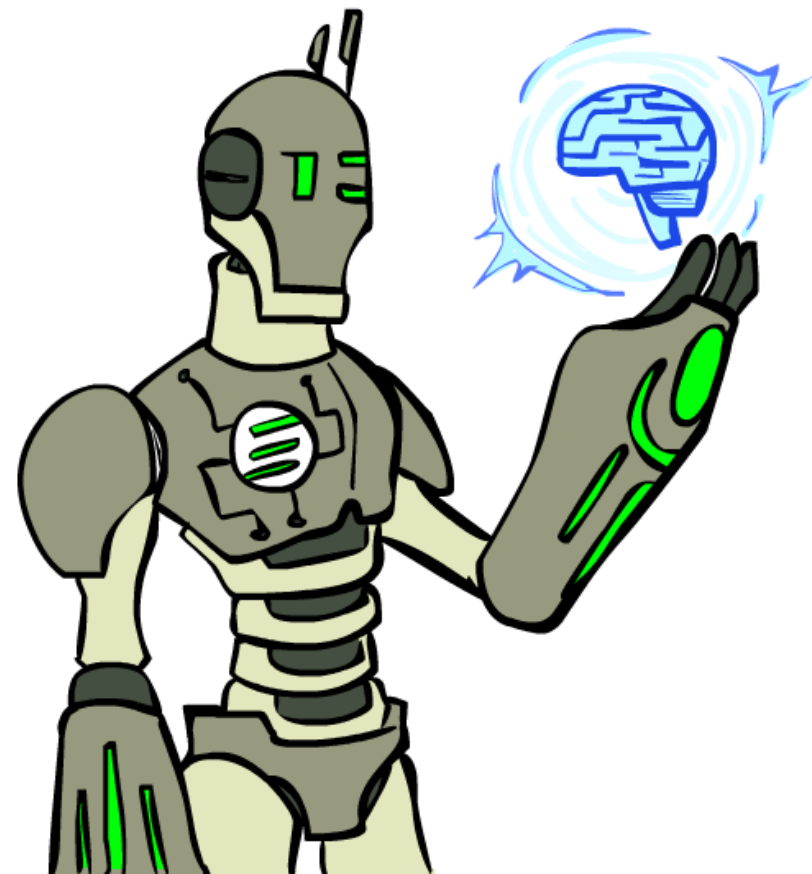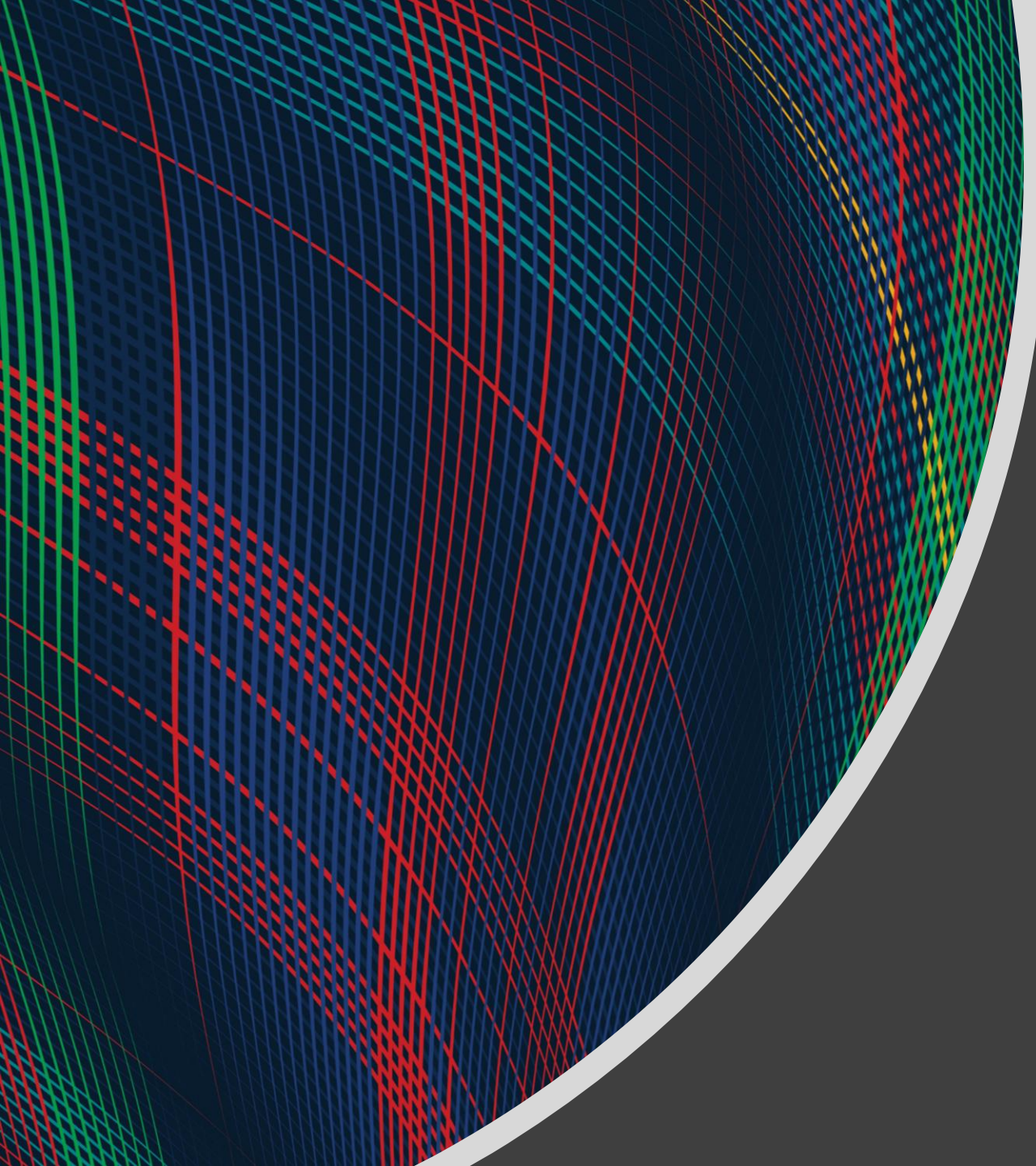# As you walk in

