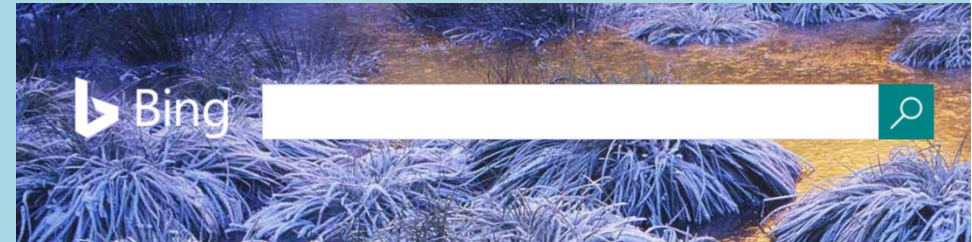
Societal Impacts of ML

What ethical responsibilities do we have as machine learning experts?

Question: What are the possible societal impacts of machine learning for each case below?

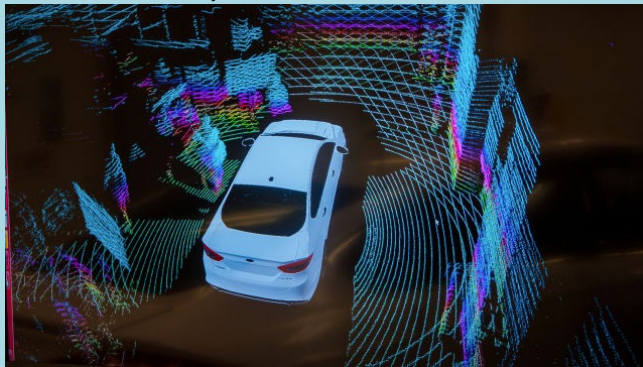
Answer:

1) Search results for news are optimized for ad revenue.



<http://bing.com/>

<http://arstechnica.com/>



2) An autonomous vehicle is permitted to drive unassisted on the road.

3) A doctor is prompted by an intelligent system with a plausible diagnosis for her patient.



<https://flic.kr/p/HNJUzV>

Societal Impacts of ML



A 72-year-old congressman goes back to school, pursuing a degree in AI



By [Meagan Flynn](#)

December 28, 2022 at 6:00 a.m. EST



Rep. Don Beyer (D-Va.) is pursuing a master's degree in machine learning at George Mason University with hopes of one day applying his AI knowledge to his legislative work. (Craig Hudson for The Washington Post)

Normally Don Beyer doesn't bring his multivariable calculus textbook to work, but his final exam was coming up that weekend.

"And I'm running out of time," he said, plopping the textbook and a scribbled notebook filled with esoteric-looking calculations on a coffee table in his office, "because I have all these—"

His phone was ringing. "I'll be there," Beyer told a colleague wondering when he would be returning to the House floor for votes.

It seemed study time would have to wait.

That's been the story of the year for Beyer (D-Va.), who has been moonlighting as a student at George Mason University in pursuit of a master's degree in machine learning while balancing his duties as a congressman. Beyer — a science wonk, economist and former car salesman — has been taking one class per semester in a slow but steady march toward the degree, with hopes of one day applying his artificial-intelligence knowledge to his legislative work as the technology evolves further.

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

Topics

- Foundations
 - Probability
 - Optimization
- Classification
 - KNN
 - Logistic Regression
 - Perceptron
- Regression
 - Linear Regression
- Important Concepts
 - Kernels
 - Regularization and Overfitting
 - Experimental Design
- Unsupervised Learning
 - K-means
 - PCA
- Neural Networks
 - Feedforward Neural Nets
 - Basic architectures
 - Backpropagation
- Deep Learning
 - CNNs
 - RNNs
 - Transformers
- Reinforcement Learning
 - Value Iteration / Policy Iteration
 - Q-Learning
 - Deep Q-Learning
- Learning Theory
 - PAC Learning
- Societal Impacts of ML
- Other Learning Paradigms
 - Matrix Factorization
 - Ensemble Methods

DEFINING LEARNING PROBLEMS

Well-Posed Learning Problems

Three components $\langle T, P, E \rangle$:

1. Task, T
2. Performance measure, P
3. Experience, E

Definition of learning:

A computer program **learns** if its performance at task T , as measured by P , improves with experience E .

Example Learning Problems

Learning to beat the masters at **chess**

1. Task, T : *play chess well*

2. Performance measure, P :

- *number of moves / time*
- *# pieces vs. opponent*
- *win rate*

- *measure of strategy quality?*
- *ELO rating*

3. Experience, E :

- *play chess matches against humans*
- *historical records of grandmaster games*
- *self-play*


