# Sets, Data Types, and Functions

Matt Gormley
Lecture 2
August 29, 2018
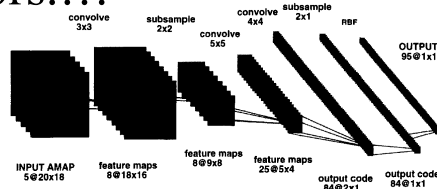
# Computer Vision

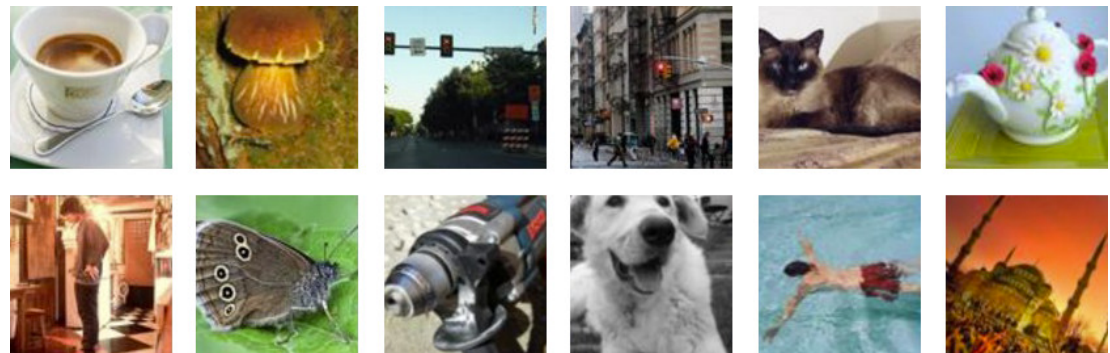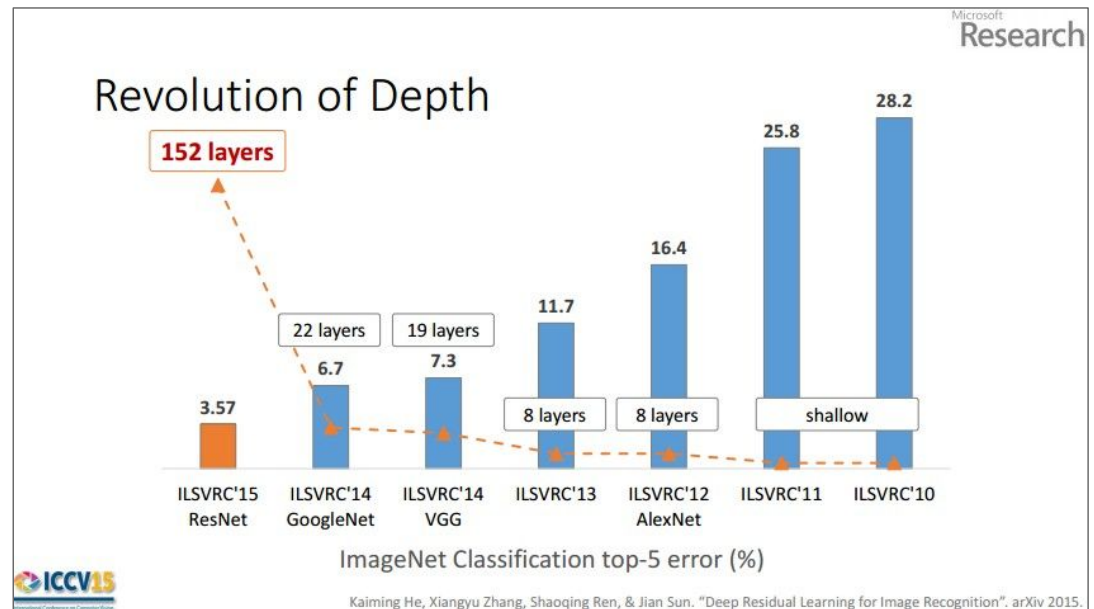## 4. Learning to recognize images

| 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…."



