

CPSC 340 and 532M: Machine Learning and Data Mining

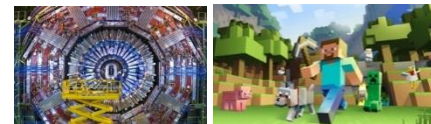
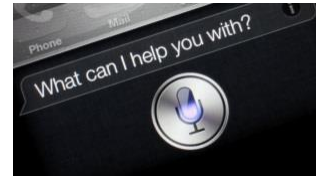
Mark Schmidt

University of British Columbia, Fall 2019

www.cs.ubc.ca/~schmidtm/Courses/340-F19

Big Data Phenomenon

- We are **collecting and storing data** at an unprecedented rate.
- Examples:
 - YouTube, Facebook, MOOCs, news sites.
 - Credit cards transactions and Amazon purchases.
 - Transportation data (Google Maps, Waze, Uber)
 - Gene expression data and protein interaction assays.
 - Maps and satellite data.
 - Large hadron collider and surveying the sky.
 - Phone call records and speech recognition results.
 - Video game worlds and user actions.

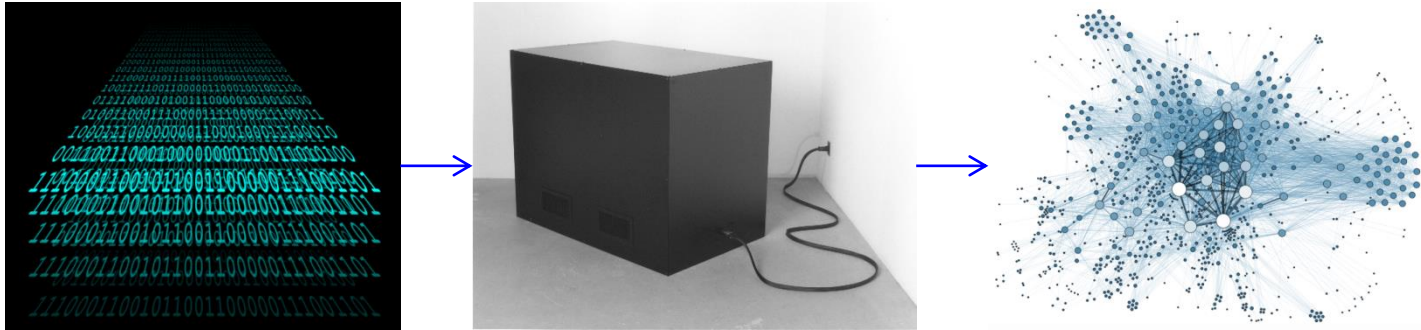


Big Data Phenomenon

- What do you do with all this data?
 - Too much data to search through it manually.
- But there is valuable information in the data.
 - How can we use it for fun, profit, and/or the greater good?
- Data mining and machine learning are key tools we use to make sense of large datasets.

Data Mining

- Automatically **extract useful knowledge** from large datasets.



- Usually, to help with human decision making.

Machine Learning

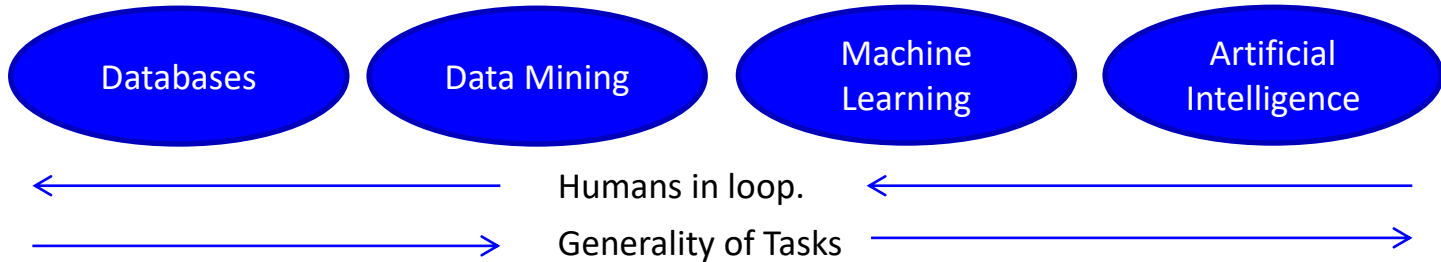
- Using computer to automatically **detect patterns in data and use these to make predictions** or decisions.