Example Learning Problems

Learning to **respond to voice commands (Siri)**

1. Task, T :
2. Performance measure, P :
3. Experience, E :

Example Learning Problems

Learning to **respond to voice commands (Siri)**

1. Task, T : 

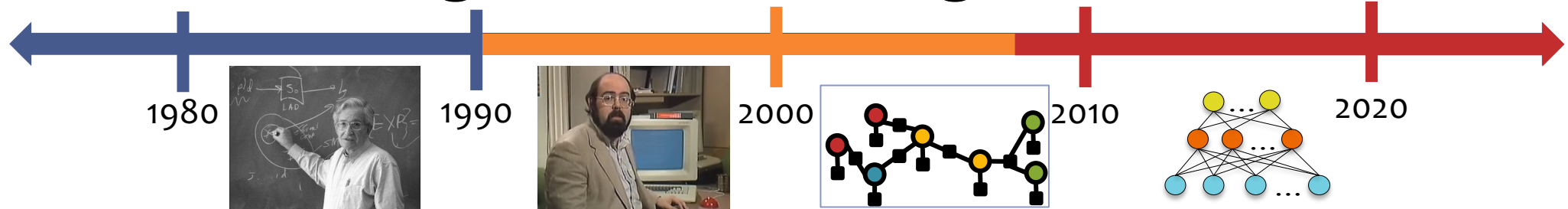
Given a transcribed sentence x predict the command y

Example:

x = "Give me directions to Starbucks"

y = DIRECTIONS(here, nearest(Starbucks))

Capturing the Knowledge of Experts



Solution #1: Expert Systems

- Over 20 years ago, we had rule-based systems:
 - Put a bunch of linguists in a room
 - Have them think about the structure of their native language and write down the rules they devise

Introspection...

`x = "Give me directions to Starbucks"`

`x = "Send Jill a txt asking for directions"`

`x = "Play the best hit music by TXT"`

`x = "How do I get to Pitt's Department of Music"`

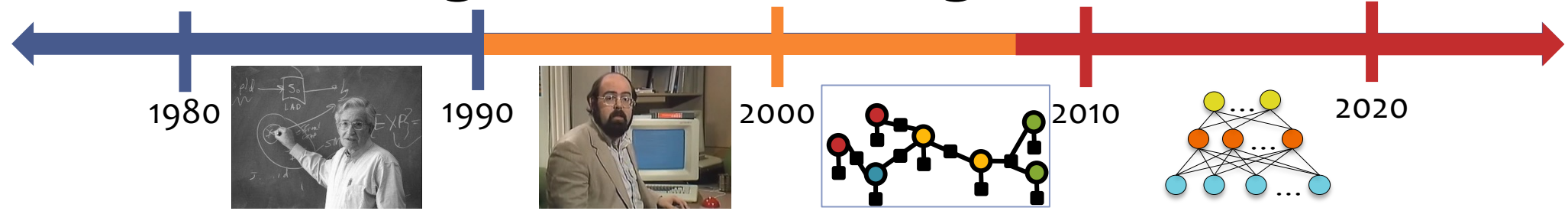
Rules...

~~`if "directions" in x:`~~
~~`—— type = DIRECTIONS()`~~

~~`if "txt" in x:`~~
~~`—— type = TXTMSG()`~~
~~`elif "directions" in x:`~~
~~`—— type = DIRECTIONS()`~~

~~`if "music" in x:`~~
~~`—— type = MUSIC()`~~
~~`elif "txt" in x:`~~
~~`—— type = TXTMSG()`~~
~~`elif "directions" in x:`~~
~~`—— type = DIRECTIONS()`~~

Capturing the Knowledge of Experts



Solution #1: Expert Systems

- Over 20 years ago, we had rule-based systems:
 - Put a bunch of linguists in a room
 - Have them think about the structure of their native language and write down the rules they devise

Introspection...

x = "Give me directions to Starbucks"

x = "How do I get to Starbucks?"

x = "Where is the nearest Starbucks?"

x = "I need directions to Starbucks"

x = "Is there a Starbucks nearby?"

x = "Starbucks now!"

Rules...

if x matches "give me directions to Z":
cmd = DIRECTIONS(here, nearest(Z))

if x matches "how do i get to Z":
cmd = DIRECTIONS(here, nearest(Z))

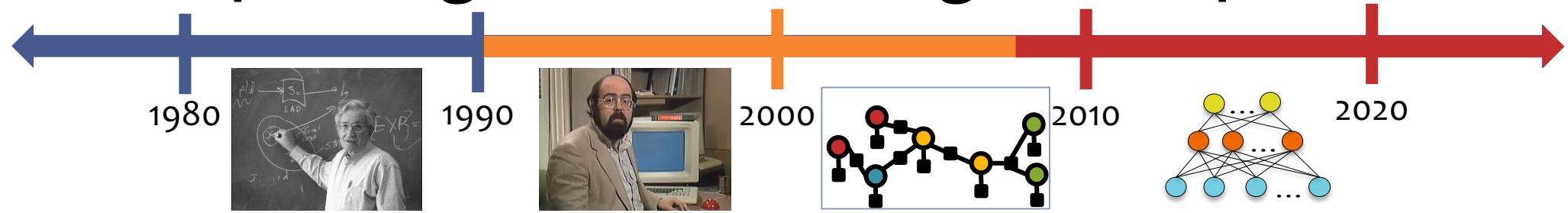
if x matches "where is the nearest Z":
cmd = DIRECTIONS(here, nearest(Z))

if x matches "I need directions to Z":
cmd = DIRECTIONS(here, nearest(Z))

if x matches "Is there a Z nearby":
cmd = DIRECTIONS(here, nearest(Z))

if x matches "Z now!":
cmd = DIRECTIONS(here, nearest(Z))