(LeCun et al., 1995)

Images from https://blog.openai.com/generative-models/

# SYLLABUS HIGHLIGHTS

# Syllabus Highlights

The syllabus is located on the course webpage:

http://www.cs.cmu.edu/~mgormley/courses/606-607-f18

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

# 606/607 Syllabus Highlights

- **Grading**: 55% homework, 10% in-class quizzes, 30% final exam, 5% participation
- **Final Exam**:
  - 606: Mini-I final exam week, date TBD
  - 607: Mini-II final exam week, date TBD
- **In-Class Quizzes**: always announced ahead of time
- **Homework**: 4 assignments with written / programming portions
  - 2 grace days for the unexpected
  - Late submissions: 80% day 1, 60% day 2, 40% day 3, 20% day 4
  - No submissions accepted after 4 days w/o extension
  - Extension requests: see syllabus

- **Recitations**: Fridays, same time/place as lecture (optional, interactive sessions)
- **Readings**: required, online, recommended for after lecture
- **Technologies**: Piazza (discussion), Gradescope (homework), Canvas (gradebook only)
- **Academic Integrity**:
  - Collaboration encouraged, but must be documented
  - Solutions must always be written independently
  - No re-use of found code / past assignments
  - Severe penalties (i.e. failure)
- **Office Hours**: posted on Google Calendar on "People" page

# 606/607 Syllabus Highlights

- **Grading**: 55% homework, 10% in-class quizzes, 30% final exam, 5% participation
- **Final Exam**:
  - 606: Mini-I final exam week, date TBD
  - 607: Mini-II final exam week, date TBD
- **In-Class Quizzes**: always announced ahead of time
- **Homework**: 4 assignments with written / programming portions
  - 2 grace days for the unexpected
  - Late submissions: 80% day 1, 60% day 2, 40% day 3, 20% day 4
  - No submissions accepted after 4 days w/o extension
  - Extension requests: see syllabus

- **Recitations**: Fridays, same time/place as lecture (optional, interactive sessions)
- **Readings**: required, online, recommended for after lecture
- **Technologies**: Piazza (discussion), Gradescope (homework), Canvas (gradebook only)
- **Academic Integrity**:
  - Collaboration encouraged, but must be documented
  - Solutions must always be written independently
  - No re-use of found code / past assignments
  - Severe penalties (i.e. failure)
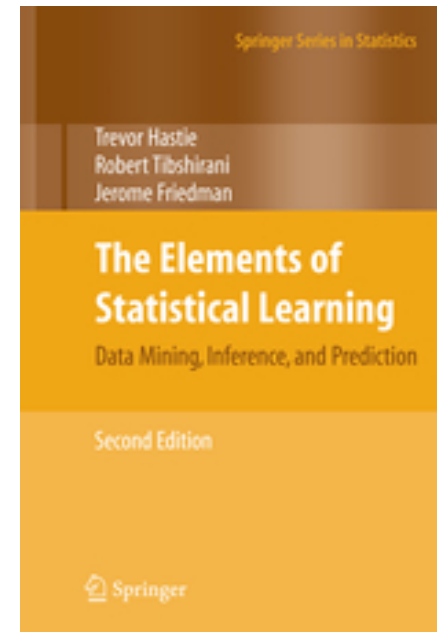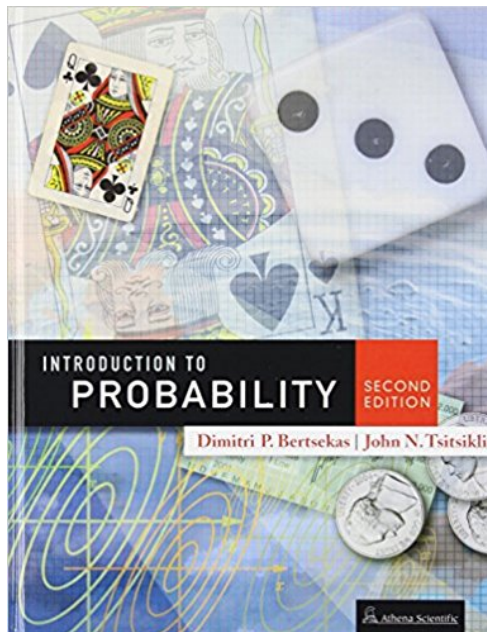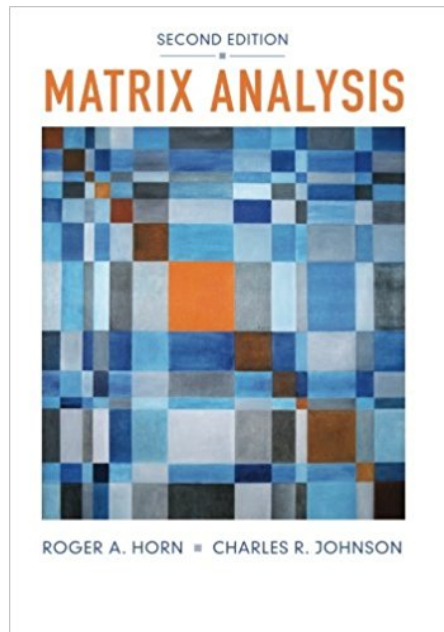- **Office Hours**: posted on Google Calendar on "People" page