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

Images: ai.berkeley.edu

# 10-607
# Computational Foundations for Machine Learning
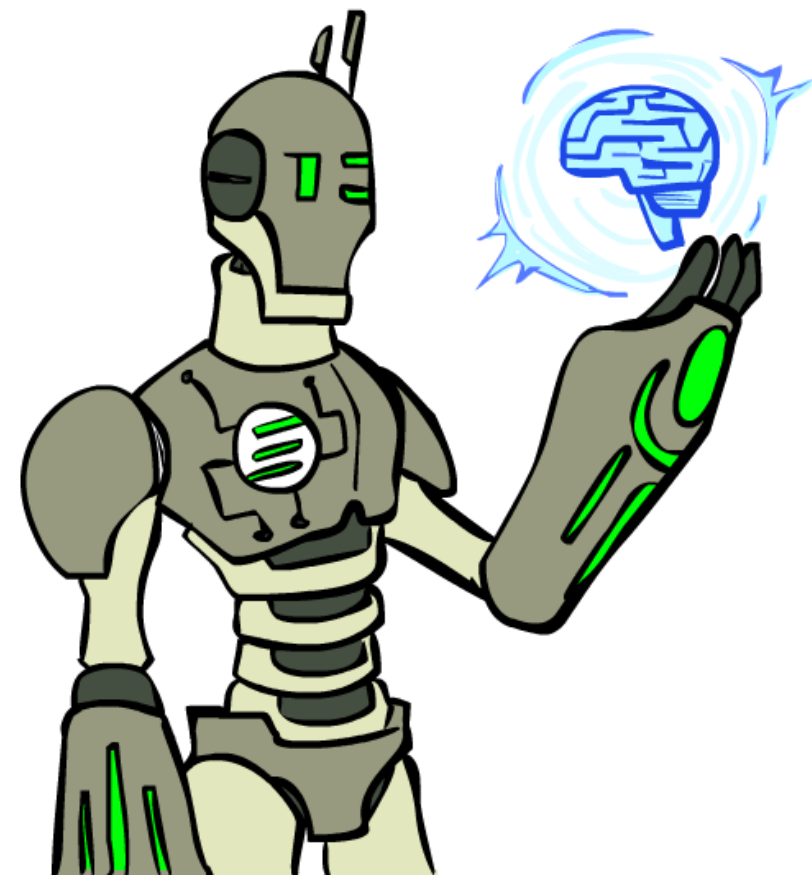
Instructor: Pat Virtue

# 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:

what's bash?

you don't know what bash is?!
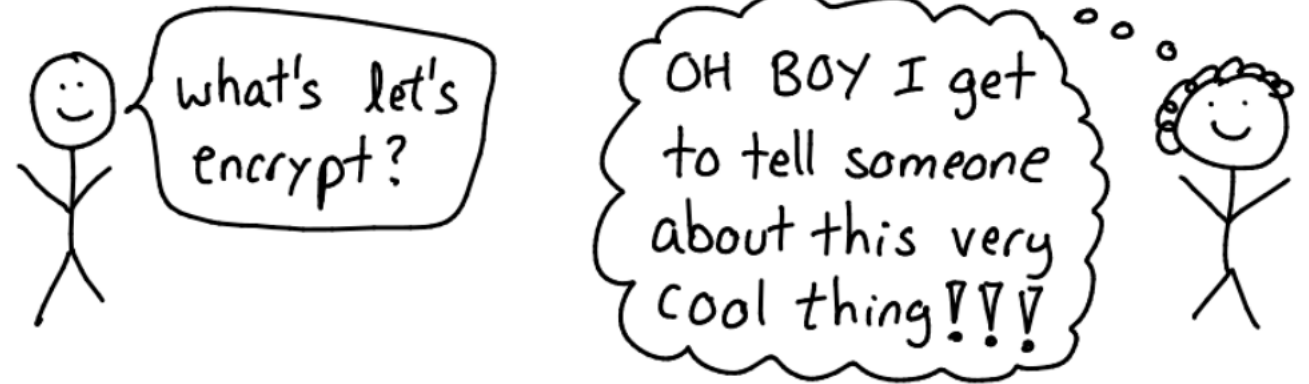
oooo oh I guess I was supposed to know that?

# 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:

what's let's encrypt?

OH BOY I get to tell someone about this very cool thing!!!

And it gets easier with practice!

# Two-column Proof

A: Socrates is human

A ⟹ B: human the mortal

Give an explicit justification for each statement based on previous statements

Prove Socrates is mortal

Justification

1   $A$

Assumption

2   $A \Rightarrow B$

Assumption

(1) & (2)

3   $B$

Modus Ponens   Notation Alert!

Modus Ponens

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

# Warm-up Exercise

## Propositional logic inference rules

- *modus ponens*: from premesis $p$ 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 $p \vee r$ and $p \vee r$ (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

| | | |
|---|---|---|
| 1 | $a \wedge b$ | Assumption |
| 2 | $a$ | $\wedge$ Elim (1) |
| 3 | $b$ | $\wedge$ Elim (1) |
| 4 | $b \wedge a$ | $\wedge$ Intro (2) (3) |
| 5 | $a \wedge b \Rightarrow b \wedge a$ | $\Rightarrow$ Intro 1-4 |

# Warm-up Exercise

Use the propositional logic inference rules provided to prove:
$$(a \land b) \Rightarrow (b \land a)$$

