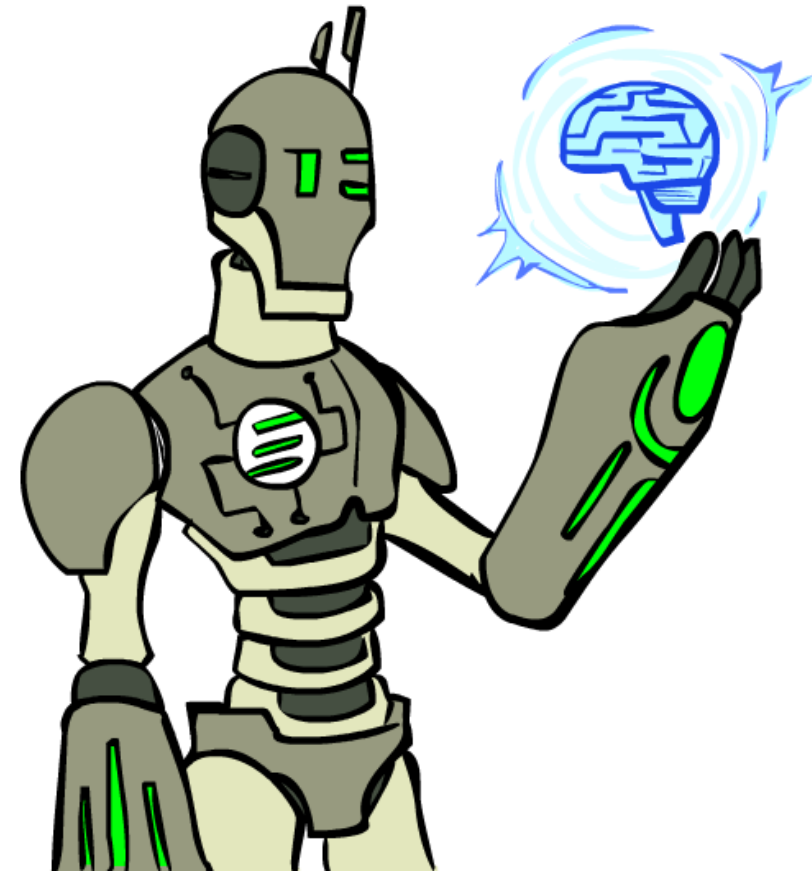
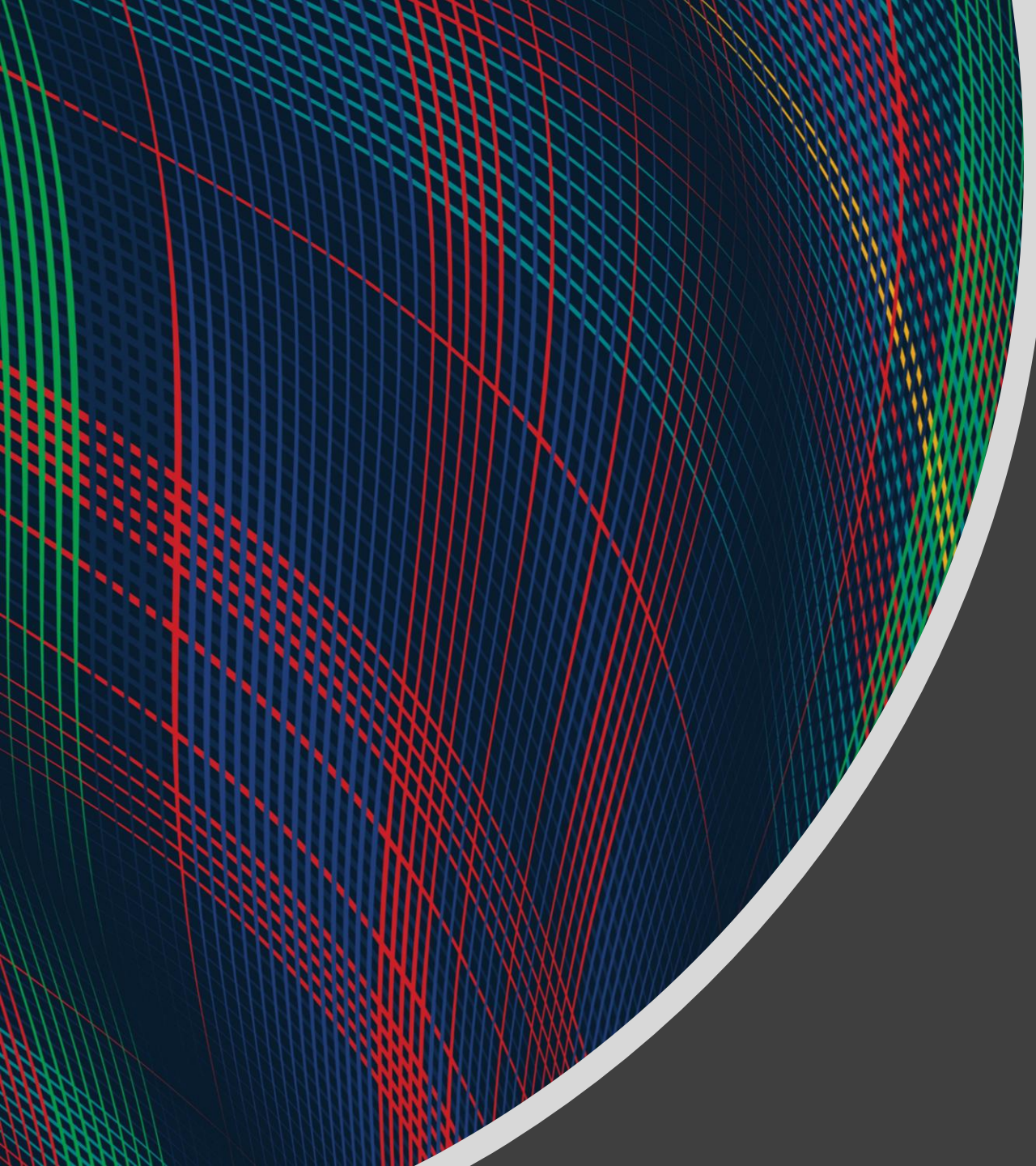


As you walk in

Welcome!

- 1) Sit at a table next to another student
- 2) Make name plate
 - Fold paper in half
 - Write preferred name
 - Below write you favorite fictional AI/robot



An abstract graphic on the left side of the slide, featuring a sphere-like shape composed of a dense grid of intersecting red, green, and blue lines. The lines are curved and follow the contour of the sphere, creating a complex, woven pattern. The sphere is set against a dark gray background.

10-607 Computational Foundations for Machine Learning

Instructor: Pat Virtue

Today

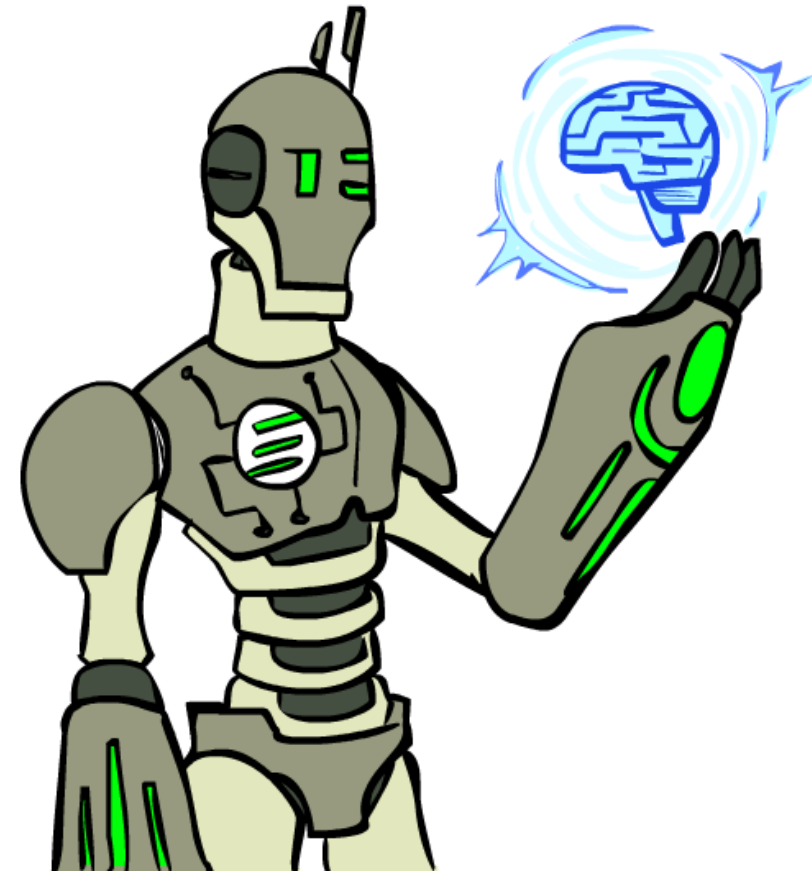
Course Info

Warm-up exercise

Propositional Logic and Proofs

ML and 606/607 Intro

More Course Info



Course Team

Instructor



Pat
Virtue
pvirtue

Teaching Assistants



Kellen
Gibson
kagibson



Ian
Char
ichar

Course Team

Students!!



Team Tips

Try not to act surprised

Here's a thing that happens a lot:



Team Tips

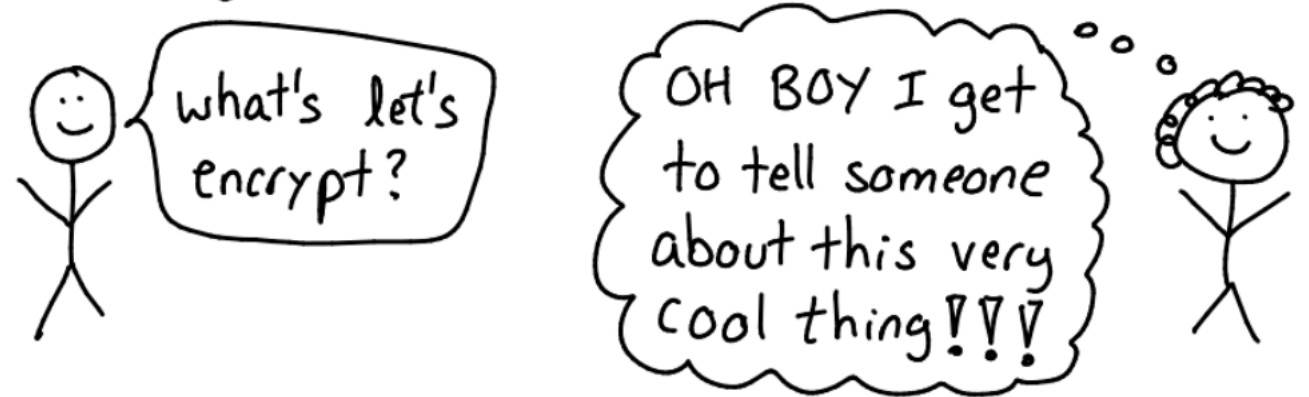
Try not to act surprised

Here's a cool simple trick !

Don't act surprised when someone doesn't know something you thought they knew

(even if you are a little surprised !)
It doesn't help.

Then you get to have fun times like this:



And it gets easier with practice ! ☺☺☺

Two-column Proof

Give an explicit justification for each statement based on previous statements

Prove Socrates is mortal

Notation Alert!

Modus Ponens

$$\frac{\alpha \Rightarrow \beta, \quad \alpha}{\beta}$$

Warm-up Exercise

Propositional logic inference rules

- *modus ponens*: from premissis p and $p \Rightarrow q$, conclude q
- \wedge introduction: if we separately prove p and q , then that constitutes a proof of $p \wedge q$.
- \wedge elimination: from $p \wedge q$ we can conclude either of p and q separately.
- \vee introduction: from p we can conclude $p \vee q$ for any q .
- \vee elimination (also called proof by cases): if we know $p \vee q$ (the cases) and we have both $p \Rightarrow r$ and $q \Rightarrow r$ (the case-specific proofs), then we can conclude r .
- T introduction: we can conclude T from no assumptions.
- F elimination: from F we can conclude an arbitrary formula p .
- Associativity: both \wedge and \vee are associative: it doesn't matter how we parenthesize an expression like $a \wedge b \wedge c \wedge d$. (So in fact we often just leave the parentheses out in such cases. But when having \vee and \wedge together, it's a good idea to keep the parentheses.)
- Distributivity: \wedge and \vee distribute over one another; for example, $a \vee (b \wedge c)$ is equivalent to $(a \vee b) \wedge (a \vee c)$ and $a \wedge (b \vee c)$ is equivalent to $(a \wedge b) \vee (a \wedge c)$.
- Commutativity: both \wedge and \vee are commutative (symmetric in the order of their arguments), so we can re-order their arguments however we please. For example, $a \wedge b \wedge c$ is equivalent to $c \wedge b \wedge a$.

Warm-up Exercise

Use the propositional logic inference rules provided to prove:

$$(a \wedge b) \Rightarrow (b \wedge a)$$

However, you cannot use the commutativity rule.

Write your proof in two-column format, i.e., give an explicit justification for each statement based on previous statements

Warm-up Exercise

Use the propositional logic inference rules provided to prove:

$$(a \wedge b) \Rightarrow (b \wedge a)$$

However, you cannot use the commutativity rule.

Write your proof in two-column format, i.e., give an explicit justification for each statement based on previous statements

Proof by Cases

Today

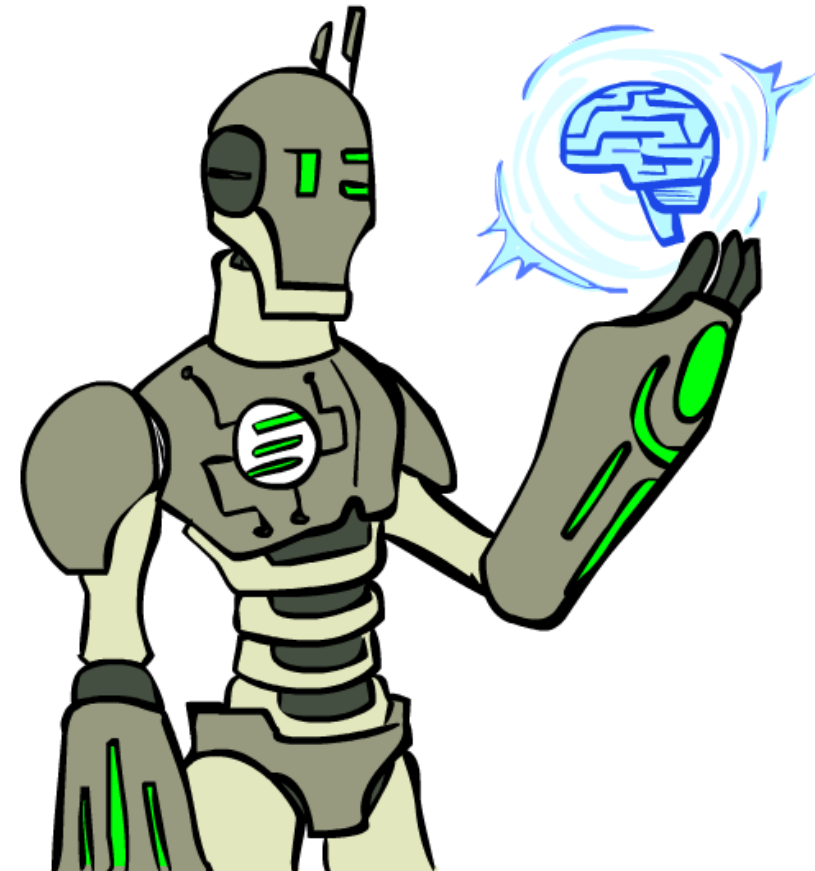
Course Info

Warm-up exercise

Propositional Logic and Proofs

ML and 606/607 Intro

More Course Info



Analysis: Perceptron

Perceptron Mistake Bound

Theorem 0.1 (Block (1962), Novikoff (1962)).