- Most useful when:
 - We want to automate something a human can do.
 - We want to do things a human can't do (look at 1 TB of data).

Data Mining vs. Machine Learning

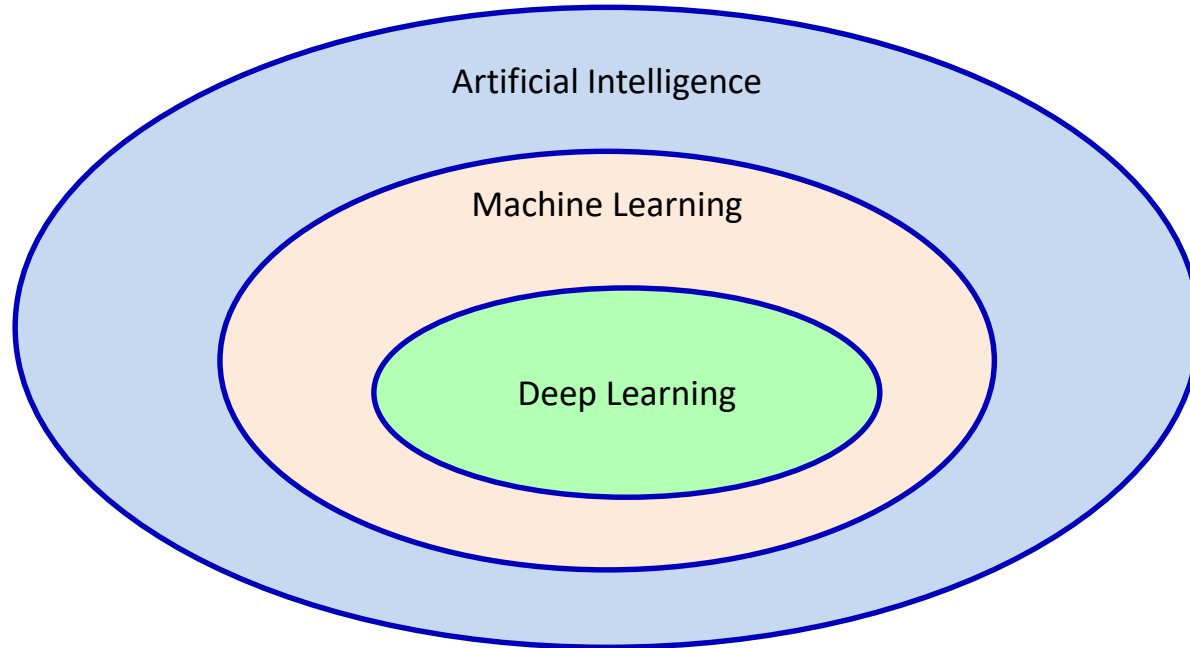
- Data mining and machine learning are very similar:
 - Data mining often viewed as closer to databases.
 - Machine learning often viewed as closer AI.



- Both are similar to statistics, but more emphasis on:
 - Large datasets and computation.
 - Predictions (instead of descriptions).
 - Flexible models (that work on many problems).