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

Goal is to prove $X$

1.    $\phi \vee \psi$       Assume

2.   Case 1: $\phi$
   - a)    _____
   - b)    _____
   - c)   $\phi \Rightarrow X$     _____

3.   Case 2: $\psi$
   - a)    _____
   - b)    _____
   - c)    _____
   - d) $\psi \Rightarrow X$     _____

4. $X$              by $\vee$ elimination

# Today

Course Info

Warm-up exercise

**Propositional Logic and Proofs**

ML and 606/607 Intro

More Course Info

Images: ai.berkeley.edu

# Analysis: Perceptron

**Perceptron Mistake Bound**
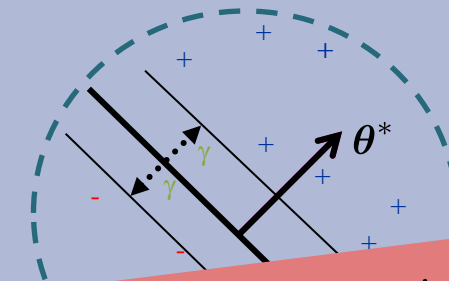
**Theorem 0.1** (Block (1962), Novikoff (1962)).

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

*Suppose:*

1. *Finite size inputs: $||x^{(i)}|| \leq R$*
2. *Linearly separable data: $\exists \boldsymbol{\theta}^* \text{ s.t. } ||\boldsymbol{\theta}^*|| = 1$ 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*
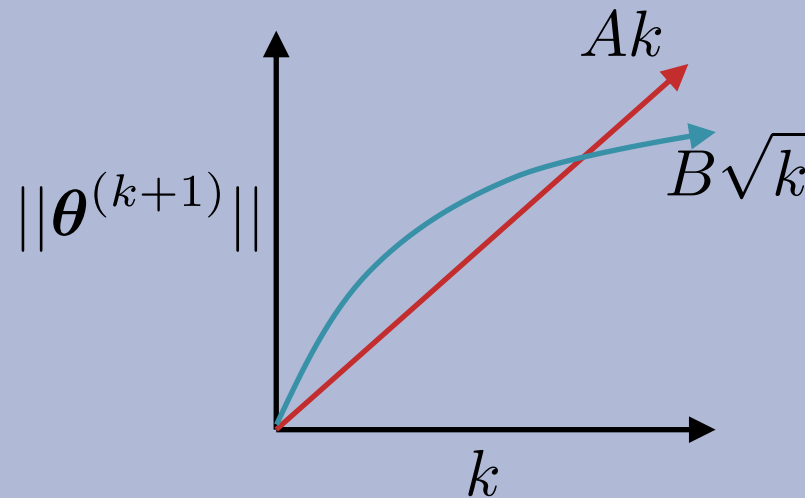
$$k \leq (R/\gamma)^2$$



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

Figure from Nina Balcan

# 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

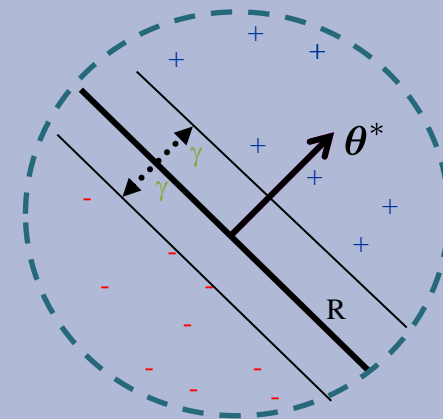**Theorem 0.1** (Block (1962), Novikoff (1962)).
*Given dataset:* $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^{N}$.
*Suppose:*

1. *Finite size inputs:* $||x^{(i)}|| \leq R$
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($\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \ldots\}$)
2:      $\boldsymbol{\theta} \leftarrow \mathbf{0}, k = 1$          ▷ Initialize parameters
3:      **for** $i \in \{1, 2, \ldots\}$ **do**          ▷ For each example
4:          **if** $y^{(i)}(\boldsymbol{\theta}^{(k)} \cdot \mathbf{x}^{(i)}) \leq 0$ **then**      ▷ If mistake
5:            $\boldsymbol{\theta}^{(k+1)} \leftarrow \boldsymbol{\theta}^{(k)} + y^{(i)}\mathbf{x}^{(i)}$     ▷ 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, $Ak \leq ||\boldsymbol{\theta}^{(k+1)}||$

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

      by Perceptron algorithm update

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

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

      by assumption

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

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

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

      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 $(\theta^{(1)})^2 = 0$

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

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: $P \lor (\neg Q \land R)$;     $X_1 \Leftrightarrow (Raining \Rightarrow Sunny)$
▪ Possible world: {P=true, Q=true, R=false, S=true} or 1101
▪ Semantics: $\alpha \land \beta$ is true in a world iff is $\alpha$ true and $\beta$ is true (etc.)

## First-order logic

▪ Syntax: $\forall x \, \exists y \, P(x,y) \land \neg Q(Joe, f(x)) \Rightarrow f(x)=f(y)$
▪ Possible world: Objects $o_1$, $o_2$, $o_3$; P holds for $<o_1,o_2>$; Q holds for $<o_3>$; $f(o_1)=o_1$; Joe=$o_3$; etc.
▪ Semantics: $\phi(\sigma)$ is true in a world if $\sigma=o_j$ and $\phi$ 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

| A | B | A $\Rightarrow$ B |
|---|---|-----|
| T | T | T |
| T | F | F |
| F | T | T |
| F | F | T |

# Poll 1

If we know that $A \lor B$ and $\neg B \lor C$ are true,

what do we know about $A \lor C$?

i.    $A \lor C$ is guaranteed to be true

ii.   $A \lor C$ is guaranteed to be false

iii.  We don't have enough information to say anything
     definitive about $A \lor C$

# Poll 1

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

| $A$ | $B$ | $C$ | $A \lor B$ | $\neg B \lor C$ | $A \lor C$ |
|---|---|---|---|---|---|
| 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 $A \lor B$ and $\neg B \lor C$ are true, what do we know about $A \lor C$?

| $A$ | $B$ | $C$ | $A \lor B$ | $\neg B \lor C$ | $A \lor C$ |
|-----|-----|-----|------------|------------------|------------|
| 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 $A \lor B$ and $\neg B \lor C$ are true,

what do we know about $A \lor C$?

i.      $A \lor C$ is guaranteed to be true

ii.     $A \lor C$ is guaranteed to be false

iii.    We don't have enough information to say anything
        definitive about $A \lor C$

# Poll 2

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

| $A$ | $B$ | $C$ | $A \lor B$ | $\neg B \lor C$ | $A \lor C$ |
|---|---|---|---|---|---|
| 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 $A \lor B$ and $\neg B \lor C$ 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 $\{X_1, X_2, ..., X_n\}$
- (we often add True and False for convenience)

$X_i$ is a sentence

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

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

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

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

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

And p.s. there are no other sentences!

# Notes on Operators

$\boldsymbol{\alpha} \lor \boldsymbol{\beta}$  is <u>inclusive or</u>, not exclusive

# Truth Tables

**α** ∨ **β**  is <u>inclusive or</u>, not exclusive

| **α** | **β** | **α ∧ β** |
|:---:|:---:|:---:|
| F | F | F |
| F | T | F |
| T | F | F |
| T | T | T |

| **α** | **β** | **α ∨ β** |
|:---:|:---:|:---:|
| F | F | F |
| F | T | T |
| T | F | T |
| T | T | T |

# Notes on Operators

$\boldsymbol{\alpha} \vee \boldsymbol{\beta}$  is <u>inclusive or</u>, not exclusive

$\boldsymbol{\alpha} \Rightarrow \boldsymbol{\beta}$  is equivalent to  $\neg\boldsymbol{\alpha} \vee \boldsymbol{\beta}$

- Says who?