# 606/607 Syllabus Highlights

- **Grading**: 55% homework, 10% in-class quizzes, 30% final exam, 5% participation
- **Final Exam**:
  - 606: Mini-I final exam week, date TBD
  - 607: Mini-II final exam week, date TBD
- **In-Class Quizzes**: always announced ahead of time
- **Homework**: 4 assignments with written / programming portions
  - 2 grace days for the unexpected
  - Late submissions: 80% day 1, 60% day 2, 40% day 3, 20% day 4
  - No submissions accepted after 4 days w/o extension
  - Extension requests: see syllabus

- **Recitations**: Fridays, same time/place as lecture (optional, interactive sessions)
- **Readings**: required, online, recommended for after lecture
- **Technologies**: Piazza (discussion), Gradescope (homework), Canvas (gradebook only)
- **Academic Integrity**:
  - Collaboration encouraged, but must be documented
  - Solutions must always be written independently
  - No re-use of found code / past assignments
  - Severe penalties (i.e. failure)
- **Office Hours**: posted on Google Calendar on "People" page

# Lectures

- You should ask lots of questions
  - Interrupting (by raising a hand) to ask your question is strongly encouraged
  - Asking questions later 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

These are optional, but highly recommended as an alternate presentation of the material

# Expected Background

## 10-606 (Math Background 4 ML)

You should be familiar with some of the following...

- Calculus:
  - can take scalar derivatives
  - can solve scalar integrals
- Linear Algebra:
  - know basic vector operations
  - seen matrix multiplication
- Probability:
  - seen the basics: conditioning, Bayes Rule, etc.
- Programming:
  - know some Python
    **OR**
    have sufficient programming background to pick up the basics of Python

But we'll offer practice to make sure you can catch up on your weaker areas

## 10-607 (CS Background 4 ML)

You should...

- be comfortable with all the topics listed for 10-606
- ideally, have the mathematical maturity of someone who completed 10-606 **because** it will aide in understanding the motivating examples from machine learning

That said, the content of 10-607 is designed stand alone

# Q&A

**Q:** Is this course right for me?

**A:**
- If you're a Master's or PhD and you lack some of the prerequisite material for 10-601/701, this is definitely the right course for you!
- If you're a Master's or PhD and you studied the prerequisite material for 10-601/701… but it was a long time ago, this is certainly the right course for you.
- If you're an undergrad, I would recommend the usual prereq sequence required for 10-601/701.
- For ugrad/MS/PhD: If you tried taken an Intro ML course here and felt a bit lost in all the math/CS, this is a great place to start.

**Q:** What calculus textbook would you recommend?

**A:** Good question… I'm still working on that one.

# Q&A

**Q:** What do I do if I'm feeling a bit lost in the math and CS in 10-606/607?

**A:** Let me know ASAP! (Do so in office hours, Piazza note, the middle of class, etc.) Our goal is to provide you with a learning environment in which you thrive. We'll certainly make adjustments if we need to.