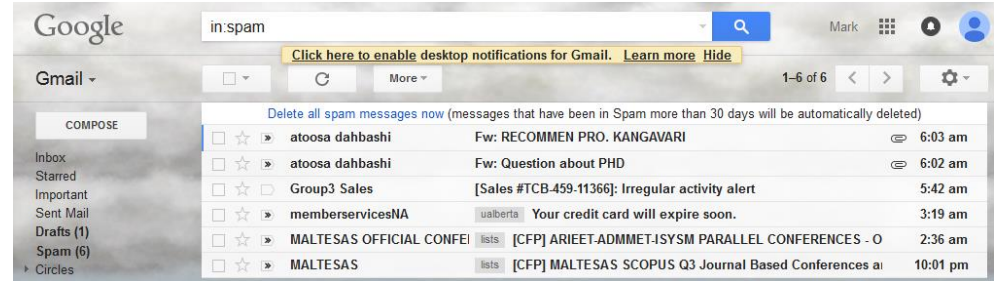
Deep Learning vs. Machine Learning vs. AI

- Traditional we've viewed ML as a subset of AI.
 - And “deep learning” as a subset of ML.



Applications

- Spam filtering:
- Credit card fraud detection:
- Product recommendation:



| Transaction Date | Posted Date | Transaction Details | Debit | Credit |
|------------------|---------------|-------------------------------------|---------|--------|
| Aug. 27, 2015 | Aug. 28, 2015 | BEAN AROUND THE WORLD VANCOUVER, BC | \$10.95 | |

Customers Who Bought This Item Also Bought

Page 1 of 20

Pattern Recognition and Machine Learning (Information Science and...)
 Christopher Bishop
 ★★★★★☆ 115
 Hardcover
 \$60.76 ✓Prime

Learning From Data
 Yaser S. Abu-Mostafa
 ★★★★★☆ 88
 Hardcover

The Elements of Statistical Learning: Data Mining, Inference, and Prediction...
 Trevor Hastie
 ★★★★★☆ 50
 Hardcover
 \$62.82 ✓Prime

Probabilistic Graphical Models: Principles and Techniques (Adaptive...)
 Daphne Koller
 ★★★★★☆ 28
 Hardcover
 \$91.66 ✓Prime

Foundations of Machine Learning (Adaptive Computation and...)
 Mehryar Mohri
 ★★★★★☆ 8
 Hardcover
 \$65.68 ✓Prime

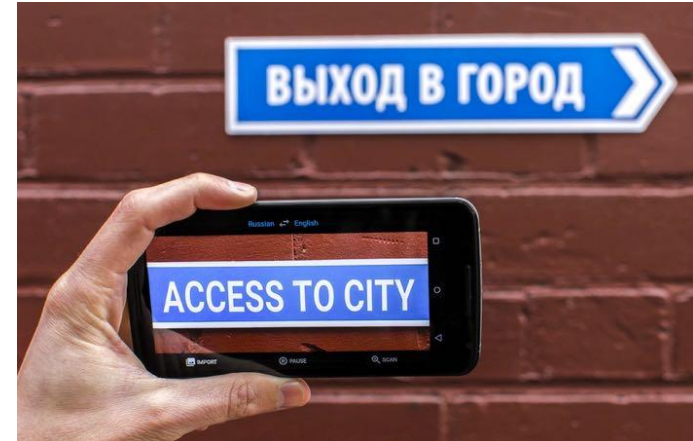
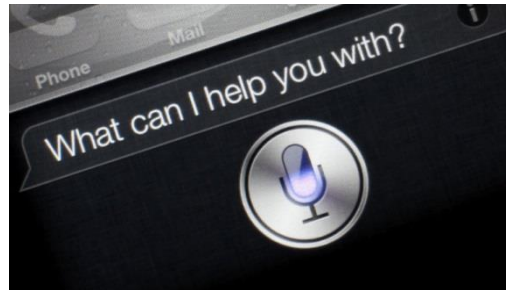
Applications

- Motion capture:



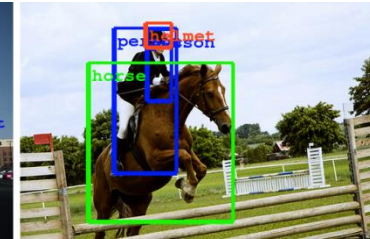
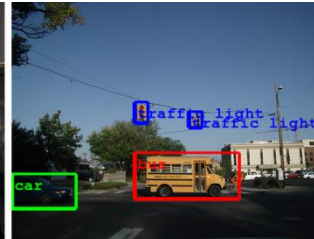
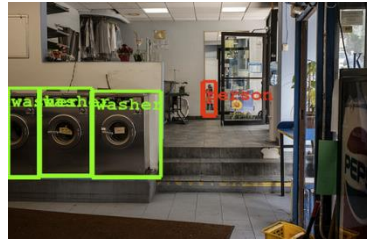
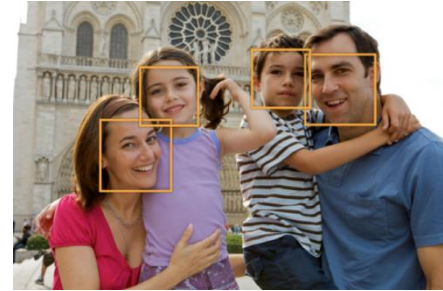
- Optical character recognition and machine translation:

- Speech recognition:

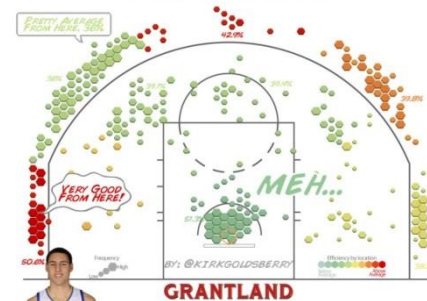


Applications

- Face detection:
- Object detection:
- Sports analytics:

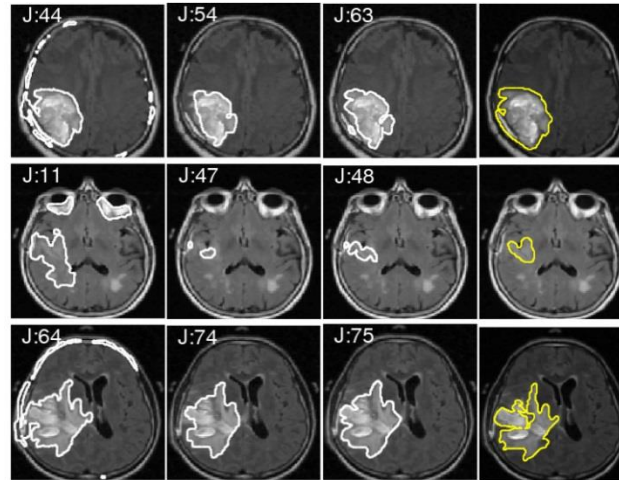


KLAY THOMPSON



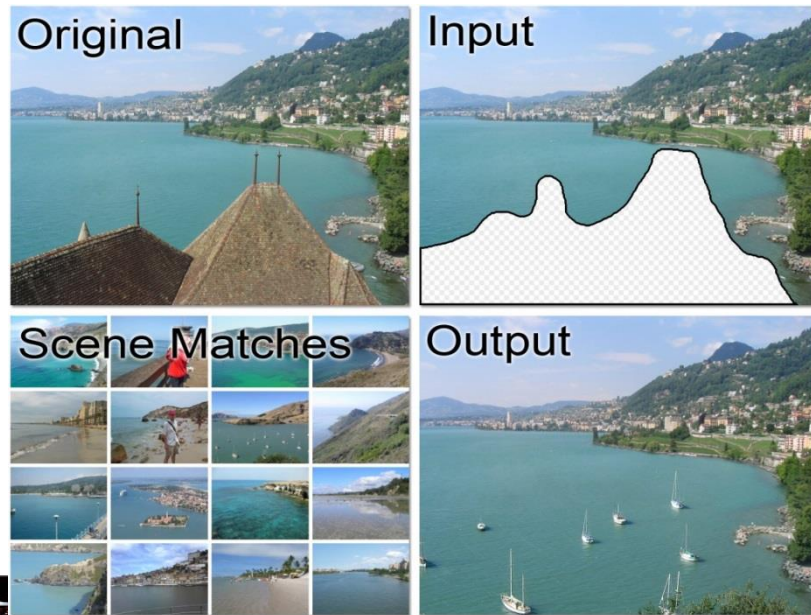
Applications

- Personal Assistants:
- Medical imaging:
- Self-driving cars:



Applications

- Scene completion:



- Image annotation:



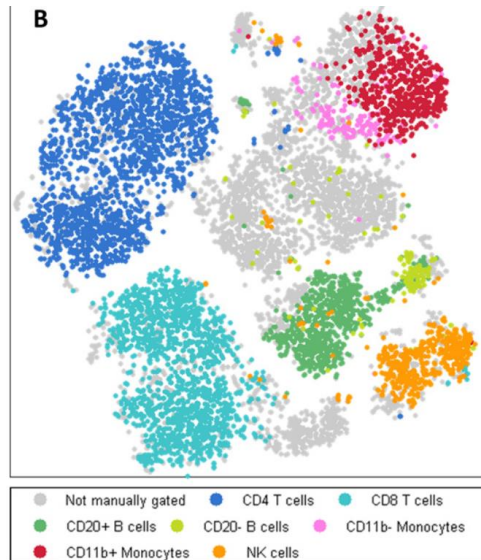
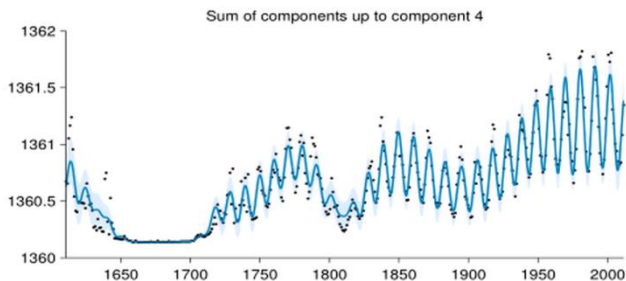
Applications

- Discovering new cancer subtypes:

- Automated Statistician:

2.4 Component 4 : An approximately periodic function with a period of 10.8 years. This function applies until 1643 and from 1716 onwards

This component is approximately periodic with a period of 10.8 years. Across periods the shape of this function varies smoothly with a typical lengthscale of 36.9 years. The shape of this function within each period is very smooth and resembles a sinusoid. This component applies until 1643 and from 1716 onwards.



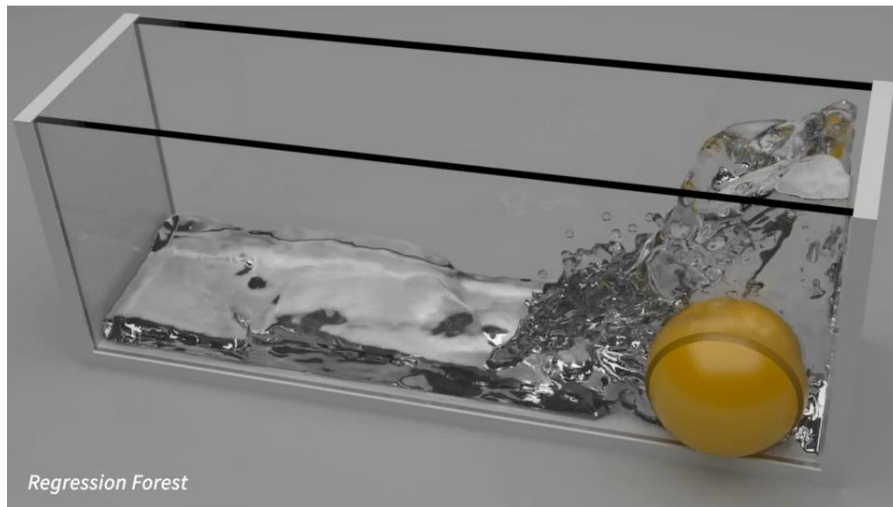
Applications

- Mimicking artistic styles:



Applications

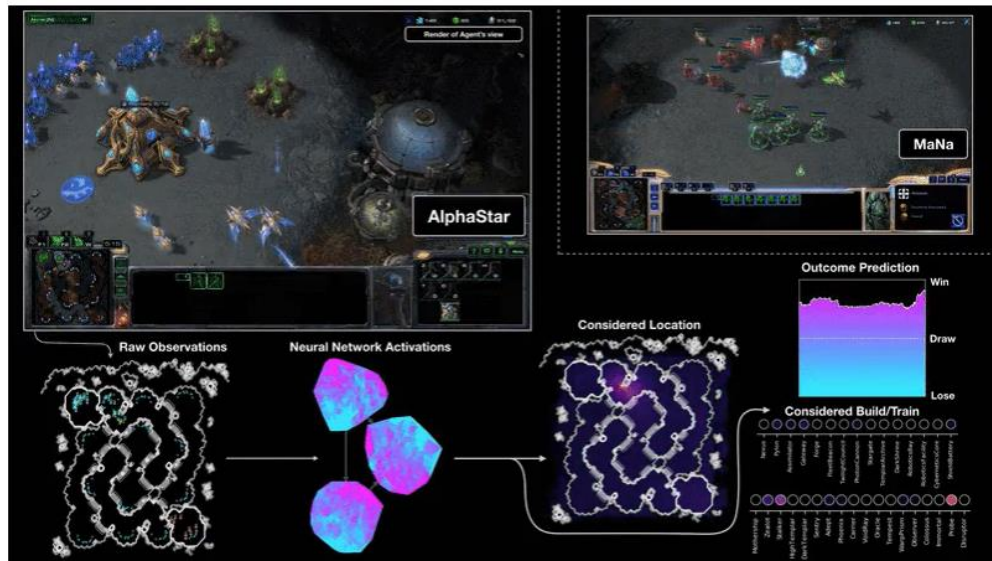
- Fast physics-based animation:



- Mimicking art style in [video](#).
- Recent work on generating text/music/voice/poetry/dance.

Applications

- Beating humans in Go and Starcraft:



- Summary:
 - There is a lot you can do with a bit of statistics and a lot data/computation.
- We are in exciting times.
 - Major recent progress in fields like speech recognition and computer vision.
 - Things are changing a lot on the timescale of 3-5 years.
 - NeurIPS conference sold out in ~11 minutes last year.
 - A bubble in ML investments (most “AI” companies are just doing ML).
- But it is important to know the **limitations** of what you are doing.
 - “The combination of some data and an aching desire for an answer does not ensure that a reasonable answer can be extracted from a given body of data.” – John Tukey
 - A huge number of people applying ML are just “**overfitting**”.
 - Or don’t understand the assumptions needed for them to work.
 - Their **methods do not work** when they are released “into the wild”.

(pause)

Reasons NOT to take this class

- Compared to typical CS classes, there is a **lot more math**:
 - Requires linear algebra, probability, and multivariate calculus (at once).
 - “I think the prerequisites for this course should require that students have obtained at least 75% (or around there) in the required math courses. As someone who who did not excel at math, I felt severely under prepared and struggled immensely in this course, especially seeing that I have taken CPSC courses in the past with similar math requirements, but were not nearly as math heavy as CPSC340.”
- If you’ve only taken a few math courses (or have low math grades),
this course will ruin your life for the next 4 months.
- It’s better to **improve your math, then take this course later.**
 - A good reference covering the relevant math is [here](#) (Chapters 1-3 and 5-6).

Reasons NOT to take this class

- This is not a class on “how to use scikit-learn or TensorFlow or PyTorch”.
 - You will need to **implement things from scratch, and modify existing code.**
- Instead, this is a 300-level computer science course:
 - You are **expected to be able to quickly understand and write code.**
 - You are **expected to be able to analyze algorithms in big-O notation.**
- We’re going to use the **Julia programming language.**
 - A relatively new language that is quickly gaining popularity for machine learning.
 - You are **expected to be able to learn a programming language on your own.**
- If you only have limited programming experience,
this course will ruin your life for the next 4 months.
- It’s better to **get programming experience, then take this course later.**
 - Take CPSC 310 and/or 320 instead, then take this course later.