# Truth Tables

**α ⇒ β** is equivalent to **¬α ∨ β**

| α | β | α ⇒ β | ¬α | ¬α ∨ β |
|---|---|-------|-----|--------|
| 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 <u>inclusive or</u>, 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

**α ⇔ β** is equivalent to (**α ⇒ β**) ∧ (**β ⇒ α**)

| **α** | **β** | **α ⇔ β** | **α ⇒ β** | **β ⇒ α** | **(α⇒β) ∧ (β⇒α)** |
|:---:|:---:|:---:|:---:|:---:|:---:|
| 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:

$$(\boldsymbol{\alpha} \Leftrightarrow \boldsymbol{\beta}) \Leftrightarrow [(\boldsymbol{\alpha} \Rightarrow \boldsymbol{\beta}) \wedge (\boldsymbol{\beta} \Rightarrow \boldsymbol{\alpha})]$$

# Inference Rules

**Modus Ponens**

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

Notation Alert!

**Unit Resolution**

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

**General Resolution**

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

# Propositional Logic

Check if sentence is true in given model

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

function PL-TRUE?($\alpha$,model) returns true or false

    if $\alpha$ is a symbol then return Lookup($\alpha$, model)

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

    if Op($\alpha$) = $\wedge$ then return **and**(PL-TRUE?(Arg1($\alpha$),model),

                                    PL-TRUE?(Arg2($\alpha$),model))

    etc.

(Sometimes called "recursion over syntax")

# Today

Images: ai.berkeley.edu

# What is ML?

607

**Computer Science**

607 | 606
**Optimization**

**Machine Learning**

**Domain of Interest**

606
**Statistics**

(607) | 606
**Linear Algebra**

(607) | 606
**Probability**

606
**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<br>• Computational Complexity<br>• Recursion; Dynamic Programming<br>• Data Structures for ML Algorithms |

Slide credit: CMU MLD Matt Gormley

# Factor Graph Notation

- Variables:

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

- Factors:

$$\psi_\alpha, \psi_\beta, \psi_\gamma, \ldots$$

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

**Joint Distribution**

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



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

time

an

arr

45

# Factors are Tensors

- Factors:

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



| | s | vp | pp | … |
|---|---|---|---|---|
| s | 0 | 2 | .3 | |
| vp | 3 | 4 | 2 | |
| pp | .1 | 2 | 1 | |
| … | | | | |

| | v | n | p | d |
|---|---|---|---|---|
| v | 1 | 6 | 3 | 4 |
| n | 8 | 4 | 2 | 0.1 |
| p | 1 | 3 | 1 | 3 |
| d | 0.1 | 8 | 0 | 0 |

| v | 3 |
|---|---|
| n | 4 |
| p | 0.1 |
| d | 0.1 |

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

time

46

# Inference

Given a factor graph, two common tasks …

- Compute the most likely joint assignment,
  $$x^* = \text{argmax}_x\, p(X=x)$$

- 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**.

$p(X_i=x_i)$ = sum of $p(X=x)$ 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: $p(\boldsymbol{x}) = \frac{1}{Z} \prod_{\alpha} \psi_\alpha(\boldsymbol{x_\alpha})$



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

48

# 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



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

49

# Marginals by Sampling on Factor Graph



Estimate the marginals as:

|   |     |
|---|-----|
| n | 4/6 |
| v | 2/6 |

|   |     |
|---|-----|
| n | 3/6 |
| v | 3/6 |

|   |     |
|---|-----|
| p | 4/6 |
| v | 2/6 |

| d | 6/6 |
|---|-----|

| n | 6/6 |
|---|-----|

| | $X_1$ | $X_2$ | $X_3$ | $X_4$ | $X_5$ |
|---|---|---|---|---|---|
| Sample 1: | n | v | p | d | n |
| Sample 2: | n | n | v | d | n |
| Sample 3: | n | v | p | d | n |
| Sample 4: | v | n | p | d | n |
| Sample 5: | v | n | v | d | n |
| Sample 6: | n | v | p | d | n |

$X_0$ — $\psi_0$ — $X_1$ — $\psi_2$ — $X_2$ — $\psi_4$ — $X_3$ — $\psi_6$ — $X$

&lt;START&gt;