# MOTIVATION: SETS & TYPES

Sets, Types, and Functions show up **everywhere** in Machine Learning

# (Negative) Gradients



These are the **negative** gradients that
Gradient **Descent** would follow.

# Convexity

Suppose we wish to define convexity of a function...

We could draw a picture.



...but that's a bit informal.

So instead, we could offer a mathematical definition.

> Function $f : \mathbb{R}^M \to \mathbb{R}$ is **convex**
> if $\forall \ \mathbf{x}_1 \in \mathbb{R}^M, \mathbf{x}_2 \in \mathbb{R}^M, 0 \le t \le 1$:
>
> $$f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \le tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$

...but even this definition requires some carefully defined objects (the vectors, a function, the set of reals, the set of real-valued vectors of length M, etc.)

Slide adapted from William Cohen

# Data for ML

What is the object we talk about more in Machine Learning than anything else?

Our data!

$$\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^{N}$$

The data consists of a set of tuples

$$\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \ldots, (\mathbf{x}^{(N)}, y^{(N)})\}$$

# The importance of sets…

- Gaussian Discriminant Analysis and Gaussian Mixture Models are almost identical
- There's really only one practical difference
- Can you spot it?

See next two slides…

# Gaussian Discriminant Analysis

**Data:** $\mathcal{D} = \{(\mathbf{x}^{(i)}, \mathbf{z}^{(i)})\}_{i=1}^{N}$ where $\mathbf{x}^{(i)} \in \mathbb{R}^M$ and $z^{(i)} \in \{1, \ldots, K\}$

**Generative Story:**
$$z \sim \text{Categorical}(\boldsymbol{\phi})$$
$$\mathbf{x} \sim \text{Gaussian}(\boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$$

**Model:** Joint: $\quad p(\mathbf{x}, z; \boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = p(\mathbf{x}|z; \boldsymbol{\mu}, \boldsymbol{\Sigma}) p(z; \boldsymbol{\phi})$

**Log-likelihood:**

$$\ell(\boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \log \prod_{i=1}^{N} p(\mathbf{x}^{(i)}, z^{(i)}; \boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$= \sum_{i=1}^{N} \log p(\mathbf{x}^{(i)}|z^{(i)}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \log p(z^{(i)}; \boldsymbol{\phi})$$

# Gaussian Mixture-Model

**Data:** $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^{N}$ where $\mathbf{x}^{(i)} \in \mathbb{R}^{M}$

**Generative Story:**

$$z \sim \text{Categorical}(\boldsymbol{\phi})$$

$$\mathbf{x} \sim \text{Gaussian}(\boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$$

**Model:**

Joint: $\quad p(\mathbf{x}, z; \boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = p(\mathbf{x}|z; \boldsymbol{\mu}, \boldsymbol{\Sigma})p(z; \boldsymbol{\phi})$

Marginal: $\quad p(\mathbf{x}; \boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{z=1}^{K} p(\mathbf{x}|z; \boldsymbol{\mu}, \boldsymbol{\Sigma})p(z; \boldsymbol{\phi})$

**(Marginal) Log-likelihood:**

$$\ell(\boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \log \prod_{i=1}^{N} p(\mathbf{x}^{(i)}; \boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$= \sum_{i=1}^{N} \log \sum_{z=1}^{K} p(\mathbf{x}^{(i)}|z; \boldsymbol{\mu}, \boldsymbol{\Sigma})p(z; \boldsymbol{\phi})$$

# Notation for ML

Machine Learning is notorious for requiring lots of notation… and not always being terribly consistent about it!

# PRELIMINARIES: SETS AND TYPES

# Sets

*Chalkboard*

- Definitions: Set, element of, equality, subset
- Example: Sets of sets
- Set builder notation
- Python list/set comprehentions
- Exercise: Set builder notation
- Definitions: Union, intersection, difference, complement
- Exercise: Set complement
- Tuples and set product
- Exercise: Set product

# Data Types and Functions

*Chalkboard*

- Data types, structs, unions
- Tagged unions
- Exercise: Tagged unions
- Functions
- Anonymous functions
- Exercises: Functions