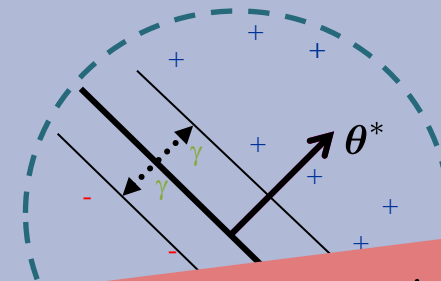
Given dataset: $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$.

Suppose:

1. Finite size inputs: $\|\mathbf{x}^{(i)}\| \leq R$
2. Linearly separable data: $\exists \boldsymbol{\theta}^*$ s.t. $\|\boldsymbol{\theta}^*\| = 1$ and $y^{(i)}(\boldsymbol{\theta}^* \cdot \mathbf{x}^{(i)}) \geq \gamma, \forall i$

Then: The number of mistakes made by the Perceptron algorithm on this dataset is

$$k \leq (R/\gamma)^2$$



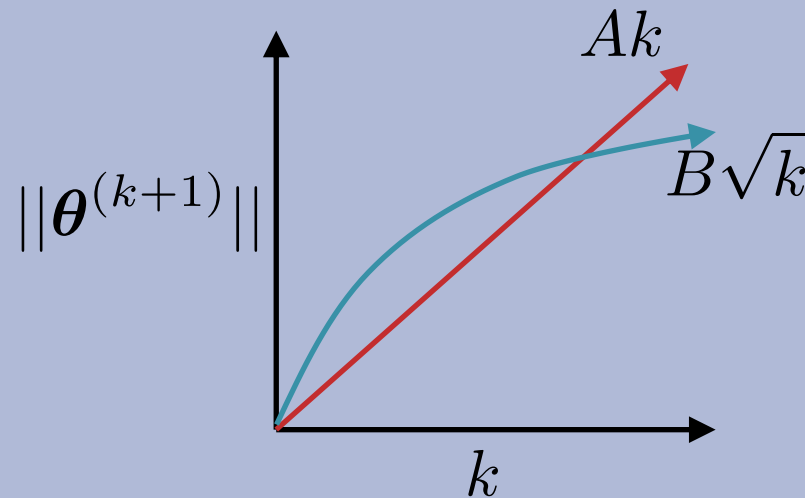
Note: This is just motivation – we'll cover the background need to understand these topics later!

Analysis: Perceptron

Proof of Perceptron Mistake Bound:

We will show that there exist constants A and B s.t.

$$Ak \leq ||\boldsymbol{\theta}^{(k+1)}|| \leq B\sqrt{k}$$



Note: This is just motivation – we'll cover the background need to understand these topics later!

Analysis: Perceptron

Theorem 0.1 (Block (1962), Novikoff (1962)).

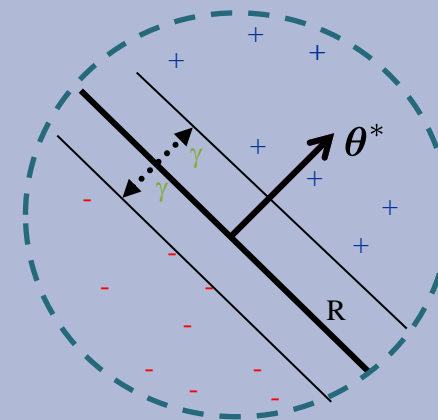
Given dataset: $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$.

Suppose:

1. Finite size inputs: $\|\mathbf{x}^{(i)}\| \leq R$
2. Linearly separable data: $\exists \boldsymbol{\theta}^*$ s.t. $\|\boldsymbol{\theta}^*\| = 1$ and $y^{(i)}(\boldsymbol{\theta}^* \cdot \mathbf{x}^{(i)}) \geq \gamma, \forall i$

Then: The number of mistakes made by the Perceptron algorithm on this dataset is

$$k \leq (R/\gamma)^2$$



Algorithm 1 Perceptron Learning Algorithm (Online)

```
1: procedure PERCEPTRON( $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots\}$ )
2:    $\boldsymbol{\theta} \leftarrow \mathbf{0}, k \leftarrow 1$                                  $\triangleright$  Initialize parameters
3:   for  $i \in \{1, 2, \dots\}$  do                                        $\triangleright$  For each example
4:     if  $y^{(i)}(\boldsymbol{\theta}^{(k)} \cdot \mathbf{x}^{(i)}) \leq 0$  then              $\triangleright$  If mistake
5:        $\boldsymbol{\theta}^{(k+1)} \leftarrow \boldsymbol{\theta}^{(k)} + y^{(i)} \mathbf{x}^{(i)}$      $\triangleright$  Update param
6:        $k \leftarrow k + 1$ 
7:   return  $\boldsymbol{\theta}$ 
```

Note: This is just motivation – we'll cover the background need to understand these topics later!

Analysis: Perceptron

Proof of Perceptron Mistake Bound:

Part 1: for some A , $Ak \leq \|\theta^{(k+1)}\|$

$$\theta^{(k+1)} \cdot \theta^* = (\theta^{(k)} + y^{(i)} \mathbf{x}^{(i)}) \theta^*$$

by Perceptron algorithm update

$$= \theta^{(k)} \cdot \theta^* + y^{(i)} (\theta^* \cdot \mathbf{x}^{(i)})$$

$$\geq \theta^{(k)} \cdot \theta^* + \gamma$$

by assumption

$$\Rightarrow \theta^{(k+1)} \cdot \theta^* \geq k\gamma$$

by induction on k since $\theta^{(1)} = \mathbf{0}$

$$\Rightarrow \|\theta^{(k+1)}\| \geq k\gamma$$

since $\|\mathbf{w}\| \times \|\mathbf{u}\| \geq \mathbf{w} \cdot \mathbf{u}$ and $\|\theta^*\| = 1$

Cauchy-Schwartz inequality

Note: This is just motivation – we'll cover the background need to understand these topics later!

Analysis: Perceptron

Proof of Perceptron Mistake Bound:

Part 2: for some B, $\|\boldsymbol{\theta}^{(k+1)}\| \leq B\sqrt{k}$

$$\|\boldsymbol{\theta}^{(k+1)}\|^2 = \|\boldsymbol{\theta}^{(k)} + y^{(i)}\mathbf{x}^{(i)}\|^2$$

by Perceptron algorithm update

$$= \|\boldsymbol{\theta}^{(k)}\|^2 + (y^{(i)})^2 \|\mathbf{x}^{(i)}\|^2 + 2y^{(i)}(\boldsymbol{\theta}^{(k)} \cdot \mathbf{x}^{(i)})$$

$$\leq \|\boldsymbol{\theta}^{(k)}\|^2 + (y^{(i)})^2 \|\mathbf{x}^{(i)}\|^2$$

since k th mistake $\Rightarrow y^{(i)}(\boldsymbol{\theta}^{(k)} \cdot \mathbf{x}^{(i)}) \leq 0$

$$= \|\boldsymbol{\theta}^{(k)}\|^2 + R^2$$

since $(y^{(i)})^2 \|\mathbf{x}^{(i)}\|^2 = \|\mathbf{x}^{(i)}\|^2 = R^2$ by assumption and $(y^{(i)})^2 = 1$

$$\Rightarrow \|\boldsymbol{\theta}^{(k+1)}\|^2 \leq kR^2$$

by induction on k since $(\boldsymbol{\theta}^{(1)})^2 = 0$

$$\Rightarrow \|\boldsymbol{\theta}^{(k+1)}\| \leq \sqrt{k}R$$

Note: This is just motivation – we'll cover the background need to understand these topics later!

Analysis: Perceptron

Proof of Perceptron Mistake Bound:

Part 3: Combining the bounds finishes the proof.

$$k\gamma \leq ||\boldsymbol{\theta}^{(k+1)}|| \leq \sqrt{k}R$$
$$\Rightarrow k \leq (R/\gamma)^2$$

The total number of mistakes
must be less than this

Note: This is just motivation – we'll cover the background need to understand these topics later!

Logic Language

Natural language?

Propositional logic

- Syntax: $P \vee (\neg Q \wedge R)$; $X_1 \Leftrightarrow (\text{Raining} \Rightarrow \text{Sunny})$
- Possible world: $\{P=\text{true}, Q=\text{true}, R=\text{false}, S=\text{true}\}$ or 1101
- Semantics: $\alpha \wedge \beta$ is true in a world iff α is true and β is true (etc.)

First-order logic

- Syntax: $\forall x \exists y P(x,y) \wedge \neg Q(\text{Joe}, f(x)) \Rightarrow f(x)=f(y)$
- Possible world: Objects o_1, o_2, o_3 ; P holds for $\langle o_1, o_2 \rangle$; Q holds for $\langle o_3 \rangle$; $f(o_1)=o_1$; $\text{Joe}=o_3$; etc.
- Semantics: $\phi(\sigma)$ is true in a world if $\sigma=o_j$ and ϕ holds for o_j ; etc.

Propositional Logic

Propositional Logic

Symbol:

- Variable that can be true or false
- We'll try to use capital letters, e.g. A , B , $P_{1,2}$
- Often include True and False

Operators:

- $\neg A$: not A
- $A \wedge B$: A and B (conjunction)
- $A \vee B$: A or B (disjunction) Note: this is not an “exclusive or”
- $A \Rightarrow B$: A implies B (implication). If A then B
- $A \Leftrightarrow B$: A if and only if B (biconditional)

Sentences

Poll 1

If we know that $A \vee B$ and $\neg B \vee C$ are true, what do we know about $A \vee C$?

- i. $A \vee C$ is guaranteed to be true
- ii. $A \vee C$ is guaranteed to be false
- iii. We don't have enough information to say anything definitive about $A \vee C$

Poll 1

If we know that $A \vee B$ and $\neg B \vee C$ are true, what do we know about $A \vee C$?

A	B	C	$A \vee B$	$\neg B \vee C$	$A \vee C$
false	false	false	false	true	false
false	false	true	false	true	true
false	true	false	true	false	false
false	true	true	true	true	true
true	false	false	true	true	true
true	false	true	true	true	true
true	true	false	true	false	true
true	true	true	true	true	true

Poll 1

If we know that $A \vee B$ and $\neg B \vee C$ are true, what do we know about $A \vee C$?

A	B	C	$A \vee B$	$\neg B \vee C$	$A \vee C$
false	false	false	false	true	false
false	false	true	false	true	true
false	true	false	true	false	false
false	true	true	true	true	true
true	false	false	true	true	true
true	false	true	true	true	true
true	true	false	true	false	true
true	true	true	true	true	true

Poll 1

If we know that $A \vee B$ and $\neg B \vee C$ are true, what do we know about $A \vee C$?

- i. $A \vee C$ is guaranteed to be true
- ii. $A \vee C$ is guaranteed to be false
- iii. We don't have enough information to say anything definitive about $A \vee C$

Poll 2

If we know that $A \vee B$ and $\neg B \vee C$ are true, what do we know about A ?

- i. A is guaranteed to be true
- ii. A is guaranteed to be false
- iii. We don't have enough information to say anything definitive about A

Poll 2

If we know that $A \vee B$ and $\neg B \vee C$ are true, what do we know about A ?

A	B	C	$A \vee B$	$\neg B \vee C$	$A \vee C$
false	false	false	false	true	false
false	false	true	false	true	true
false	true	false	true	false	false
false	true	true	true	true	true
true	false	false	true	true	true
true	false	true	true	true	true
true	true	false	true	false	true
true	true	true	true	true	true

Poll 2

If we know that $A \vee B$ and $\neg B \vee C$ are true, what do we know about A ?

- i. A is guaranteed to be true
- ii. A is guaranteed to be false
- iii. We don't have enough information to say anything definitive about A

Propositional Logic

Symbol:

- Variable that can be true or false
- We'll try to use capital letters, e.g. A , B , $P_{1,2}$
- Often include True and False

Operators:

- $\neg A$: not A
- $A \wedge B$: A and B (conjunction)
- $A \vee B$: A or B (disjunction) Note: this is not an “exclusive or”
- $A \Rightarrow B$: A implies B (implication). If A then B
- $A \Leftrightarrow B$: A if and only if B (biconditional)

Sentences

Propositional Logic Syntax

Given: a set of proposition symbols $\{X_1, X_2, \dots, X_n\}$

- (we often add **True** and **False** for convenience)

X_i is a sentence

If α is a sentence then $\neg\alpha$ is a sentence

If α and β are sentences then $\alpha \wedge \beta$ is a sentence

If α and β are sentences then $\alpha \vee \beta$ is a sentence

If α and β are sentences then $\alpha \Rightarrow \beta$ is a sentence

If α and β are sentences then $\alpha \Leftrightarrow \beta$ is a sentence

And p.s. there are no other sentences!

Notes on Operators

$\alpha \vee \beta$ is inclusive or, not exclusive

Truth Tables

$\alpha \vee \beta$ is inclusive or, not exclusive

α	β	$\alpha \wedge \beta$
F	F	F
F	T	F
T	F	F
T	T	T

α	β	$\alpha \vee \beta$
F	F	F
F	T	T
T	F	T
T	T	T

Notes on Operators

$\alpha \vee \beta$ is inclusive or, not exclusive

$\alpha \Rightarrow \beta$ is equivalent to $\neg\alpha \vee \beta$

- Says who?

Truth Tables

$\alpha \Rightarrow \beta$ is equivalent to $\neg\alpha \vee \beta$

α	β	$\alpha \Rightarrow \beta$	$\neg\alpha$	$\neg\alpha \vee \beta$
F	F	T	T	T
F	T	T	T	T
T	F	F	F	F
T	T	T	F	T

Notes on Operators

$\alpha \vee \beta$ is inclusive or, not exclusive

$\alpha \Rightarrow \beta$ is equivalent to $\neg\alpha \vee \beta$

- Says who?