$\psi_1$    $\psi_3$    $\psi_5$

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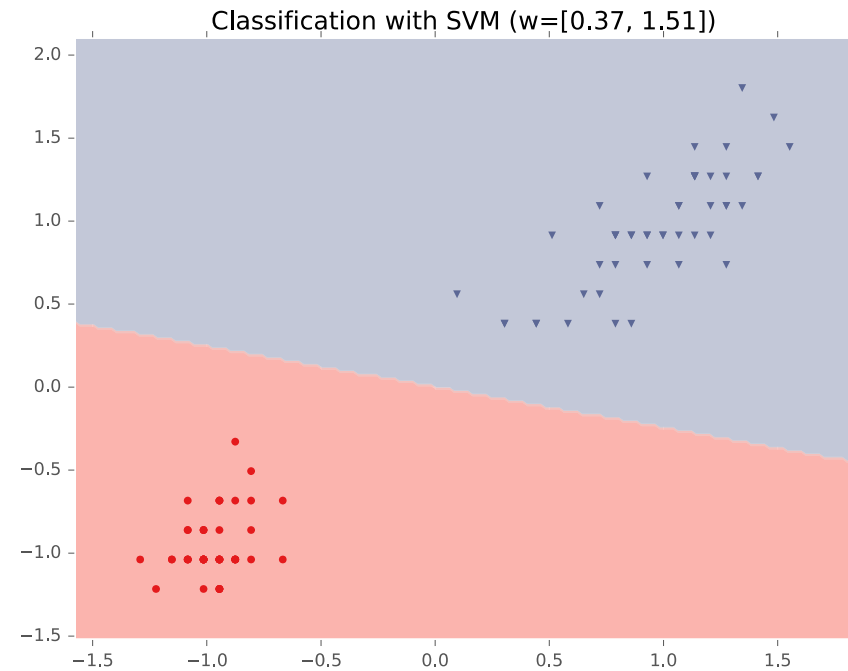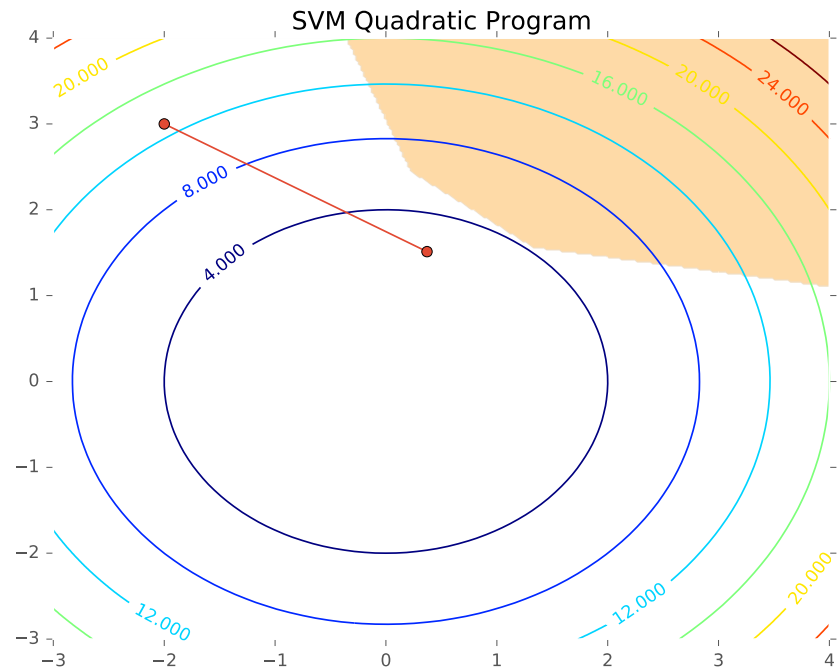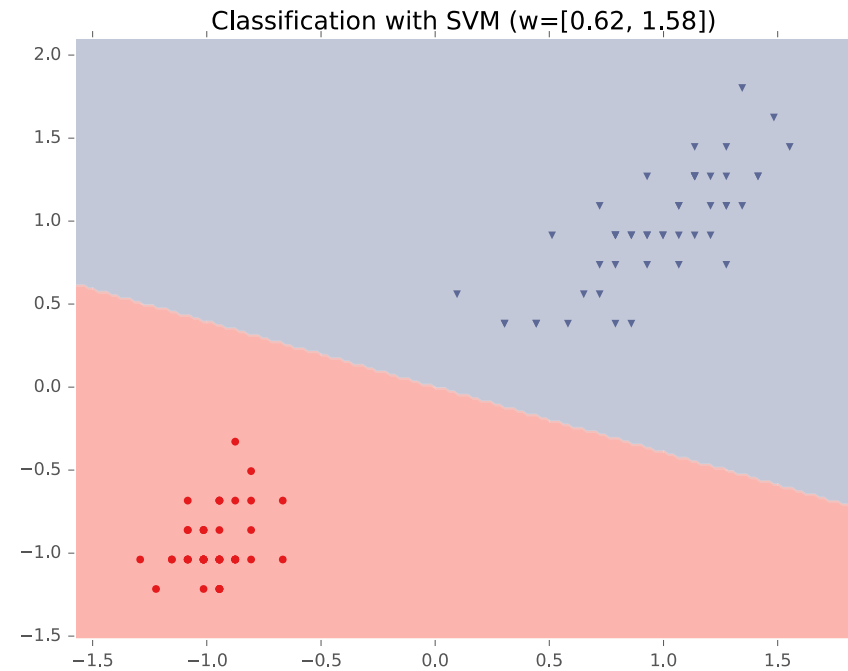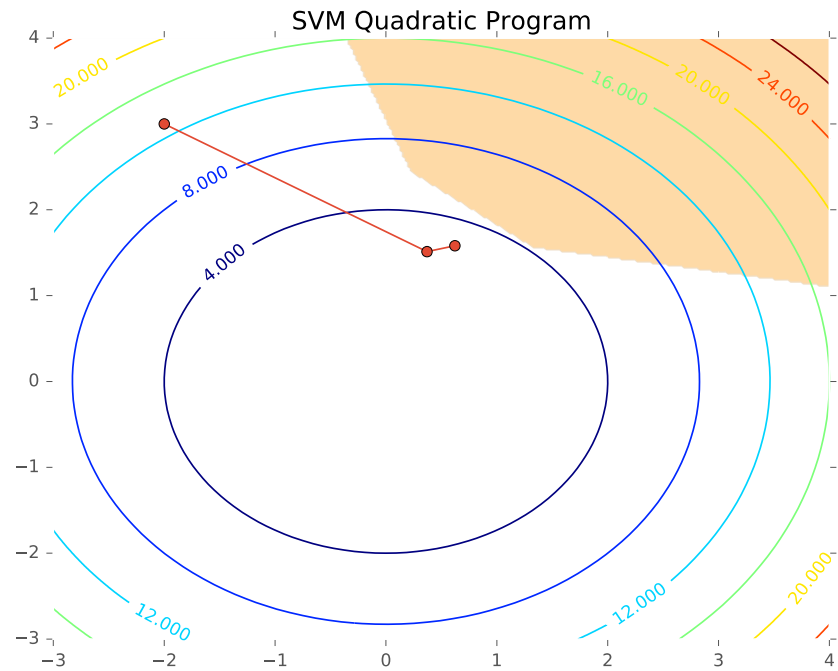
Note: This is just motivation – we'll cover the math need to understand these topics later!