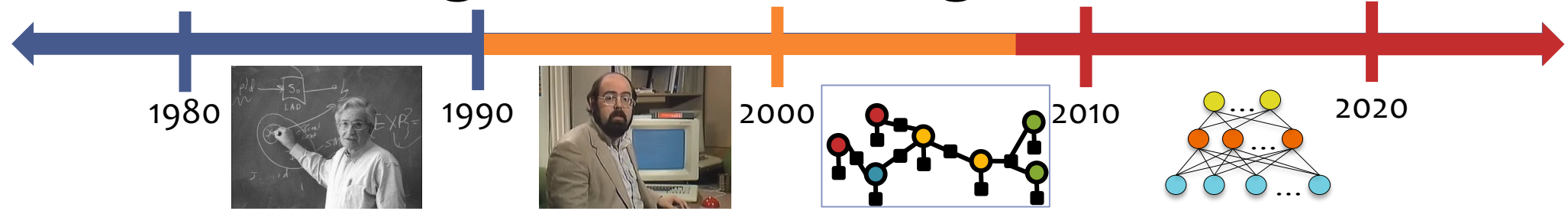
Capturing the Knowledge of Experts



Solution #2: Annotate Data and Learn

- Experts:
 - **Very good** at answering questions about specific cases
 - **Not very good** at telling **HOW** they do it
- 1990s: So why not just have them tell you what they do on **SPECIFIC CASES** and then let **MACHINE LEARNING** tell you how to come to the same decisions that they did

Capturing the Knowledge of Experts



Solution #2: Annotate Data and Learn

1. Collect raw sentences $\{x^{(1)}, \dots, x^{(n)}\}$
2. Experts annotate their meaning $\{y^{(1)}, \dots, y^{(n)}\}$

$x^{(1)}$: How do I get to Starbucks?

$y^{(1)}$: DIRECTIONS(here, nearest(Starbucks))

$x^{(3)}$: Send a text to John that I'll be late

$y^{(3)}$: TXTNSG(John, I'll be late)

$x^{(2)}$: Show me the closest Starbucks

$y^{(2)}$: MAP(nearest(Starbucks))

$x^{(4)}$: Set an alarm for seven in the morning

$y^{(4)}$: SETALARM(7:00AM)

Example Learning Problems

Learning to **respond to voice commands (Siri)**

1. Task, T :
predicting action from speech
2. Performance measure, P :
percent of correct actions taken in user pilot study
3. Experience, E :
examples of (speech, action) pairs

Problem Formulation

Often, the same task can be formulated in more than one way.

Example: Loan applications

- creditworthiness/score (regression)
- probability of default (density estimation)
- loan decision (classification)

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)

Well-posed Learning Problems



In-Class Exercise

1. Select a **task**, T
2. Identify **performance measure**, P
3. Identify **experience**, E
4. Report ideas back to rest of class

Example Tasks

- Identify objects in an image
- Translate from one human language to another
- Recognize speech
- Assess risk (e.g. in loan application)
- Make decisions (e.g. in loan application)
- Assess potential (e.g. in admission decisions)
- Categorize a complex situation (e.g. medical diagnosis)
- Predict outcome (e.g. medical prognosis, stock prices, inflation, temperature)
- Predict events (default on loans, quitting school, war)
- Plan ahead under perfect knowledge (chess)
- Plan ahead under partial knowledge (poker, bridge)

In-Class Exercise

1. Select a **task**, T
2. Identify **performance measure**, P
3. Identify **experience**, E
4. Report ideas back to rest of class

Well-posed Learning Problems

task, T	performance measure, P	experience, E
detect tumors	# detected, correctness string	scans from multiple hospitals, expert labels
product recommended	did they buy click	also features viewing history other customer decisions

In-Class Exercise

1. Select a **task**, T
2. Identify **performance measure**, P
3. Identify **experience**, E
4. Report ideas back to rest of class


Well-posed Learning Problems

task, T	performance measure, P	experience, E

(without any math!)

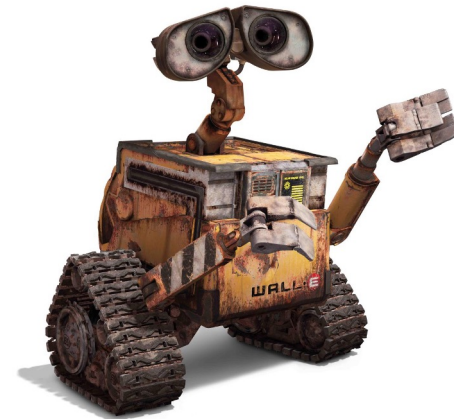
SUPERVISED LEARNING

Building a Trash Classifier

- Suppose the  ask CMU to build a robot for collecting trash along Pittsburgh's rivers
- You are tasked with building a classifier that detects whether an object is a piece of trash (+) or not a piece of trash (-)
- The robot can detect an object's *color*, *sound*, and *weight*
- You manually annotate the following dataset based on objects you find



trash?	color	sound	weight
+	green	crinkly	high
-	brown	crinkly	low
-	grey	none	high
+	clear	none	low
-	green	none	low



WARNING!

Like many fields, Machine Learning is riddled with copious amounts of technical jargon!