$\alpha \Leftrightarrow \beta$ is equivalent to $(\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha)$

- Prove it!

Truth Tables

$\alpha \Leftrightarrow \beta$ is equivalent to $(\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha)$

α	β	$\alpha \Leftrightarrow \beta$	$\alpha \Rightarrow \beta$	$\beta \Rightarrow \alpha$	$(\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha)$
F	F	T	T	T	T
F	T	F	T	F	F
T	F	F	F	T	F
T	T	T	T	T	T

Equivalence: it's true in all models. Expressed as a logical sentence:

$$(\alpha \Leftrightarrow \beta) \Leftrightarrow [(\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha)]$$

Inference Rules

Modus Ponens

$$\frac{\alpha \Rightarrow \beta, \quad \alpha}{\beta}$$

Notation Alert!

Unit Resolution

$$\frac{a \vee b, \quad \neg b \vee c}{a \vee c}$$

General Resolution

$$\frac{a_1 \vee \cdots \vee a_m \vee b, \quad \neg b \vee c_1 \vee \cdots \vee c_n}{a_1 \vee \cdots \vee a_m \vee c_1 \vee \cdots \vee c_n}$$

Propositional Logic

Check if sentence is true in given model

In other words, does the model *satisfy* the sentence?

function PL-TRUE?(α , model) returns true or false

if α is a symbol then return $\text{Lookup}(\alpha, \text{model})$

```
if Op( $\alpha$ ) =  $\neg$  then return not(PL-TRUE?(Arg1( $\alpha$ ),model))
```

[illegible]

etc.

(Sometimes called “recursion over syntax”)

Today

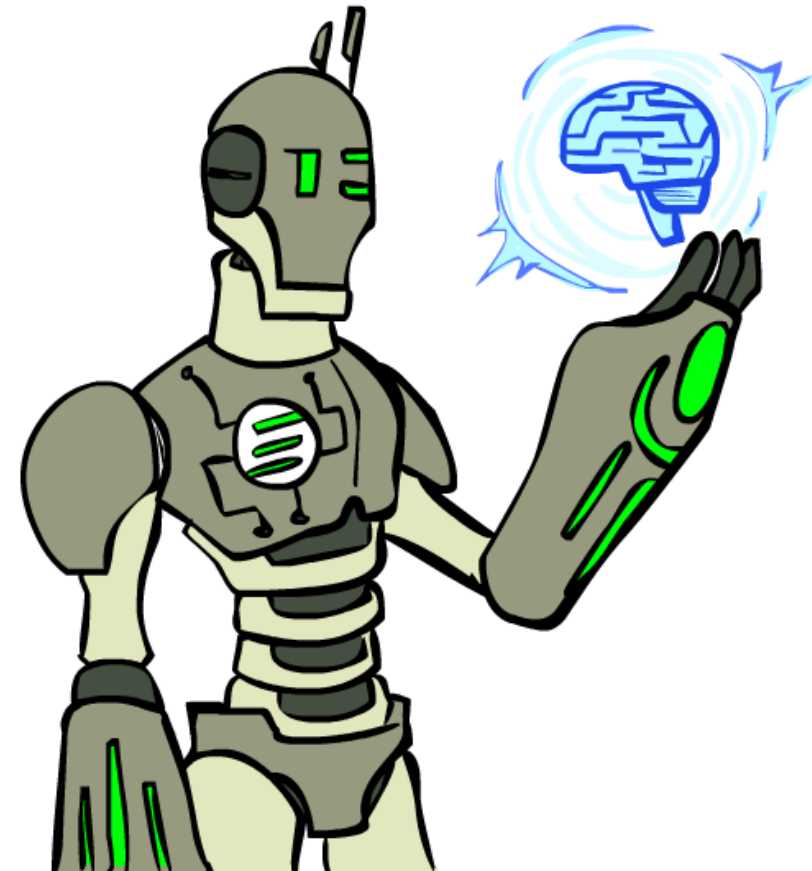
Course Info

Warm-up exercise

Propositional Logic and Proofs

ML and 606/607 Intro

More Course Info



What is ML?

Computer
Science

Machine Learning

Domain of
Interest

Optimization

Statistics

Linear Algebra

Probability

Calculus

Measure
Theory

Why Computer Science for ML?

To best understand A we need B

A

B

Why Computer Science for ML?

To best understand A we need B

A	B
Analysis of Exact Inference in Graphical Models	Computation <ul style="list-style-type: none">• Computational Complexity• Recursion; Dynamic Programming• Data Structures for ML Algorithms

Factor Graph Notation

- Variables:

$$\mathcal{X} = \{X_1, \dots, X_i, \dots, X_n\}$$

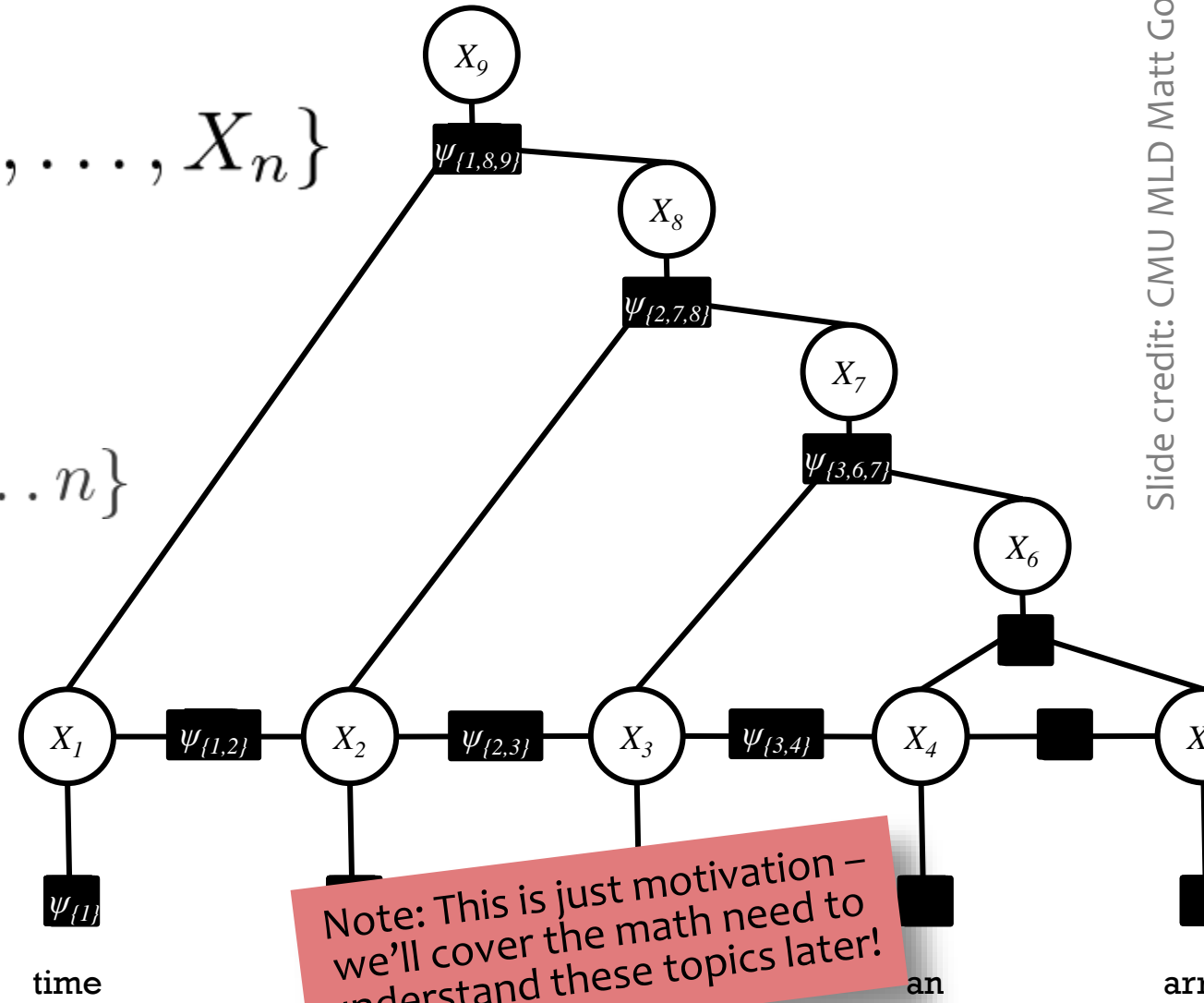
- Factors:

$$\psi_\alpha, \psi_\beta, \psi_\gamma, \dots$$

$$\text{where } \alpha, \beta, \gamma, \dots \subseteq \{1, \dots, n\}$$

Joint Distribution

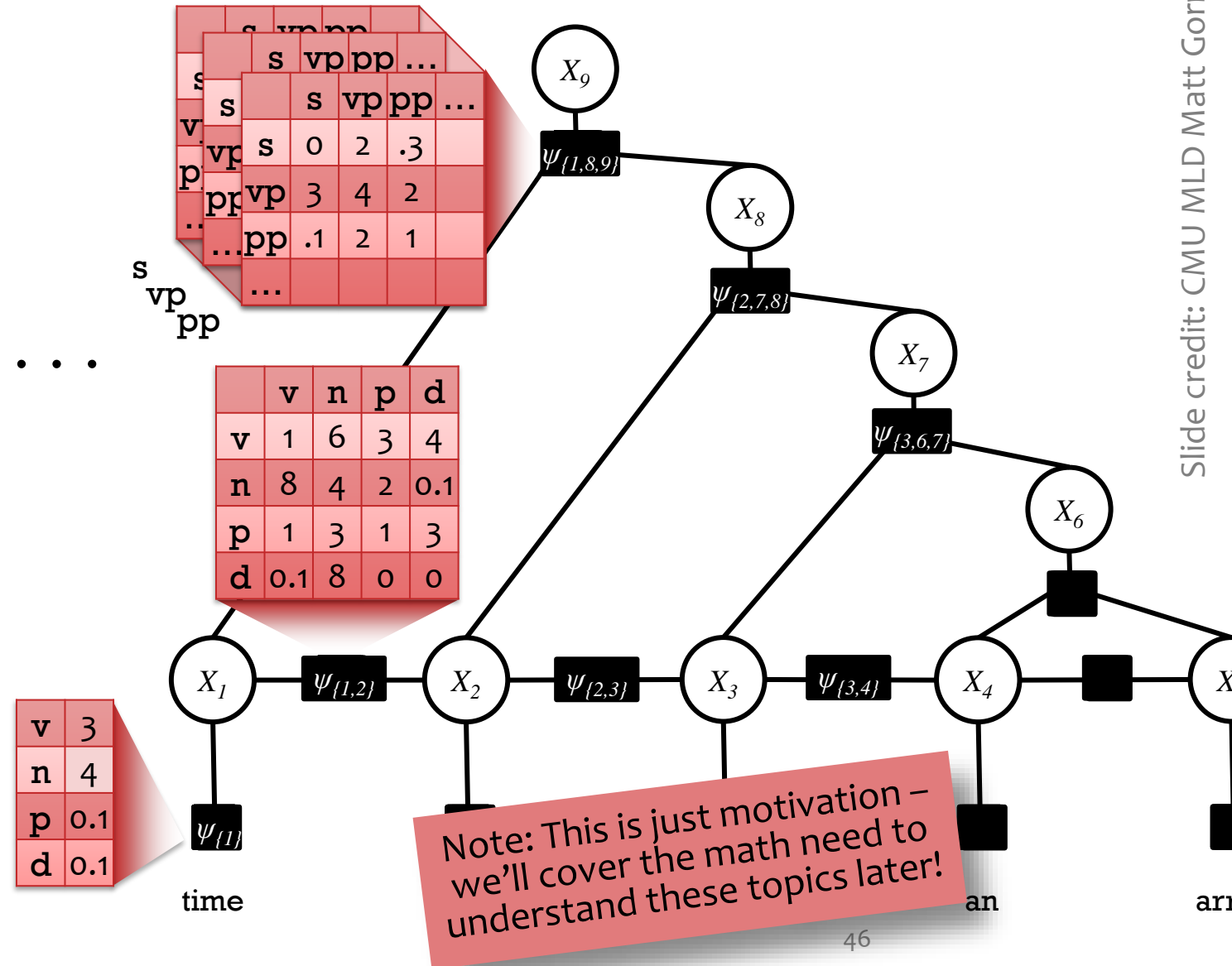
$$p(\mathbf{x}) = \frac{1}{Z} \prod_{\alpha} \psi_{\alpha}(\mathbf{x}_{\alpha})$$



Factors are Tensors

- Factors:

$\psi_\alpha, \psi_\beta, \psi_\gamma, \dots$



Inference

Given a factor graph, two common tasks ...

- Compute the most likely joint assignment,
 $\mathbf{x}^* = \operatorname{argmax}_{\mathbf{x}} p(\mathbf{X}=\mathbf{x})$
- ★ – Compute the marginal distribution of variable X_i :
 $p(X_i=x_i)$ for each value x_i

Both consider *all* joint assignments.

Both are NP-Hard in general.