Slide credit: CMU MLD Matt Gormley

50

# Why Computer Science for ML?

To best understand __A__ we need __B__

| A | B |
|---|---|
| Analysis of Exact Inference in Graphical Models | Computation<br>• Computational Complexity<br>• Recursion; Dynamic Programming<br>• Data Structures for ML Algorithms |
| Implementation Design of a Deep Learning Library | Programming & Efficiency<br>• Debugging for Machine Learning<br>• Efficient Implementation / Profiling ML Algorithms |

Slide credit: CMU MLD Matt Gormley

# Finite Difference Method

The *centered* finite difference approximation is:

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

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

**Notes:**

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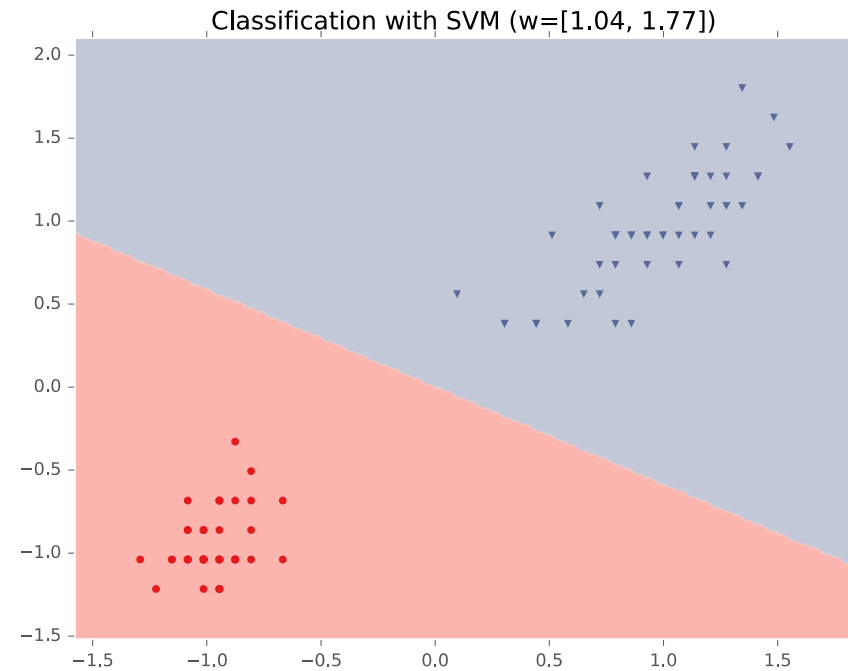- Typically only appropriate to use on small examples with an appropriately chosen epsilon



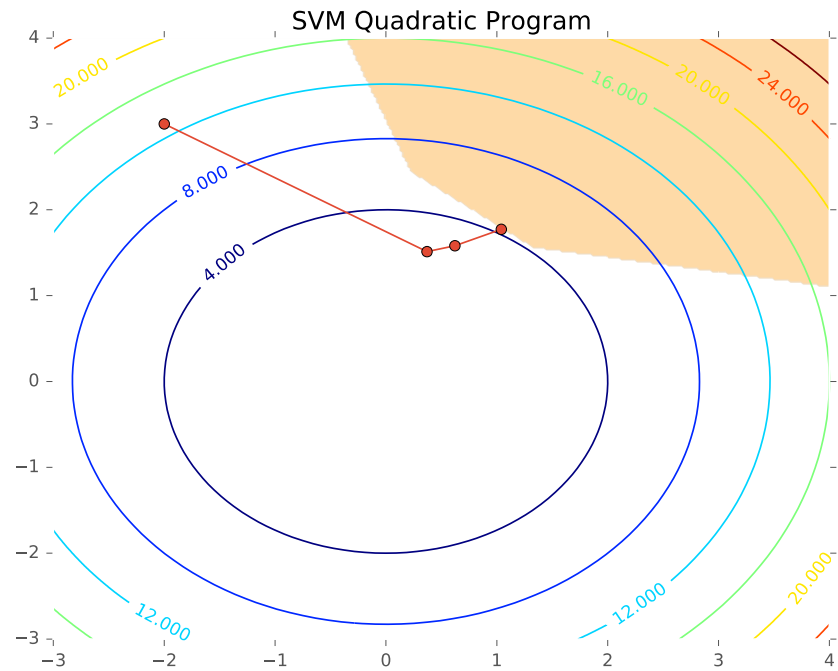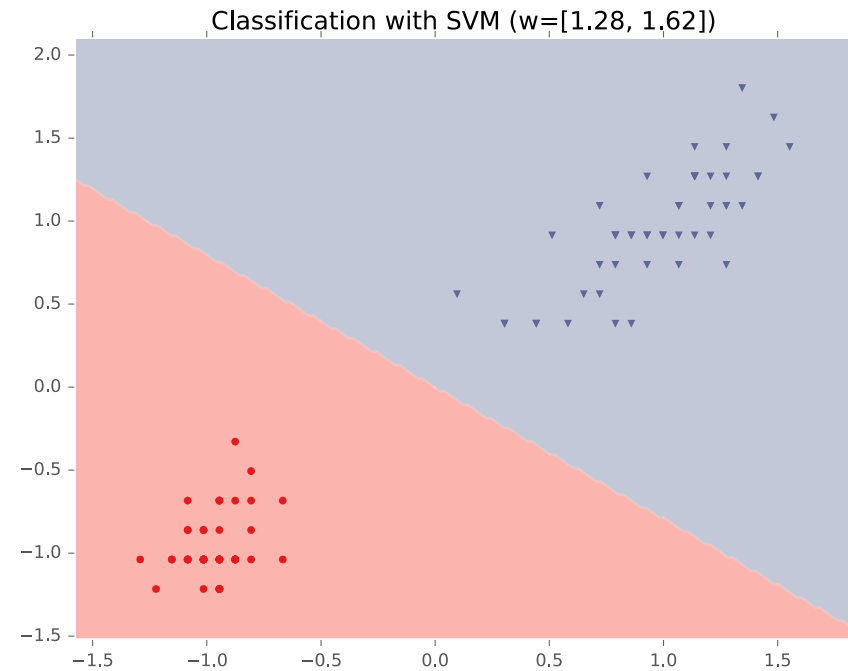Note: This is just motivation – we'll cover the math need to understand these topics later!