For many terms we'll define in this class, you'll find four or five different terms in the literature that refer to the same thing.

Supervised Binary Classification

- Def: an **example** contains a **label** (aka. **class**) and **features** (aka. **point** or **attributes**)
- Def: a **labeled dataset** consists of rows, where each row is an example
- Def: an **unlabeled dataset** only has **features**

One example:

label		features	
trash?	color	sound	weight
-	brown	none	high

Labeled Dataset:

	label	features		
index	trash?	color	sound	weight
1	-	brown	none	high
2	+	clear	crinkly	low
3	-	brown	none	low

Unlabeled Dataset:

	features		
index	color	sound	weight
1	brown	none	high
2	clear	crinkly	low
3	brown	none	low

Supervised Binary Classification

- Def: an **example** contains a **label** (aka. **class**) and **features** (aka. **point** or **attributes**)
- Def: a **labeled dataset** consists of rows, where each row is an example
- Def: an **unlabeled dataset** only has **features**
- Def: a **classifier** is a function that takes in features and predicts a label
- Def: a **training dataset** is a labeled dataset used to **learn** a classifier
- Def: a **test dataset** is a labeled dataset used to **evaluate** a classifier

Classifier
features → label

Training Dataset:

	label	features		
index	trash?	color	sound	weight
1	+	green	crinkly	high
2	-	brown	crinkly	low
3	-	grey	none	high
4	+	clear	none	low
5	-	green	none	low

Test Dataset:

	label	features		
index	trash?	color	sound	weight
1	-	brown	none	high
2	+	clear	crinkly	low
3	-	brown	none	low

Supervised Binary Classification

- Def: **predictions** are the output of a trained classifier
- Def: **error rate** is the proportion of examples on which we predicted the wrong label
- Def: a **classifier** is a function that takes in features and predicts a label
- Def: a **training dataset** is a labeled dataset used to **learn** a classifier
- Def: a **test dataset** is a labeled dataset used to **evaluate** a classifier

Classifier
features → label

Test Predictions:		(Unlabeled) Test Dataset:			
		features			
index	trash?	index	color	sound	weight
1	+	1	brown	none	high
2	+	2	clear	crinkly	low
3	-	3	brown	none	low

Supervised Binary Classification

- Def: **predictions** are the output of a trained classifier
- Def: **error rate** is the proportion of examples on which we predicted the wrong label
- Def: a **classifier** is a function that takes in features and predicts a label
- Def: a **training dataset** is a labeled dataset used to **learn** a classifier
- Def: a **test dataset** is a labeled dataset used to **evaluate** a classifier

error rate = $1/3$

Test Predictions:

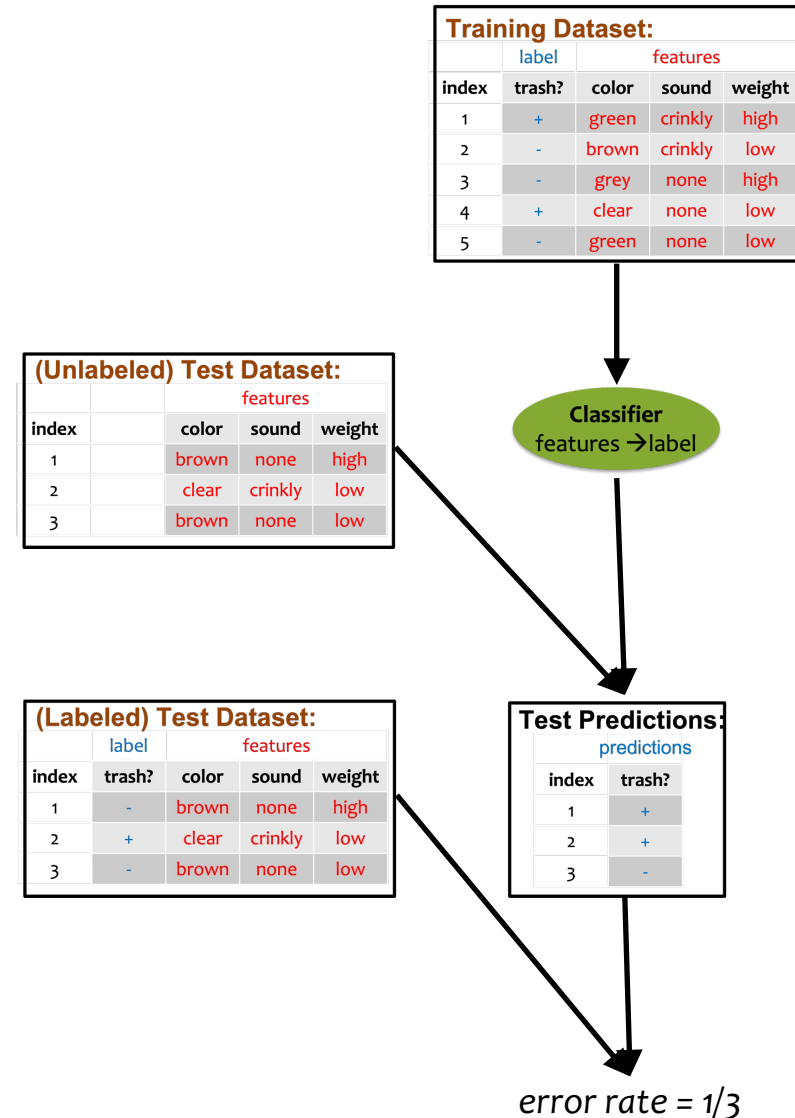
		predictions
index		trash?
1	✗	+
2	✓	+
3	✓	-