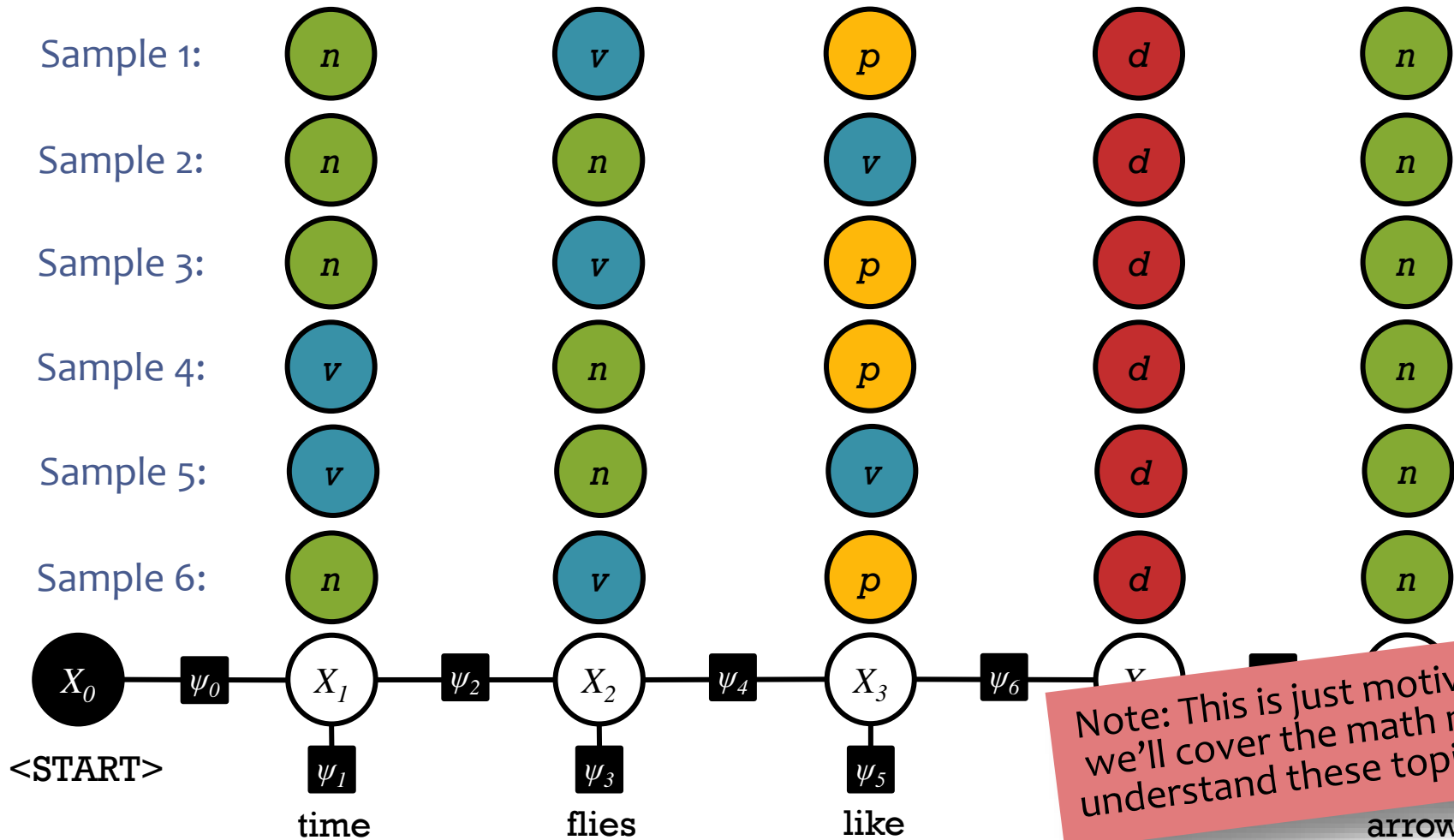
So, we turn to **approximations**.

$$p(X_i=x_i) = \text{sum of } p(\mathbf{X}=\mathbf{x}) \text{ over joint}$$

Note: This is just motivation – we'll cover the math need to understand these topics later!

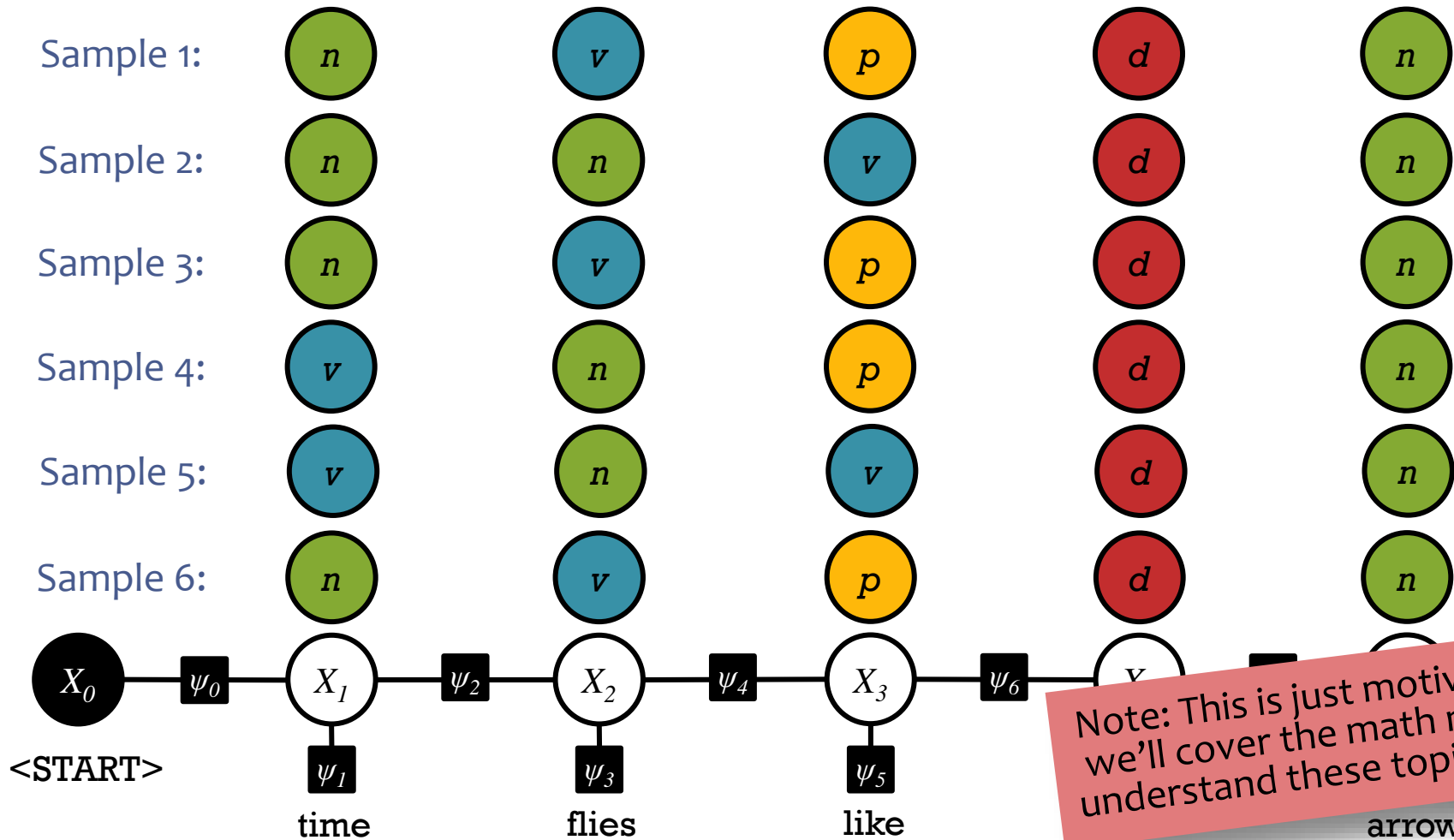
Marginals by Sampling on Factor Graph

Suppose we took many samples from the distribution over taggings: $p(\mathbf{x}) = \frac{1}{Z} \prod_{\alpha} \psi_{\alpha}(\mathbf{x}_{\alpha})$

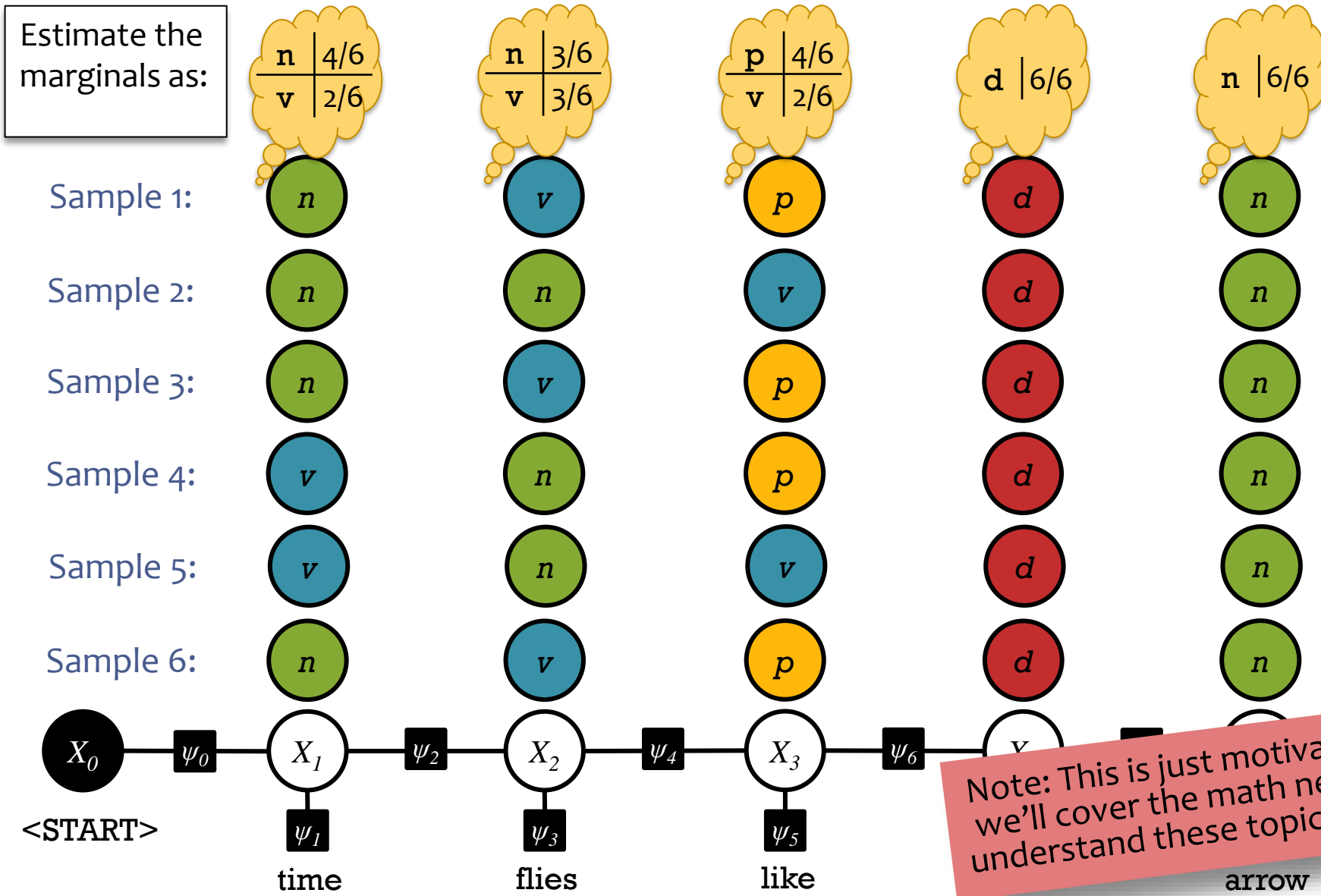


Marginals by Sampling on Factor Graph

The marginal $p(X_i = x_i)$ gives the probability that variable X_i takes value x_i in a random sample



Marginals by Sampling on Factor Graph



Why Computer Science for ML?

To best understand A we need B

A	B
Analysis of Exact Inference in Graphical Models	Computation <ul style="list-style-type: none">• Computational Complexity• Recursion; Dynamic Programming• Data Structures for ML Algorithms
Implementation Design of a Deep Learning Library	Programming & Efficiency <ul style="list-style-type: none">• Debugging for Machine Learning• Efficient Implementation / Profiling ML Algorithms

Finite Difference Method

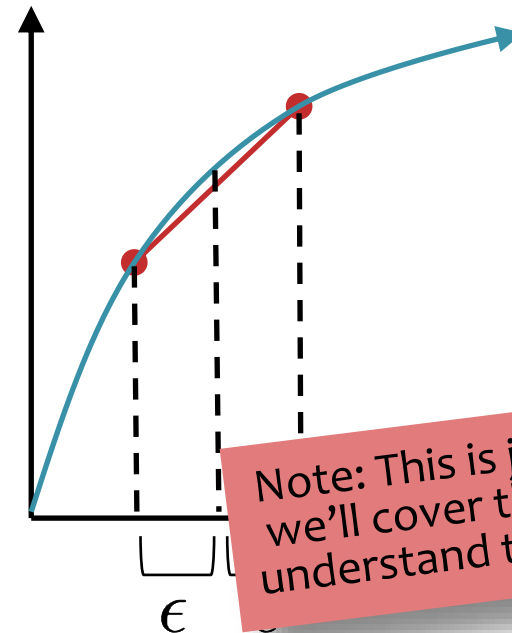
The *centered* finite difference approximation is:

$$\frac{\partial}{\partial \theta_i} J(\boldsymbol{\theta}) \approx \frac{(J(\boldsymbol{\theta} + \epsilon \cdot \mathbf{d}_i) - J(\boldsymbol{\theta} - \epsilon \cdot \mathbf{d}_i))}{2\epsilon} \quad (1)$$

where \mathbf{d}_i is a 1-hot vector consisting of all zeros except for the i th entry of \mathbf{d}_i , which has value 1.

Notes:

- Suffers from issues of floating point precision, in practice
- Typically only appropriate to use on small examples with an appropriately chosen epsilon



Note: This is just motivation – we'll cover the math need to understand these topics later!

Differentiation

Chain Rule Quiz #1:

Suppose $x = 2$ and $z = 3$, what are dy/dx and dy/dz for the function below?

$$y = \exp(xz) + \frac{xz}{\log(x)} + \frac{\sin(\log(x))}{\exp(xz)}$$

Finite Difference Solution:

```
from math import *

# Define function
def f(x, z):
    return exp(x*z) + x*z/log(x) + sin(log(x)) / exp(x*z)

# Inputs
x = 2; z = 3; e = 1e-8

# Finite difference check
dydx = (f(x+e, z) - f(x-e, z)) / (2*e)
dydz = (f(x, z+e) - f(x, z-e)) / (2*e)
print "dydx =", dydx
print "dydz =", dydz
```

Note: This is just motivation – we'll cover the math need to understand these topics later!

Automatic Differentiation – Reverse Mode (aka. Backpropagation)

Forward Computation

1. Write an **algorithm** for evaluating the function $y = f(\mathbf{x})$. The algorithm defines a **directed acyclic graph**, where each variable is a node (i.e. the “**computation graph**”)
2. Visit each node in **topological order**.
For variable u_i with inputs v_1, \dots, v_N
 - a. Compute $u_i = g_i(v_1, \dots, v_N)$
 - b. Store the result at the node

Backward Computation

1. **Initialize** all partial derivatives dy/du_j to 0 and $dy/dy = 1$.
2. Visit each node in **reverse topological order**.
For variable $u_i = g_i(v_1, \dots, v_N)$
 - a. We already know dy/du_i
 - b. Increment dy/dv_j by $(dy/du_i)(du_i/dv_j)$
(Choice of algorithm ensures computing (du_i/dv_j) is efficient)

Return partial derivatives dy/du_i for all variables

Note: This is just motivation – we’ll cover the math need to understand these topics later!

Why Computer Science for ML?