52

# 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!

53

# Training
# Backpropagation

**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 $v_1, \ldots, v_N$
   a. Compute $u_i = g_i(v_1, \ldots, v_N)$
   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 $u_i = g_i(v_1, \ldots, v_N)$
   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 $(du_i/dv_j)$

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

**Return** partial derivatives $dy/du_i$ for all variables

54

# Why Computer Science for ML?

10-607

To best understand __A__ we need __B__

| A | B |
|---|---|
| Analysis of Exact Inference in Graphical Models | Computation<br>• Computational Complexity<br>• Recursion; Dynamic Programming<br>• Data Structures for ML Algorithms |
| Implementation Design of a Deep Learning Library | Programming & Efficiency<br>• Debugging for Machine Learning<br>• Efficient Implementation / Profiling ML Algorithms |
| Optimization for Support Vector Machines (SVMs) | Optimization<br>• Unconstrained Optimization<br>• Preconditioning<br>• Constrained Optimization |

Slide credit: CMU MLD Matt Gormley

# Support Vector Machines (SVMs)

**Hard-margin SVM (Primal)**

$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|_2^2$$

$$\text{s.t. } y^{(i)}(\mathbf{w}^T\mathbf{x}^{(i)} + b) \geq 1, \quad \forall i = 1,\ldots,N$$

**Hard-margin SVM (Lagrangian Dual)**

$$\max_{\boldsymbol{\alpha}} \sum_{i=1}^{N} \alpha_i - \frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N} \alpha_i \alpha_j y^{(i)} y^{(j)} \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}$$

$$\text{s.t. } \alpha_i \geq 0, \quad \forall i = 1,\ldots,N$$

$$\sum_{i=1}^{N} \alpha_i y^{(i)} = 0$$

- 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 α$^{(i)}$ ≠ 0

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

# SVM QP



SVM Quadratic Program

Classification with SVM (w=[-2.00, 3.00])

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

# SVM QP



SVM Quadratic Program

Classification with SVM (w=[0.37, 1,51])

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

58

# SVM QP



SVM Quadratic Program

Classification with SVM (w=[0.62, 1.58])

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

59

# SVM QP



SVM Quadratic Program

Classification with SVM (w=[1.04, 1.77])

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

Slide credit: CMU MLD Matt Gormley

60

# SVM QP



SVM Quadratic Program

Classification with SVM (w=[1.28, 1.62])

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

# SVM QP



SVM Quadratic Program

Classification with SVM (w=[1.28, 1.60])

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

Slide credit: CMU MLD Matt Gormley

62

# Why Computer Science for ML?

To best understand __A__ we need __B__

| A | B |
|---|---|
| Analysis of Exact Inference in Graphical Models | Computation<br>• Computational Complexity<br>• Recursion; Dynamic Programming<br>• Data Structures for ML Algorithms |
| Implementation Design of a Deep Learning Library | Programming & Efficiency<br>• Debugging for Machine Learning<br>• Efficient Implementation / Profiling ML Algorithms |
| Optimization for Support Vector Machines (SVMs) | Optimization<br>• Unconstrained Optimization<br>• Preconditioning<br>• 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)

Slide credit: CMU MLD Matt Gormley

# 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

http://www-formal.stanford.edu/jmc/whatisai/whatisai.html

# AI Stack for CMU AI

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



| Autonomy | Human AI Interaction |
|---|---|
| Planning & Acting | |
| Decision Support | |
| Modeling | |
| Machine Learning | |
| Massive Data Management | |
| Devices | |
| Computing | |

Ethics

AI Stack

https://ai.cs.cmu.edu/about

# 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."

https://ai.cs.cmu.edu/about

# Artificial Intelligence vs Machine Learning?

# A Brief History of AI

# A Brief History of AI



AI Excitement! 1950-1970

Knowledge Based Systems 1970-1990

Statistical Approaches 1990-

Deep Learning Era 2012-

machine learning

artificial intelligence

deep learning

formal logic

0.000400%
0.000350%
0.000300%
0.000250%
0.000200%
0.000150%
0.000100%
0.000050%
0.000000%

1940  1950  1960  1970  1980  1990  2000  2010

https://books.google.com/ngrams

# 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 |
|------|-----|
| "…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…"<br><br>(Mitchell, 1997) |  |

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

# Robotics

## 2. Learning to drive an autonomous vehicle

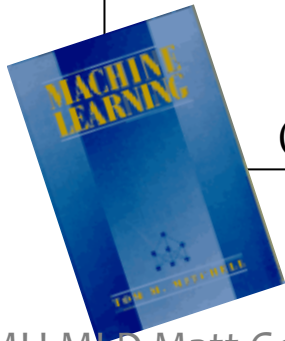| THEN | NOW |
|---|---|
| "…the ALVINN system (Pomerleau 1989) has used its learned strategies to drive unassisted at 70 miles per hour for 90 miles on public highways among other cars…" <br><br> (Mitchell, 1997) |  <br> https://www.geek.com/wp-content/uploads/2016/03/uber.jpg |

# Games / Reasoning

## 3. Learning to beat the masters at board games

| THEN | NOW |
|---|---|

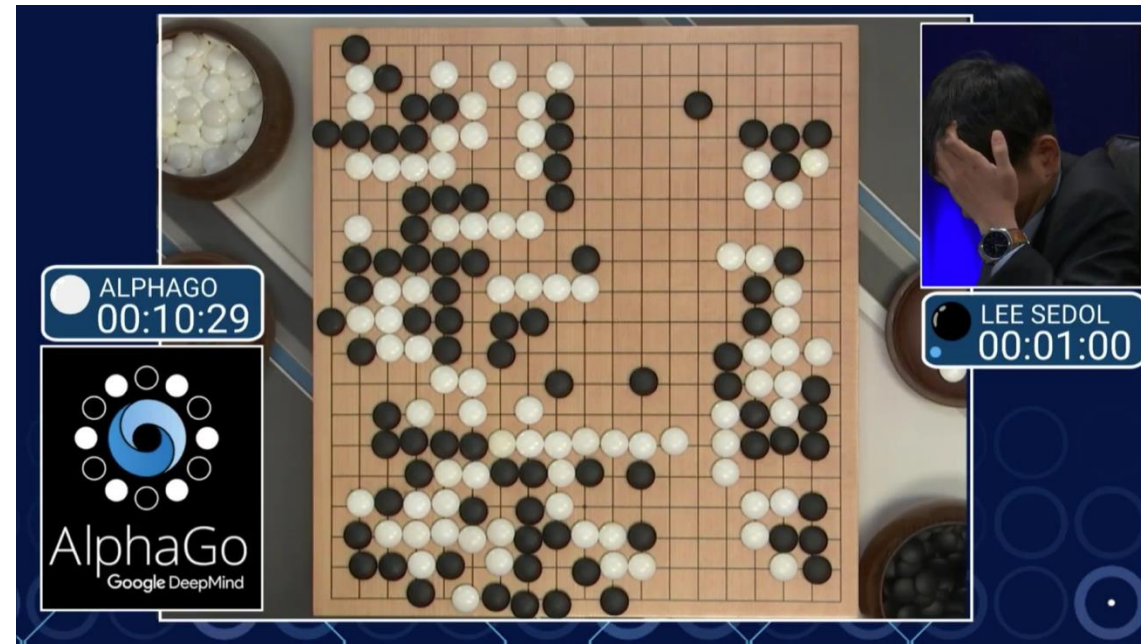"…the world's top computer program for backgammon, TD-GAMMON (Tesauro, 1992, 1995), learned its strategy by playing over one million practice games against itself…"