(Labeled) Test Dataset:

	label	features		
index	trash?	color	sound	weight
1	-	brown	none	high
2	+	clear	crinkly	low
3	-	brown	none	low

Supervised Binary Classification

- Step 1: training
 - *Given:* labeled **training dataset**
 - *Goal:* learn a **classifier** from the training dataset
- Step 2: prediction
 - *Given:* unlabeled **test dataset**
: learned classifier
 - *Goal:* **predict** a label for each instance
- Step 3: evaluation
 - *Given:* **predictions** from Step II
: labeled **test dataset**
 - *Goal:* compute the **test error rate** (i.e. error rate on the test dataset)



Supervised Binary Classification

- Step 1: training
 - Given: labeled **training dataset**
 - Goal: learn a **classifier** from the training dataset
- Step 2: prediction
 - Given: unlabeled **test dataset**
: learned classifier
 - Goal: **predict** a label for each instance
- Step 3: evaluation
 - Given: **predictions** from Step II
: labeled **test dataset**
 - Goal: compute the **test error rate** (i.e. error rate on the test dataset)

Training Dataset:				
	label	features		
index	trash?	color	sound	weight
1	+	green	crinkly	high
2	-	brown	crinkly	low
3	-	grey	none	high
4	+	clear	none	low
5	-	green	none	low

Classifier
features → label

(Unlabeled) Test Dataset:				
	features			
index		color	sound	weight
1		brown	none	high
2		clear	crinkly	low
3		brown	none	low

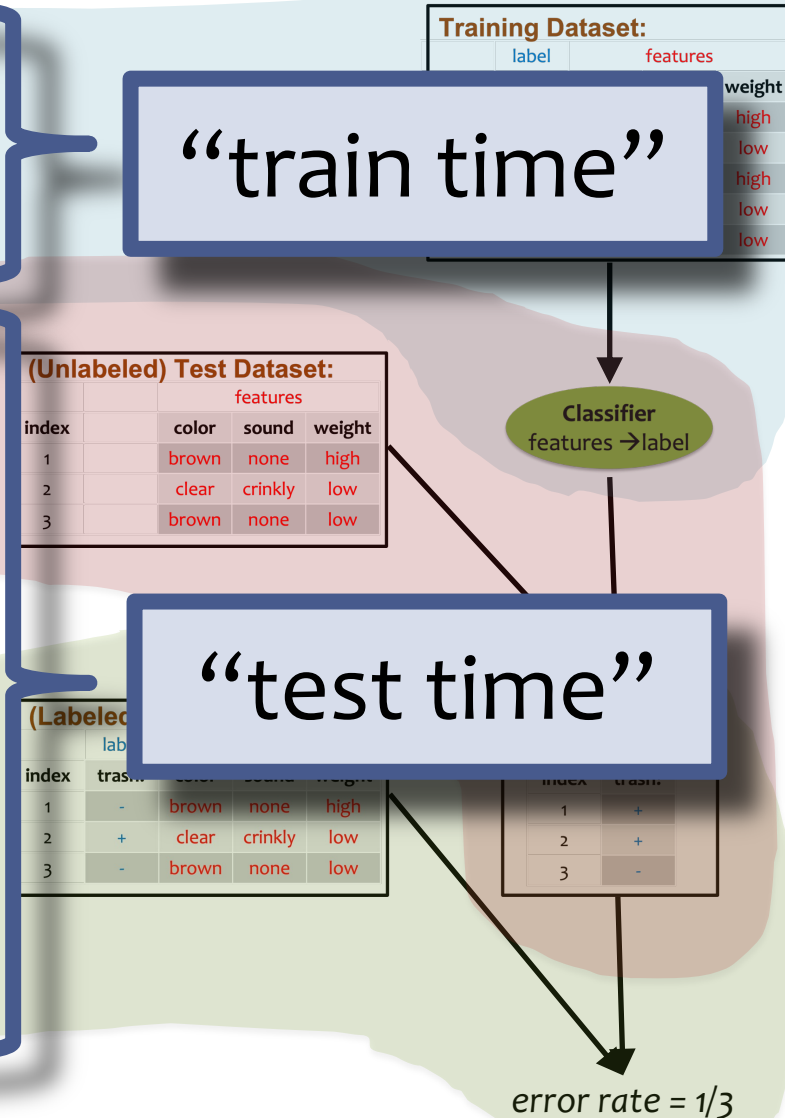
(Labeled) Test Dataset:				
	label	features		
index	trash?	color	sound	weight
1	-	brown	none	high
2	+	clear	crinkly	low
3	-	brown	none	low

Test Predictions:	
	predictions
index	trash?
1	+
2	+
3	-

error rate = 1/3

Supervised Binary Classification

- Step 1: training
 - Given: labeled **training dataset**
 - Goal: learn a **classifier** from the training dataset
- Step 2: prediction
 - Given: unlabeled **test dataset**
: learned classifier
 - Goal: **predict** a label for each instance
- Step 3: evaluation
 - Given: **predictions** from Phase II
: labeled **test dataset**
 - Goal: compute the **test error rate** (i.e. error rate on the test dataset)