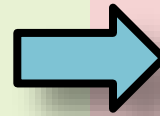
To best understand A we need B

A	B
Analysis of Exact Inference in Graphical Models	Computation <ul style="list-style-type: none">• Computational Complexity• Recursion; Dynamic Programming• Data Structures for ML Algorithms
Implementation Design of a Deep Learning Library	Programming & Efficiency <ul style="list-style-type: none">• Debugging for Machine Learning• Efficient Implementation / Profiling ML Algorithms
Optimization for Support Vector Machines (SVMs)	Optimization <ul style="list-style-type: none">• Unconstrained Optimization• Preconditioning• Constrained Optimization

Support Vector Machines (SVMs)

Hard-margin SVM (Primal)

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1, \quad \forall i = 1, \dots, N \end{aligned}$$



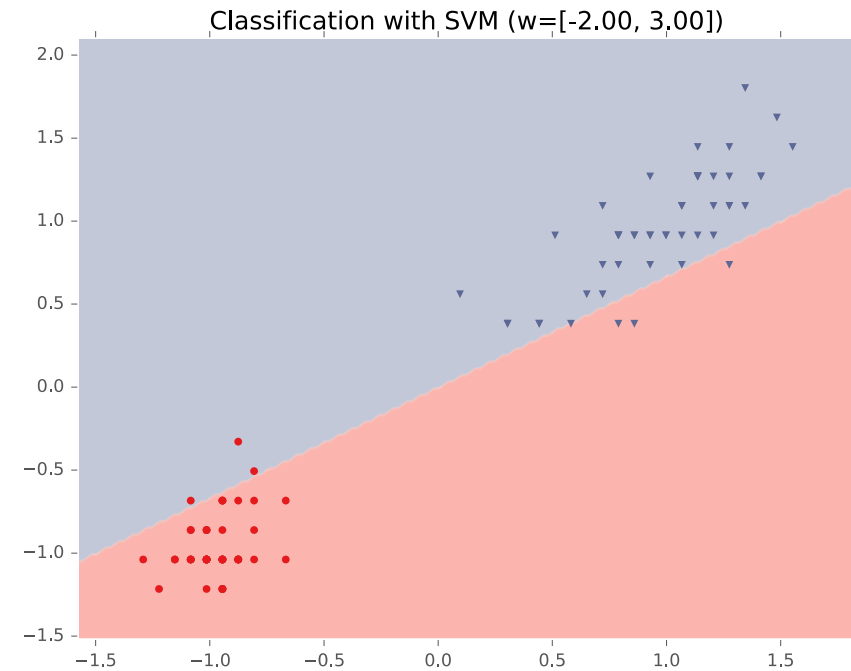
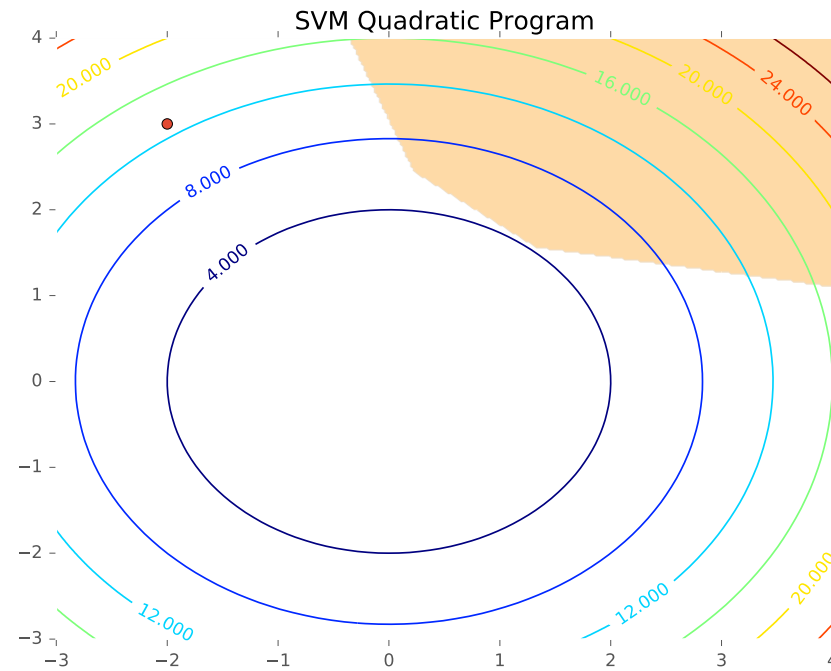
Hard-margin SVM (Lagrangian Dual)

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y^{(i)} y^{(j)} \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} \\ \text{s.t.} \quad & \alpha_i \geq 0, \quad \forall i = 1, \dots, N \\ & \sum_{i=1}^N \alpha_i y^{(i)} = 0 \end{aligned}$$

- Instead of minimizing the primal, we can maximize the dual problem
- For the SVM, these two problems give the same answer (i.e. the minimum of one is the maximum of the other)
- **Definition: support vectors** are those which $\alpha^{(i)} \neq 0$

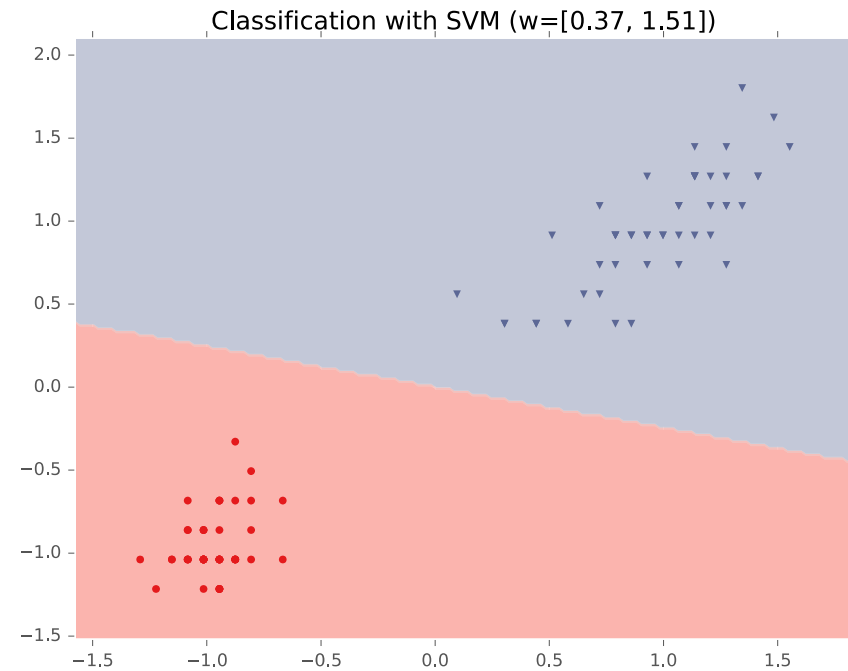
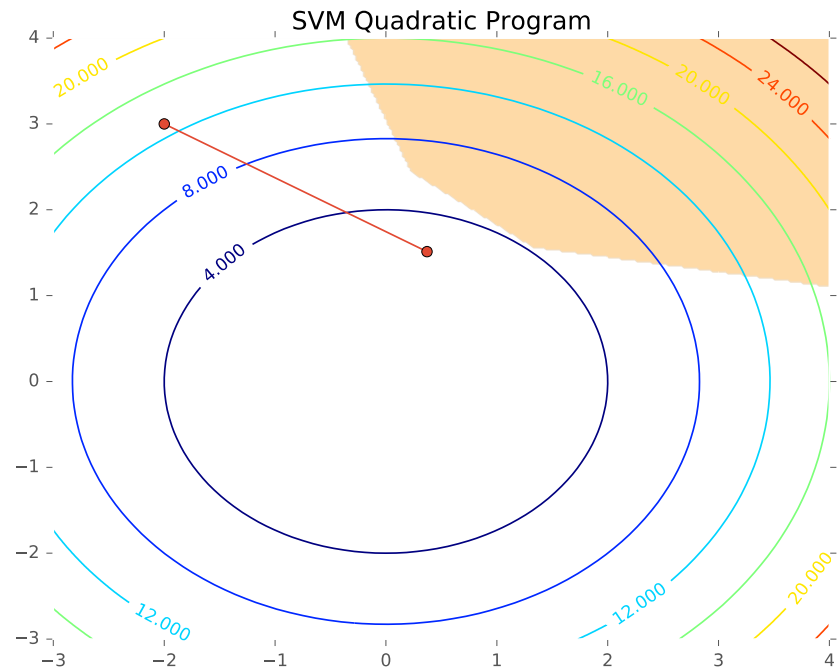
Note: This is just motivation – we'll cover the math need to understand these topics later!

SVM QP



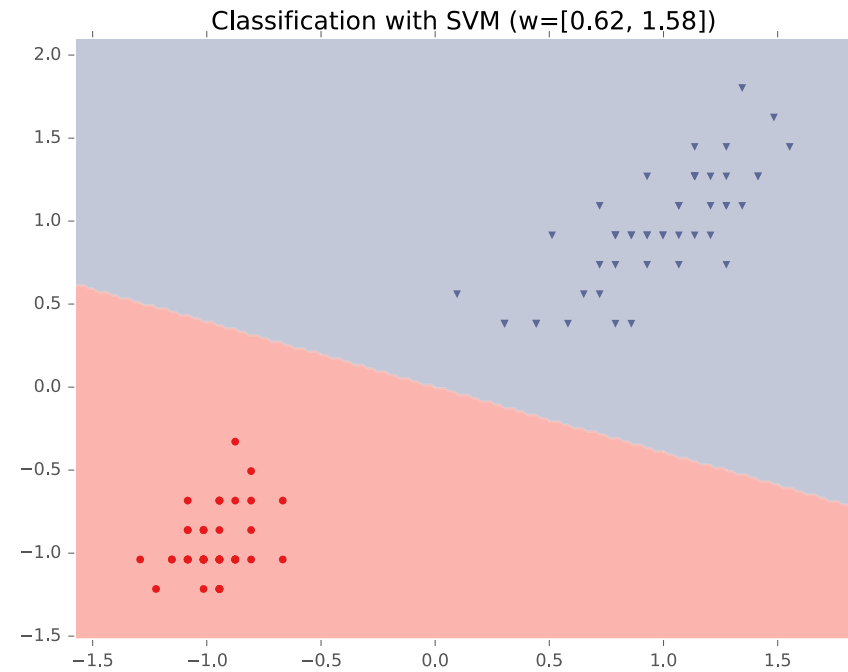
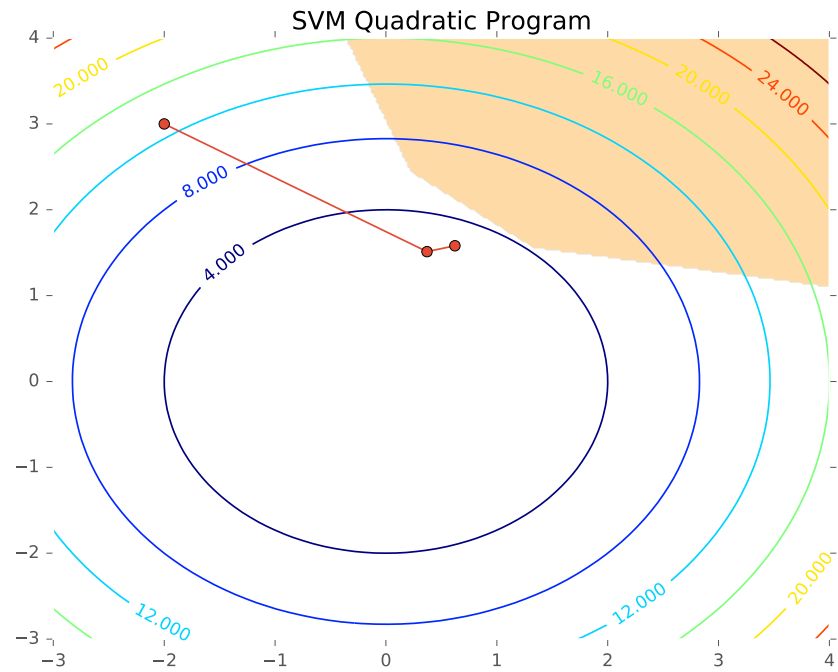
Note: This is just motivation – we'll cover the math need to understand these topics later!

SVM QP



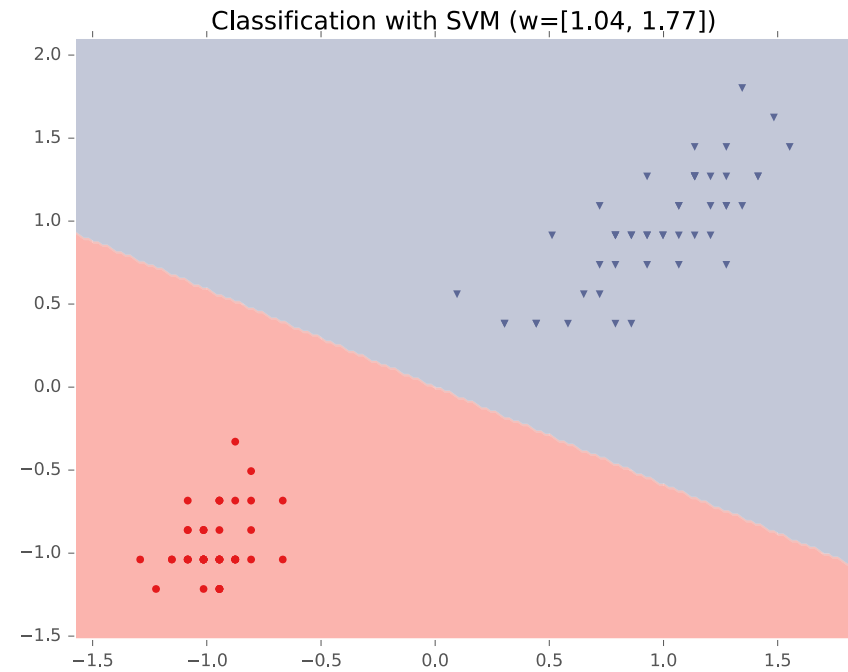
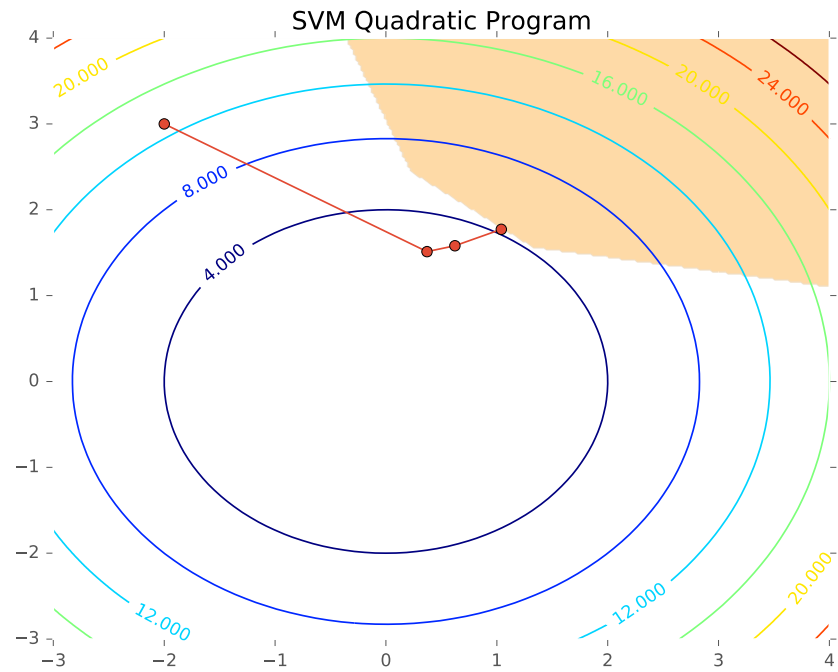
Note: This is just motivation – we'll cover the math need to understand these topics later!