(Mitchell, 1997)

# Computer Vision

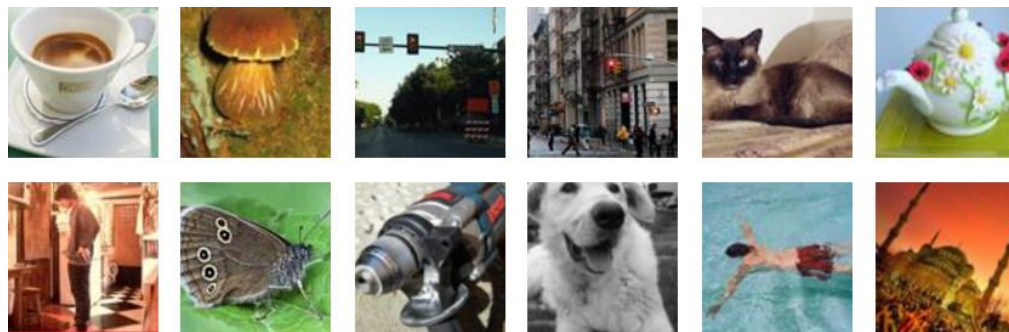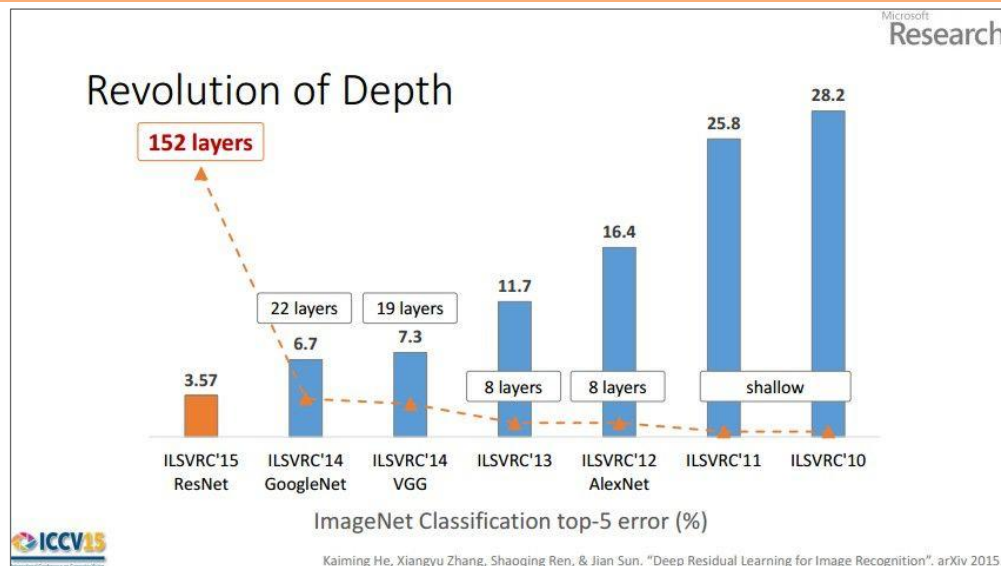## 4. Learning to recognize images

| THEN | NOW |
|------|-----|

"…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)

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



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
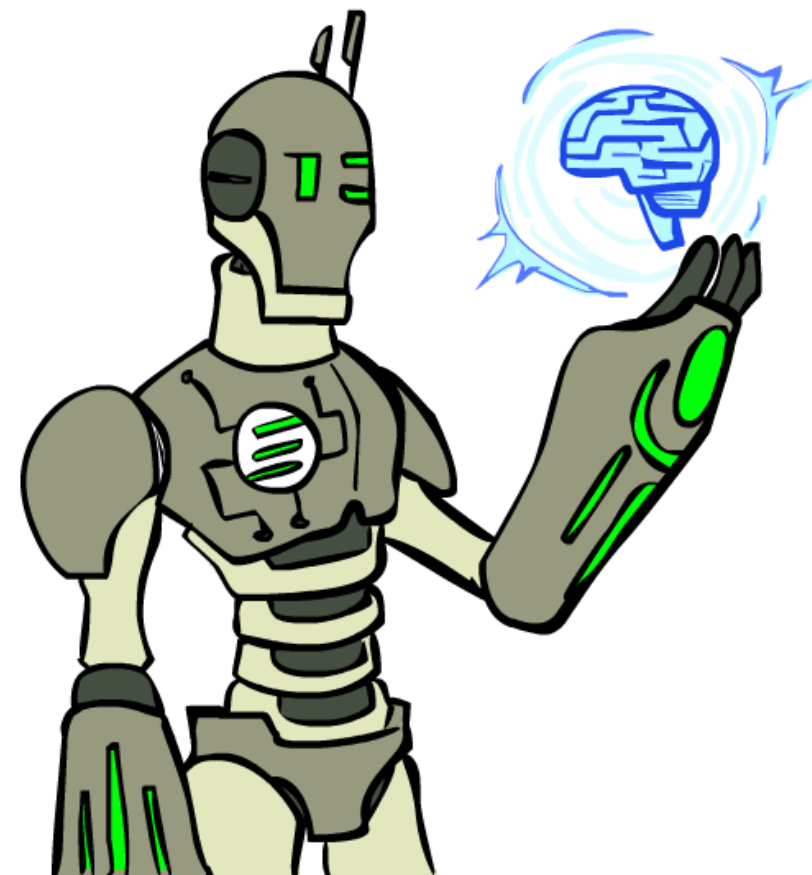  - 10-701
  - 10-715

- Prerequisites

# Today

Course Info

Warm-up exercise

Propositional Logic and Proofs

ML and 606/607 Intro

More Course Info

Images: ai.berkeley.edu

# 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