Supervised Binary Classification

-
- Step 1: training
 - Given: labeled **training dataset**
 - Goal: learn a **classifier** from the training dataset
 - Step 2: prediction
 - Given: unlabeled **test data**; learned classifier
 - Goal: **predict** a label for each instance
 - Step 3: evaluation
 - Given: **predictions** from ; labeled **test datasets**
 - Goal: compute the **test error rate** (i.e. error rate on the test dataset)

Training Dataset:				
	label	features		
index	trash?	color	sound	weight
1	+	green	crinkly	high
2	-	brown	crinkly	low
3	-	grey	none	high
4	+	clear	none	low
5	-	green	none	low

Key question in Machine Learning:

How do we learn the classifier from data?

Random Classifier

The **random classifier** takes in the features and always predicts a random label.

... this is a terrible idea. It completely **ignores the training data!**

Classifier
features → random!

Training Dataset:

index	label	features		
	trash?	color	sound	weight
1	+	green	crinkly	high
2	-	brown	crinkly	low
3	-	grey	none	high
4	+	clear	none	low
5	-	green	none	low

Test Predictions:

predictions

index	trash?
1	-
2	-
3	+

error rate = 2/3

Test Dataset:

index	label	features		
	trash?	color	sound	weight
1	-	brown	none	high
2	+	clear	crinkly	low
3	-	brown	none	low

Random Classifier

The **random classifier** takes in the features and always predicts a random label.

... this is a terrible idea. It completely **ignores the training data!**

Classifier
features → random!

Training Dataset:

index	label	features		
	trash?	color	sound	weight
1	+	green	crinkly	high
2	-	brown	crinkly	low
3	-	grey	none	high
4	+	clear	none	low
5	-	green	none	low

Test Predictions:

predictions

index	trash?
1	+
2	+
3	-

error rate = 1/3

Test Dataset:

index	label	features		
	trash?	color	sound	weight
1	-	brown	none	high
2	+	clear	crinkly	low
3	-	brown	none	low

Random Classifier

The **random classifier** takes in the features and always predicts a random label.

... this is a terrible idea. It completely **ignores the training data!**

Classifier
features → random!

Training Dataset:

index	label	features		
	trash?	color	sound	weight
1	+	green	crinkly	high
2	-	brown	crinkly	low
3	-	grey	none	high
4	+	clear	none	low
5	-	green	none	low

Test Predictions:

predictions

index	trash?
1	+
2	-
3	+

error rate = 3/3

Test Dataset:

index	label	features		
	trash?	color	sound	weight
1	-	brown	none	high
2	+	clear	crinkly	low
3	-	brown	none	low

Majority Vote Classifier

The **majority vote classifier** takes in the features and always predicts the **most common label** in the training dataset.

... this is still a pretty bad idea. It completely **ignores the features!**

Classifier
features → always predict “-”

Training Dataset:

index	label	features		
	trash?	color	sound	weight
1	+	green	crinkly	high
2	-	brown	crinkly	low
3	-	grey	none	high
4	+	clear	none	low
5	-	green	none	low

Test Predictions:

predictions

index	trash?
1	-
2	-
3	-

error rate = 1/3

Test Dataset:

index	label	features		
	trash?	color	sound	weight
1	-	brown	none	high
2	+	clear	crinkly	low
3	-	brown	none	low

Majority Vote Classifier

The **majority vote classifier** takes in the features and always predicts the **most common label** in the training dataset.

...this is still a pretty bad idea. It completely **ignores the features!**

Classifier
features → always predict “-”

The majority vote classifier even ignores the features if it's making predictions on the training dataset!

Train Predictions:

predictions

index	trash?
1	-
2	-
3	-
4	-
5	-

error rate = 2/5

Training Dataset:

index	label	features		
	trash?	color	sound	weight
1	+	green	crinkly	high
2	-	brown	crinkly	low
3	-	grey	none	high
4	+	clear	none	low
5	-	green	none	low

Majority Vote Classifier

- Step 1: training
 - Given: labeled **training dataset**
 - Goal: learn a **classifier** from the training dataset
- Step 2: prediction
 - Given: unlabeled **test dataset**
: learned classifier
 - Goal: **predict** a label for each instance
- Step 3: evaluation
 - Given: **predictions** from Step II
: labeled **test dataset**
 - Goal: compute the **test error rate** (i.e. error rate on the test dataset)

Training Dataset:				
	label	features		
index	trash?	color	sound	weight
1	+	green	crinkly	high
2	-	brown	crinkly	low
3	-	grey	none	high
4	+	clear	none	low
5	-	green	none	low

(Unlabeled) Test Dataset:				
		features		
index		color	sound	weight
1		brown	none	high
2		clear	crinkly	low
3		brown	none	low

Classifier
features → always predict “-”

(Labeled) Test Dataset:				
	label	features		
index	trash?	color	sound	weight
1	-	brown	none	high
2	+	clear	crinkly	low
3	-	brown	none	low

Test Predictions:	
	predictions
index	trash?
1	-
2	-
3	-

error rate = 1/3

SYLLABUS HIGHLIGHTS

Syllabus Highlights

The syllabus is located on the course webpage:

<http://www.cs.cmu.edu/~mgormley/courses/10601>

or

<http://mlcourse.org>

The **course policies** are **required** reading.