SVM QP



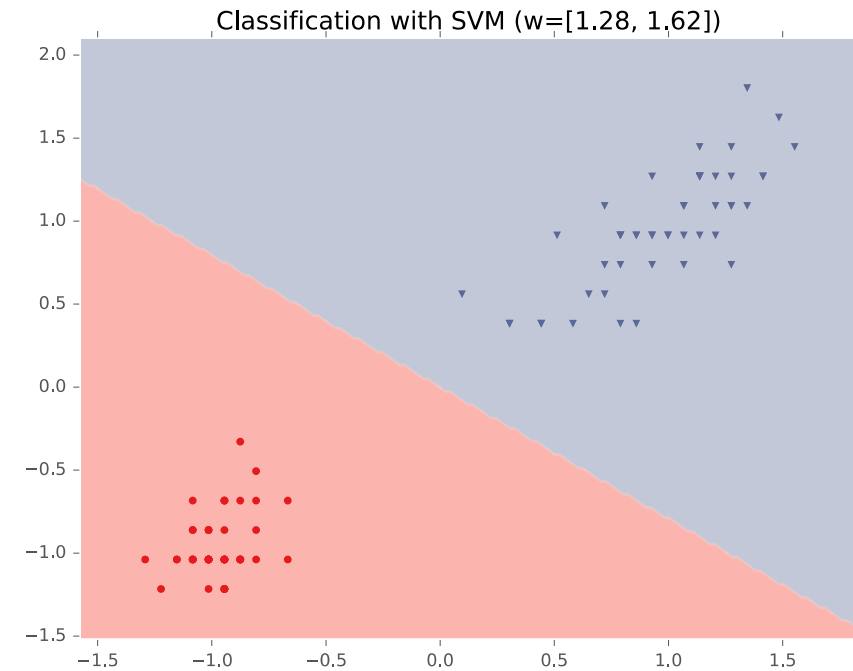
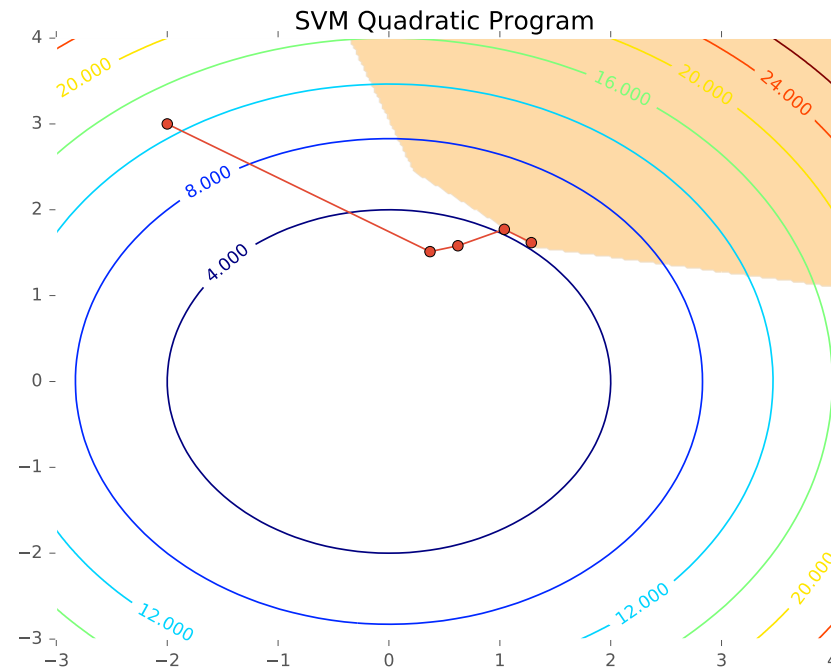
Note: This is just motivation – we'll cover the math need to understand these topics later!

SVM QP



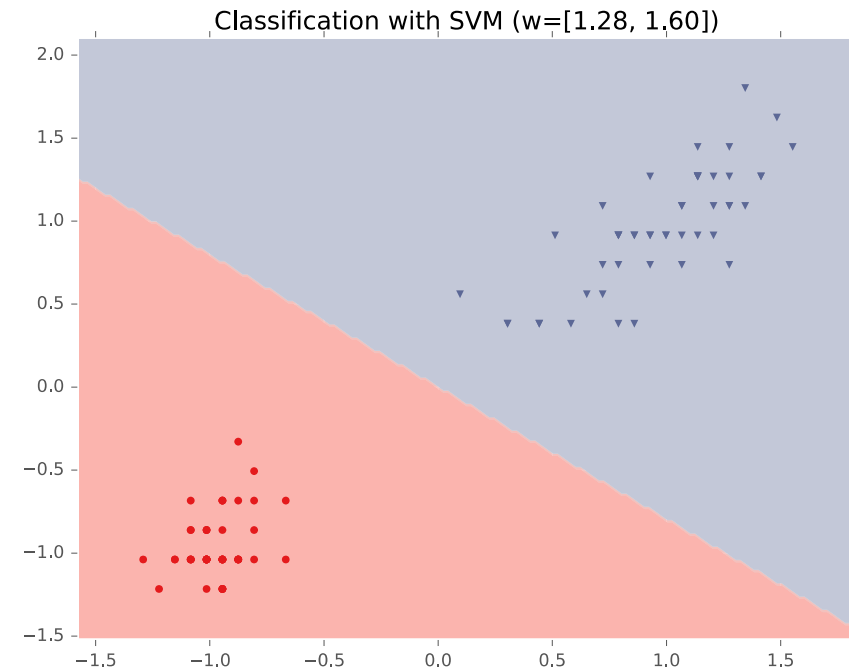
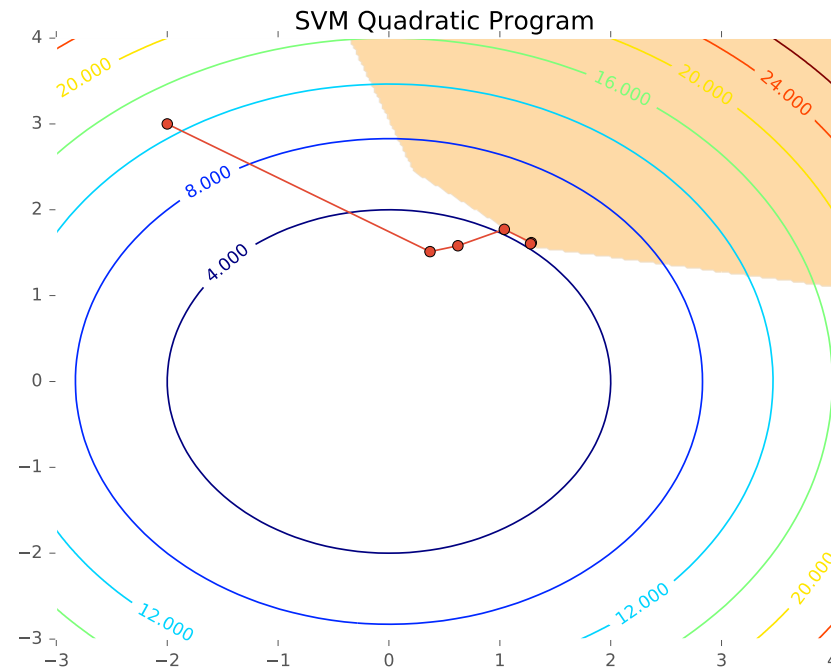
Note: This is just motivation – we'll cover the math need to understand these topics later!

SVM QP



Note: This is just motivation – we'll cover the math need to understand these topics later!

SVM QP



Note: This is just motivation –
we'll cover the math need to
understand these topics later!

Why Computer Science for ML?

To best understand A we need B

A	B
Analysis of Exact Inference in Graphical Models	Computation <ul style="list-style-type: none"> • Computational Complexity • Recursion; Dynamic Programming • Data Structures for ML Algorithms
Implementation Design of a Deep Learning Library	Programming & Efficiency <ul style="list-style-type: none"> • Debugging for Machine Learning • Efficient Implementation / Profiling ML Algorithms
Optimization for Support Vector Machines (SVMs)	Optimization <ul style="list-style-type: none"> • Unconstrained Optimization • Preconditioning • Constrained Optimization

The core content for this course is the **computer science** (Column B), but you will apply what you learn to **real problems in machine learning** (Column A)

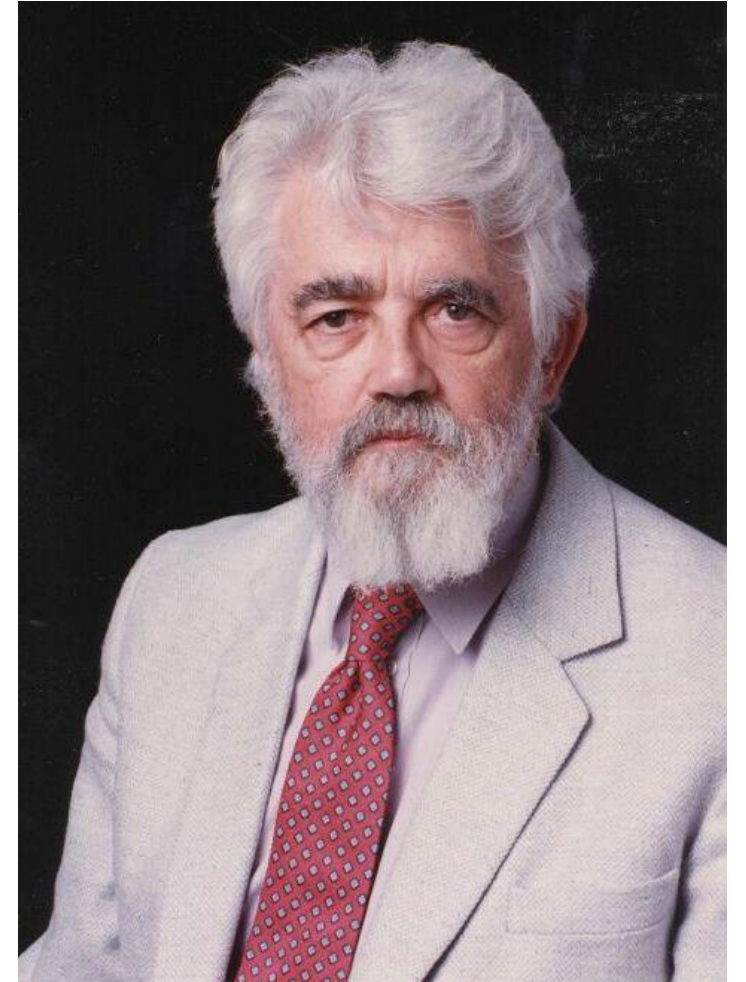
AI Definition by John McCarthy

What is artificial intelligence

- It is the science and engineering of making intelligent machines, especially intelligent computer programs

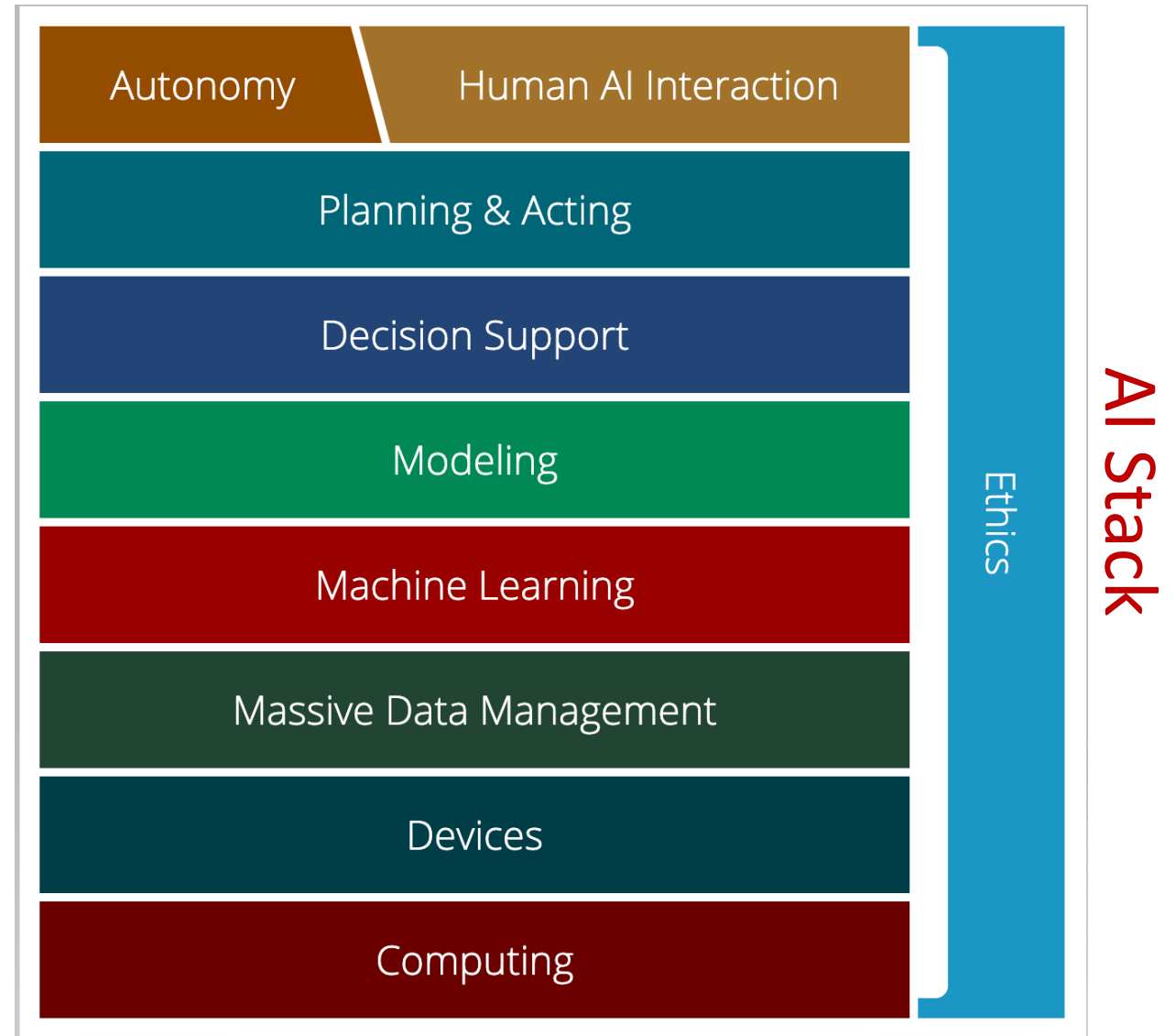
What is intelligence

- Intelligence is the computational part of the ability to achieve goals in the world



AI Stack for CMU AI

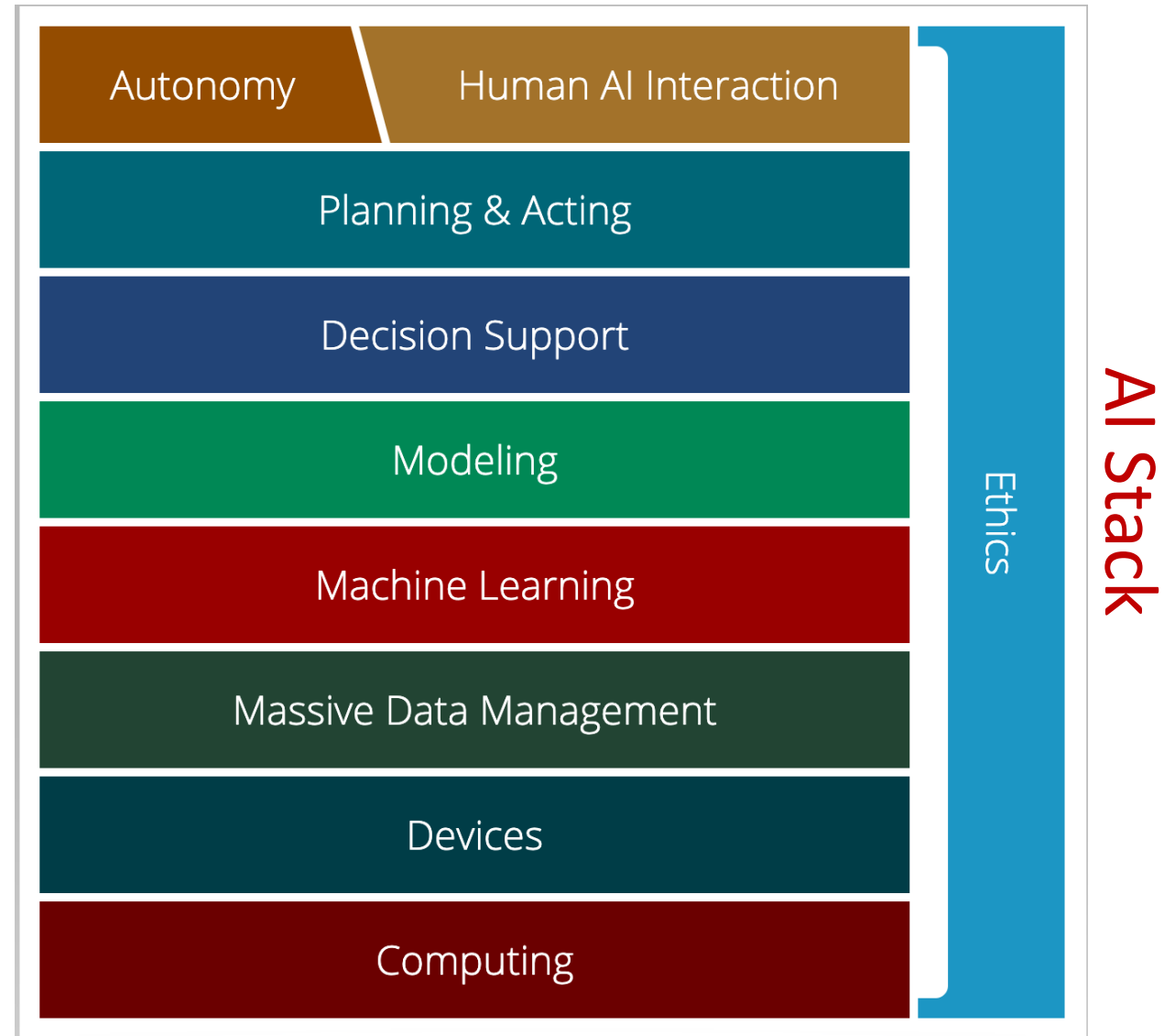
“AI must understand the human needs and it must make smart design decisions based on that understanding”



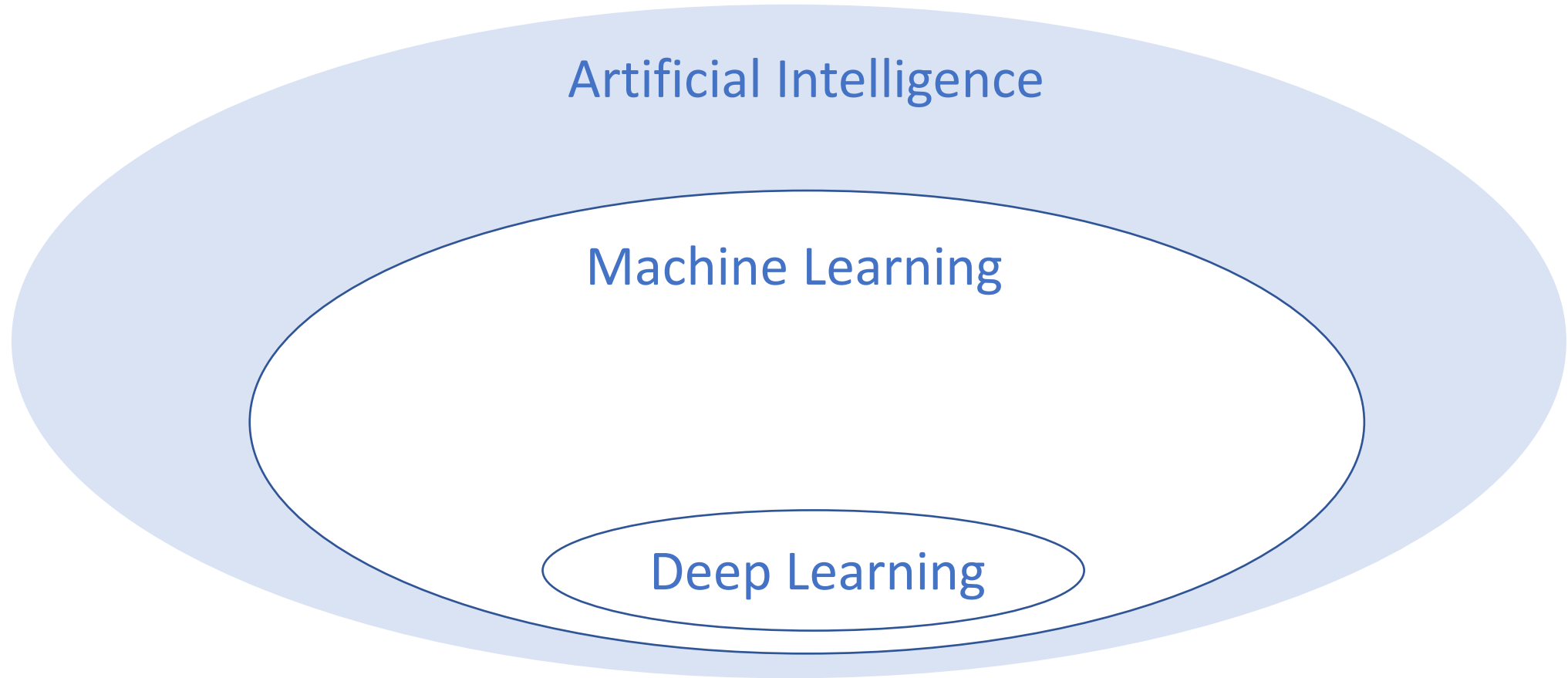
AI Stack for CMU AI

“Machine learning focuses on creating programs that learn from experience.”

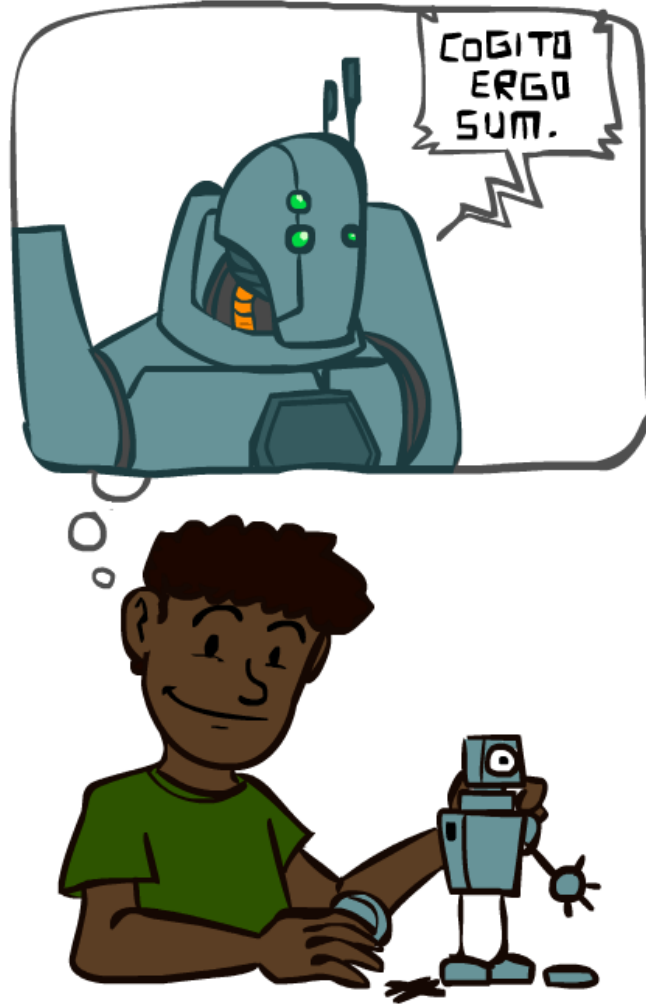
“It advances computing through exposure to new scenarios, testing and adaptation, while using pattern- and trend-detection to help the computer make better decisions in similar, subsequent situations.”



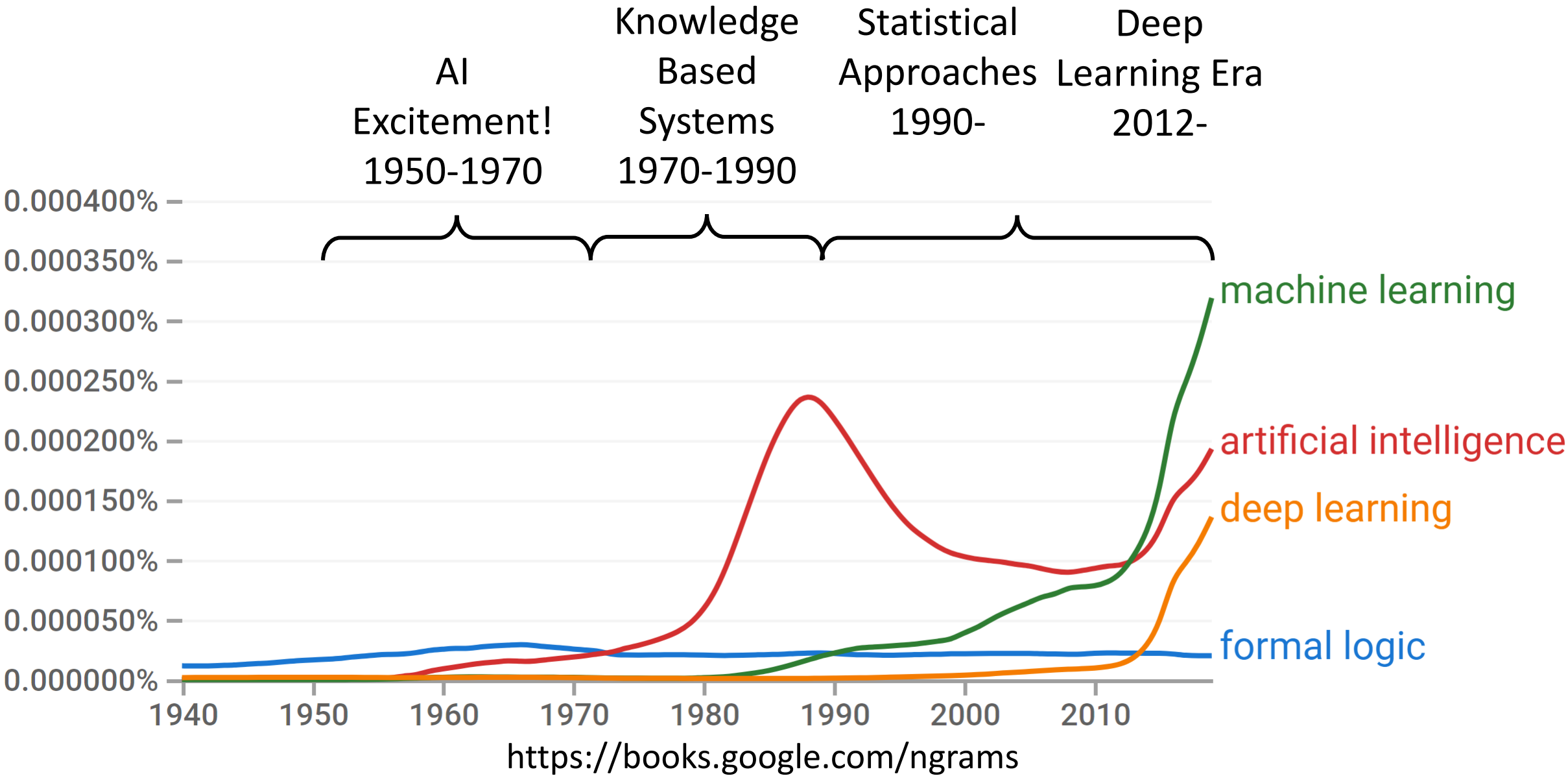
Artificial Intelligence vs Machine Learning?



A Brief History of AI



A Brief History of AI



A Brief History of AI

1940-1950: Early days

- 1943: McCulloch & Pitts: Boolean circuit model of brain
- 1950: Turing's "Computing Machinery and Intelligence"

1950—70: Excitement: Look, Ma, no hands!