Syllabus Highlights

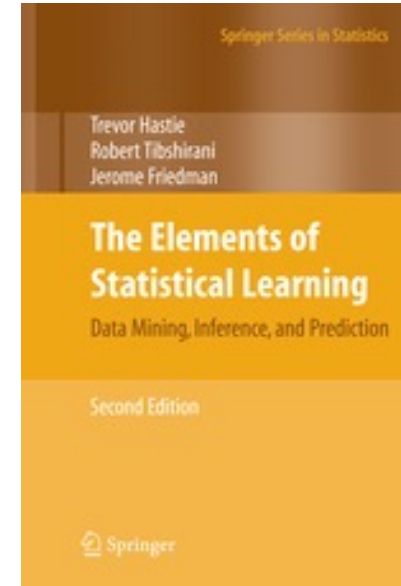
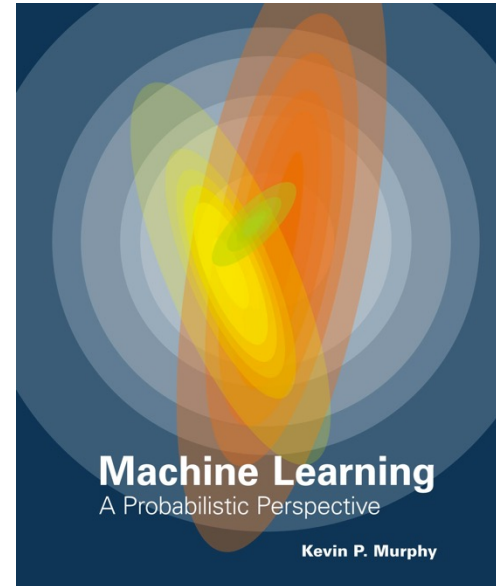
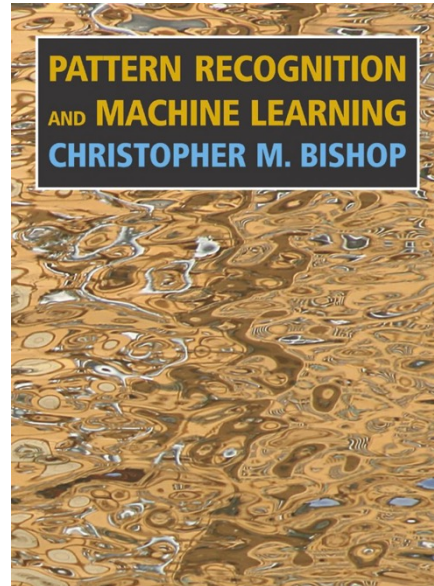
- **Grading:** 45% homework, 45% exams, 5% quizzes, 5% participation
- **Quizzes:** in-class, programming focused
- **Exam 1:** evening, Mon, Sep. 29
- **Exam 2:** evening, Thu, Nov. 6
- **Exam 3:** final exam week, date TBD by registrar
- **Homework:** 3 written and 6 written + programming (Python)
 - 6 grace days for homework assignments
 - Late submissions: 75% day 1, 50% day 2, 25% day 3
 - No submissions accepted after 3 days w/o extension; HW3, HW6, HW9 only 2 days
 - Extension requests: for emergencies, see syllabus
- **Recitations:** Fridays, same time/place as lecture (optional, interactive sessions)
- **Readings:** required, online PDFs, recommended for after lecture
- **Technologies:** Piazza (discussion), Gradescope (homework), Google Forms (polls)
- **Academic Integrity:**
 - Collaboration encouraged, but must be documented
 - Solutions must always be written independently
 - No re-use of found code / past assignments
 - Severe penalties (e.g. -100%)
- **Office Hours:** posted on Google Calendar on “Office Hours” page

Lectures



- You should ask lots of questions
 - Interrupting (by raising a hand) to ask your question is strongly encouraged
 - Asking questions later (or in real time) on Piazza is also great
- When I ask a question...
 - I want you to answer
 - Even if you don't answer, think it through as though I'm about to call on you
- Interaction improves learning (both in-class and at my office hours)

Textbooks

You are not *required* to read a textbook, but it will help immensely!



Where can I find...?

Home FAQ Syllabus People Schedule Office Hours Coursework Previous Links ▾							
Date	Lecture		Readings		Announcements		
Classification & Regression							
Mon, 1-Feb	Lecture 1: Course Overview [Slides]		<ul style="list-style-type: none">• 10601 Notation Crib Sheet. Matt Gormley (2018).• Command Line and File I/O Tutorial. 10601 Course Staff (2020).• 10601 Learning Objectives. Matt Gormley (2018).• Visual Information Theory. Christopher Olah (2015). blog.				
Wed, 3-Feb	Lecture 2: Decision Trees, Overfitting [Slides]		<ul style="list-style-type: none">• Decision Trees. Hal Daumé III (2017). CIML, Chapter 1.		HW1 out		
Fri, 5-Feb	Recitation: HW1 [Handout] [Solutions]						
Mon, 8-Feb	Lecture 3: Generalizing from examples - the Big Picture [Slides] [Poll]		<ul style="list-style-type: none">• Limits of Learning. Hal Daumé III (2017). CIML, Chapter 2.				
Wed, 10-Feb	Lecture 4: k-Nearest Neighbors [Slides] [Whiteboard] [Poll]		<ul style="list-style-type: none">• Geometry and Nearest Neighbors. Hal Daumé III (2017). CIML, Chapter 3.		HW1 due HW2 out		
Fri, 12-Feb	Recitation: HW2 [Handout] [Solutions]						
Mon, 15-Feb	Lecture 5: Model Selection [Slides] [Whiteboard] [Poll]						
Wed, 17-Feb	Lecture 6: Perceptron [Slides] [Whiteboard] [Poll]		<ul style="list-style-type: none">• The Perceptron. Hal Daumé III (2017). CIML, Chapter 4.		HW1 solution session (Thursday)		

Where can I find...?

Home FAQ Syllabus People Schedule **Office Hours** Coursework Previous Links ▾

Introduction to Machine Learning

10-301 + 10-601,
School of Computer Science
Carnegie Mellon University

10-301/601 Office Hours

Today Aug 29 – Sep 4, [Print](#) [Week](#) [Month](#) [Agenda](#)

	Sun 8/29	Mon 8/30	Tue 8/31	Wed 9/1	Thu 9/2	Fri 9/3	Sat 9/4
9am							
10am		10:10 – 11:30 10-301/601 Section A/C Mellon Institute		10:10 – 11:30 10-301/601 Section A/C Mellon Institute		10:10 – 11:30 10-301/601 Section A/C Mellon Institute	
11am		11:30 – Matt and H		11:30 – Matt and H			
12pm							
1pm		1:25p – 2:45p 10-301/601 Section B/D CUC McConomy		1:25p – 2:45p 10-301/601 Section B/D CUC McConomy		1:25p – 2:45p 10-301/601 Section B/D CUC McConomy	
2pm		2:45p – Matt and H		2:45p – Matt and H			
3pm							

Where can I find...?

[Home](#)[FAQ](#)[Syllabus](#)[People](#)[Schedule](#)[Office Hours](#)[Coursework](#)[Previous](#)[Links ▾](#)

Introduction to Machine Learning

10-301 + 10-601, §
School of Computer Science
Carnegie Mellon U

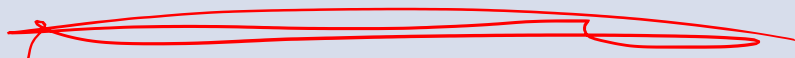
Assignments

There will be 8 homework assignments during the semester in addition to the exams. The assignments will consist of both theoretical and programming assignments. The assignments will be released via a Piazza announcement explaining where to find the handout, starter code, LaTeX template, etc.

- Homework 1: Background Material (written / programming)
[Handout](#)
- Homework 2: Decision Trees (written / programming)
[Handout](#)
- Homework 3: KNN, Perceptron, and Linear Regression (written)
[Handout](#)
- Mock Exam 1:
[Handout and Solution](#)
- Homework 4: Logistic Regression (written / programming)
[Handout](#)
- Homework 5: Neural Networks (written / programming)
[Handout](#)
- Homework 6: Neural Networks and Reinforcement Learning (written / programming)
[Handout](#)
- Homework 7: Graphical Models (written / programming)

In-Class Polls

Q: How do these In-Class Polls work?

A: Don't worry about it for today. We won't use them until the second week of class, i.e. the third lecture. 

Details are on the syllabus.

PREREQUISITES

Prerequisites

What they are:

- Significant programming experience (15-122)
 - Written programs of 100s of lines of code
 - Comfortable learning a new language
- Probability and statistics (36-217, 36-225, etc.)
- Mathematical maturity: discrete mathematics (21-127, 15-151), OR linear algebra (21240, 21241) OR calculus (21259, 21254)

Prerequisites

What if you need additional review?

- Consider first taking 10-606/607: Mathematical/Computational Foundations for Machine Learning

How to describe 606/607 to a friend

606/607 is...

a **formal** presentation of **mathematics** and **computer science**...

motivated by (carefully chosen) **real-world problems** that arise in **machine learning**...

where the **broader picture** of how those problems arise is treated **somewhat informally**.

Prerequisites

What if I'm not sure whether I meet them?

- Don't worry: we're not sure either
- However, we've designed a way to assess your background knowledge so that you know what to study!

Reminders

- **Homework 1: Background**
 - **Out: Mon, Aug 25**
 - **Due: Wed, Sep 3 at 11:59pm**
 - Two parts:
 1. written part to Gradescope
 2. programming part to Gradescope

Learning Objectives

You should be able to...

1. Formulate a well-posed learning problem for a real-world task by identifying the task, performance measure, and training experience
2. Describe common learning paradigms in terms of the type of data available, when it's available, the form of prediction, and the structure of the output prediction
3. Implement Decision Tree training and prediction (w/simple scoring function)
4. Explain the difference between memorization and generalization [CIML]
5. Identify examples of the ethical responsibilities of an ML expert