- 1950s: Early AI programs, including Samuel's checkers program, Newell & Simon's Logic Theorist, Gelernter's Geometry Engine
- 1956: Dartmouth meeting: "Artificial Intelligence" adopted

1970—90: Knowledge-based approaches

- 1969—79: Early development of knowledge-based systems
- 1980—88: Expert systems industry booms
- 1988—93: Expert systems industry busts: "AI Winter"

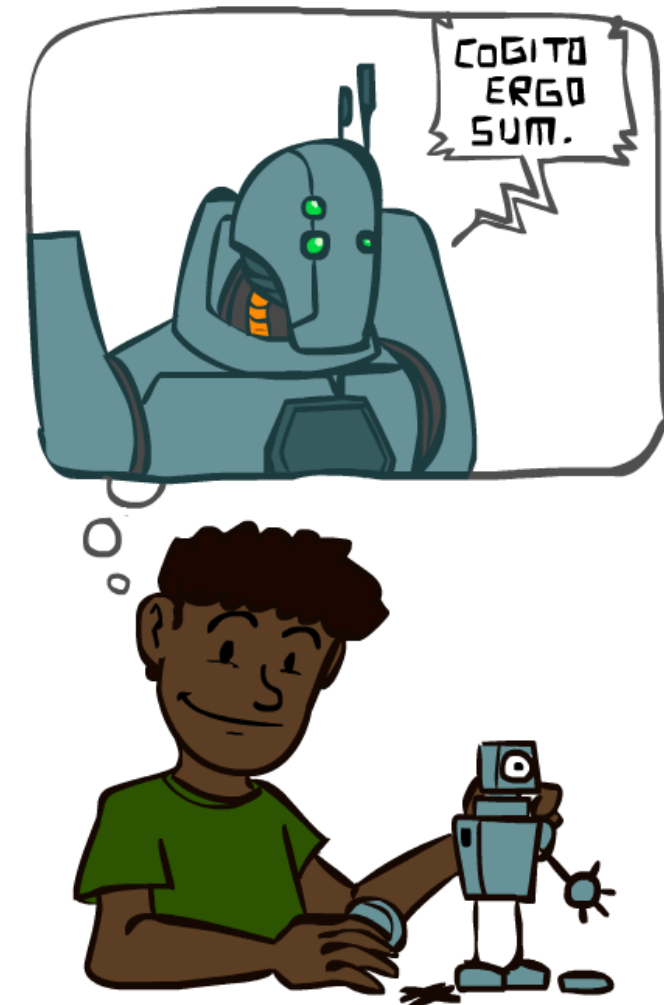
1990—: Statistical approaches

- Resurgence of probability, focus on uncertainty
- General increase in technical depth
- Agents and learning systems... "AI Spring"?

2012—: Deep learning

- 2012: ImageNet & AlexNet


Images: ai.berkeley.edu



ML Applications?


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>Source: https://www.stonemple.com/great-knowledge-box-showdown/#VoiceStudyResults</p>


Robotics

2. Learning to drive an autonomous vehicle

THEN	NOW
<p>"...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..."</p> <p>(Mitchell, 1997)</p>	 <p>waymo.com</p>


Games / Reasoning

3. Learning to beat the masters at board games

THEN	NOW
<p>"...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..."</p> <p>(Mitchell, 1997)</p>	

Computer Vision

4. Learning to recognize images

THEN	NOW
<p>"...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..."</p> <p>(LeCun et al., 1995)</p>	 <p>Images from https://blog.openai.com/generative-models/</p>

Learning Theory

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

Sample Complexity Results

Definition 5.1. 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 [K]:** $N = \frac{1}{\epsilon} \log \frac{1}{\delta} \log \frac{1}{\epsilon}$ (related to the number of hypotheses in the hypothesis space)
- Infinite [K]:** $N = \frac{1}{\epsilon^2} \log \frac{1}{\delta} \log \frac{1}{\epsilon}$ (related to the number of hypotheses in the hypothesis space)
- Agnostic:** $N = \frac{1}{\epsilon^2} \log \frac{1}{\delta} \log \frac{1}{\epsilon}$ (related to the number of hypotheses in the hypothesis space)
- Realizable:** $N = \frac{1}{\epsilon} \log \frac{1}{\delta} \log \frac{1}{\epsilon}$ (related to the number of hypotheses in the hypothesis space)

Handwritten Notes:

- Two Types of Error:**
 - Empirical Error:** $R(h) = \frac{1}{N} \sum_{i=1}^N \ell(h(x_i), y_i)$
 - Expected Error:** $R(h) = \mathbb{E}[\ell(h(x), y)]$
- Can we bound $R(h)$ in terms of $R(h)$?**
- PK: $R(h)$ is the error of h on the training set.**
- PK: $R(h)$ is the error of h on the test set.**
- PK: $R(h)$ is the error of h on the training set.**
- PK: $R(h)$ is the error of h on the test set.**
- PK: $R(h)$ is the error of h on the training set.**
- PK: $R(h)$ is the error of h on the test set.**

Speech Recognition

1. Learning to recognize spoken words

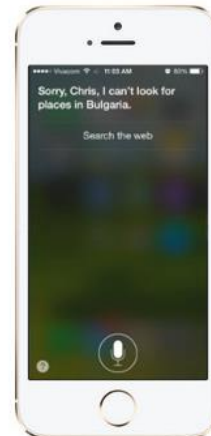
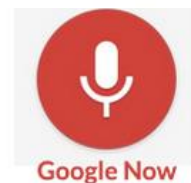
THEN

“...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)



NOW



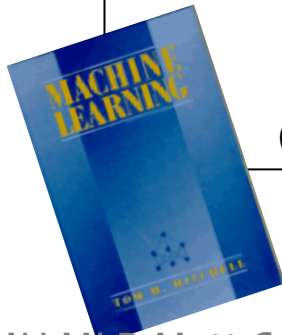
Source: <https://www.stonetemple.com/great-knowledge-box-showdown/#VoiceStudyResults>

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



<https://www.geek.com/wp-content/uploads/2016/03/uber.jpg>

Games / Reasoning

3. Learning to beat the masters at board games

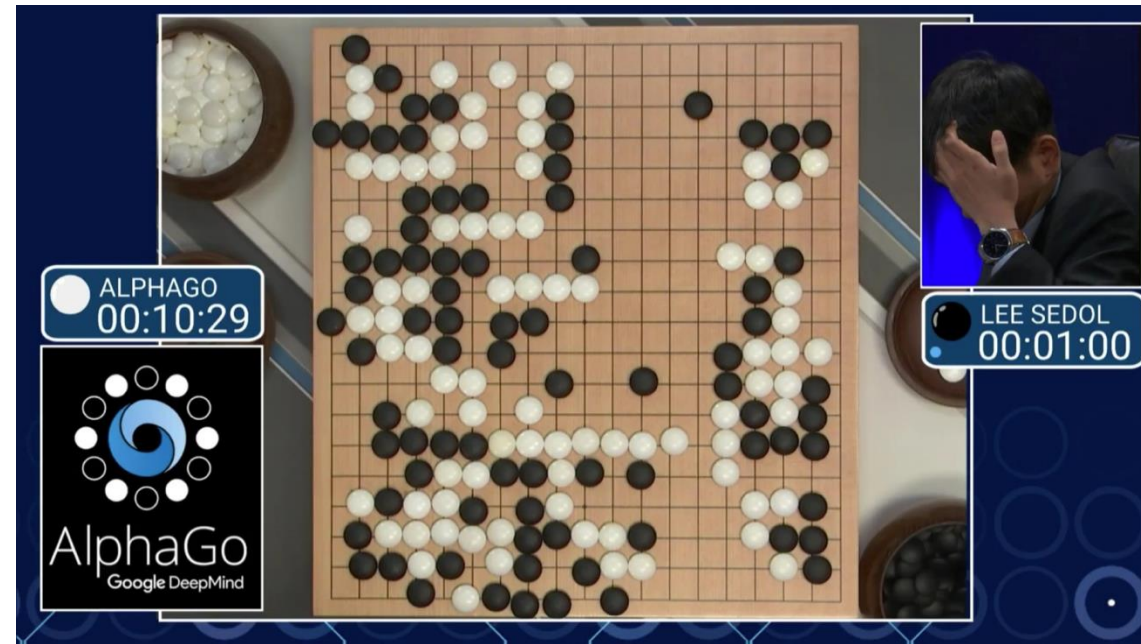
THEN

“...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

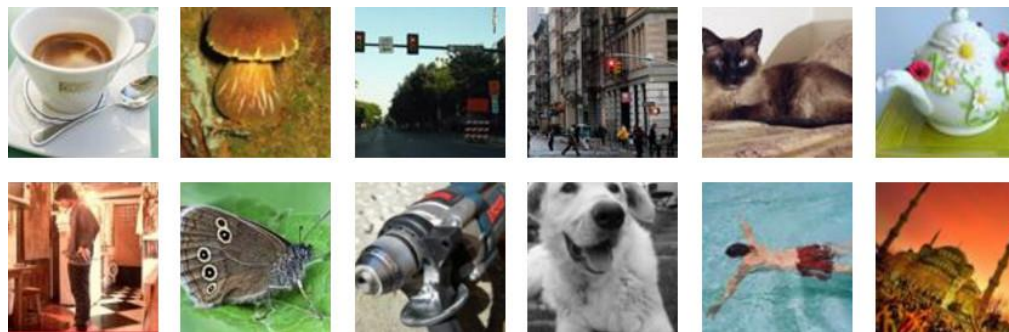
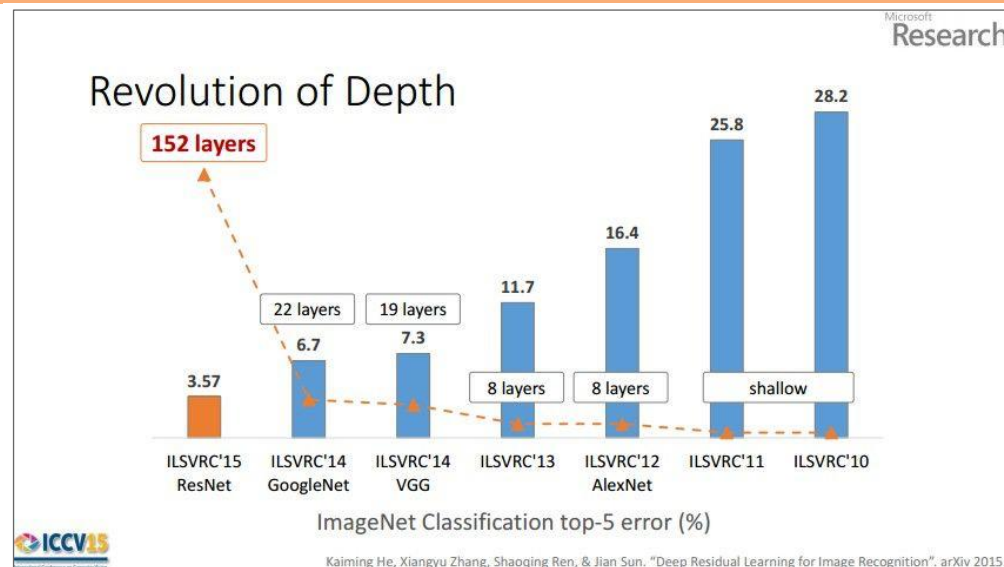
“...The recognizer is a convolution network that can be spatially replicated. From the network output, a hidden Markov model produces word scores. The entire system is globally trained to minimize word-level errors....”

Figure 2: Convolutional neural network character recognizer. This architecture is robust to local translations and distortions, with subsampling, shared weights, and local receptive fields.

number of subsampling layers and the sizes of the kernels are chosen, the sizes of all the layers, including the input, are determined unambiguously. The only architectural parameters that remain to be selected are the number of feature maps in each layer, and the information as to what feature map is connected to what other feature map. In our case, the subsampling rates were chosen as small as possible (2×2), and the kernels as small as possible in the first layer (3×3) to limit the total number of connections. Kernel sizes in the upper layers are chosen to be as small as

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

② Train Error (aka. empirical risk)
 $\hat{R}(h) = \mathbb{P}_{x \sim S}(c^*(x) \neq h(x))$ ← $S = \{x^{(1)}, \dots, x^{(N)}\}$
 $= \frac{1}{N} \sum_{i=1}^N \mathbb{I}(c^*(x^{(i)}) \neq h(x^{(i)}))$ ← known, computable
 $= \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y^{(i)} \neq h(x^{(i)}))$

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?

10-606 and 10-607

- Mini Courses
 - 10-606
 - 10-607
- Intro ML Courses
 - 10-315
 - 10-301/601
 - 10-701
 - 10-715
- Prerequisites

Today

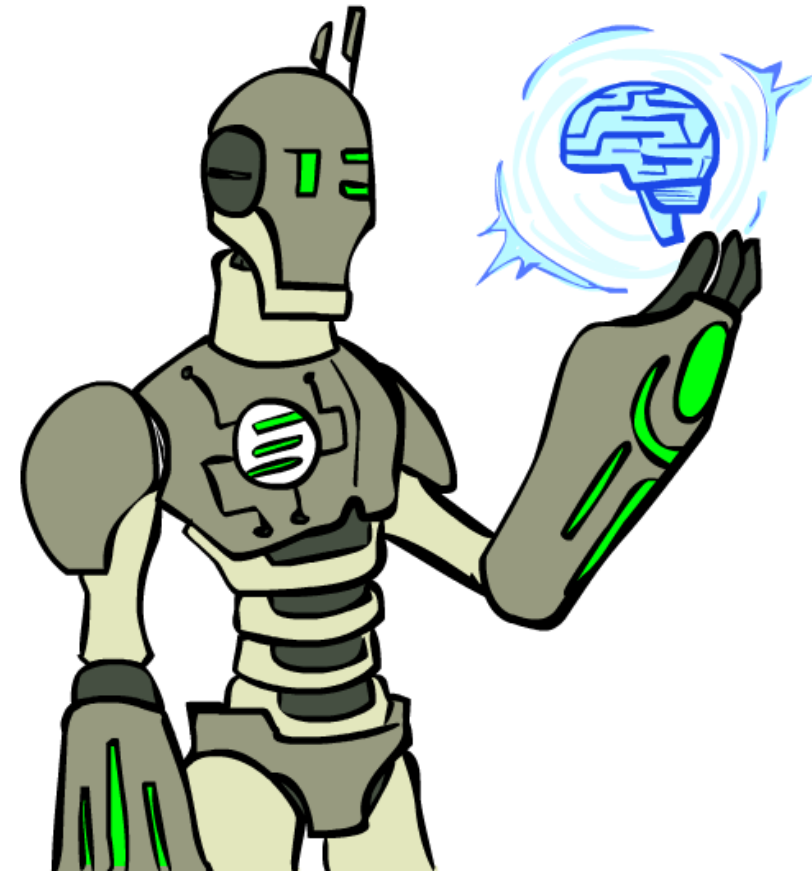
Course Info

Warm-up exercise

Propositional Logic and Proofs

ML and 606/607 Intro

More Course Info



Course Information

Website: <https://www.cs.cmu.edu/~10607>

Canvas: canvas.cmu.edu



Gradescope: gradescope.com



Communication:

piazza.com



E-mail (if piazza doesn't work):

pvirtue@andrew.cmu.edu

Course Information

Lectures

- Lectures are recorded
 - Shared with our course and ML course staff only
- Participation point earned by answering Piazza polls in lecture
- Quizzes will in lecture, announced two days ahead of time
- Slides will be posted

Recitations

- Recommended attendance
- No plans to record at this point
- No participation points in recitation
- Recitation materials are in-scope for quizzes and exams

Course Information

Office Hours

- OH calendar on course website
- OH-by-appointment requests are certainly welcome

Mental Health