Reasons NOT to take this class

- Do NOT take this course expecting a high grade with low effort.
- Many people find the assignments very long and very difficult.
 - You will need to put time and effort into learning new/difficult skills.
 - If you aren't strong at math and CS, they may take all of your time.
- Class averages have only been high because of graduate students.
 - NOT because this is an “easy” course, for most people it's not.
- From “Rate My Professors”:
 - “Lectures were dull, dry, and glossed over the material skipping over the theoretical details. Ironically, assignments were detail-heavy and LONG. Doesn't seem to care about students because some of us have 4 other classes and well, if they're all like this course, my girlfriend would have broken up with me two months ago.”

Different Sections of 340 and 532M

- I am teaching both sections of 340 this term.
 - Both sections have the same webpage, assignments, and exams.
- You are free to attend the lectures of the other section.
 - However, don't take a seat if you aren't registered and people are standing.
- Lectures will cover roughly the same set of topics.
 - You will only be tested on material that appears in both sections.
- Next term it will be taught by Frank Wood (probably in Python).
 - Frank and I are research faculty.
 - In other years 340 has been taught by Mike Gelbart (teaching faculty).

CPSC 340 vs. 532M

- One section of CPSC 340 is also **cross-listed as CPSC 532M**.
 - For graduate students who want/need graduate credit.
- Students in CPSC 532M must do a small **research project**.
 - Literature survey on an ML topic not covered in class.
 - Must be done in groups of 2-3.
 - More details later.
- Grading will be slightly different:

| Number | Assignments | Midterm | Final Exam | Survey |
|--------|-------------|---------|------------|--------|
| 340 | 30 | 20 | 50 | 0 |
| 532M | 25 | 15 | 40 | 20 |

CPSC 330 vs. CPSC 340

- There is also a **less-advanced ML course**, **CPSC 330**:
 - Taught by Mike Gelbart for the first time in January.
 - Fewer prerequisites (and probably lower workload).
 - You **can take both** for credit (if you do this then take 330 first).
 - 330 emphasizes “**when to use**” tools, 340 emphasizes “**how they work**”.
 - 330 is more like the Coursera course and other online courses.
- From a former 340 student:
 - “I took Andrew Ng's Coursera course and had a lot of fun and so I would recommend it. But before you spend any time, the Coursera course (I feel) covers only a subset of the concepts covered in this class and wouldn't be an efficient way of gaining understanding of the course material.”

CPSC 340 vs. CPSC 540

- There is also a **more-advanced ML course**, **CPSC 540**:
 - Starts where this course ends.
 - More focus on theory/implementation, less focus on applications.
 - More prerequisites and higher workload.
- For almost all students, **CPSC 340 is the better class to take**:
 - CPSC 330/340 focus on the most widely-used methods in practice.
 - It covers much more material than standard ML classes like Coursera.
 - CPSC 540 focuses on less widely-used methods and research topics.
 - It is intended as a continuation of CPSC 340.
 - You'll miss important topics if you skip CPSC 340.

Essential Links

- Please bookmark the course webpage:
 - <http://www.cs.ubc.ca/~schmidtm/Courses/340-F19>
 - Contains lecture slides, assignments, optional readings, additional notes.
- You should sign up for Piazza:
 - <http://piazza.com/ubc.ca/winterterm12019/cpsc340/home>.
 - Can be used to ask questions about lectures/assignments/exams.
 - May occasionally be used for course announcements.
- Use Piazza instead of e-mail for questions:
 - I can take a long time to respond e-mails.

Textbooks

- No required textbook.
- I'll post relevant sections out of these books as optional readings:
 - Artificial Intelligence: A Modern Approach (Russell & Norvig).
 - Introduction to Data Mining (Tan et al.).
 - The Elements of Statistical Learning (Hastie et al.).
 - Mining Massive Datasets (Leskovec et al.)
 - Machine Learning: A Probabilistic Perspective (Murphy).
- Most of these are on reserve in the ICICS reading room.
- List of related courses on the webpage, or you can use Google.

TA Cheat Sheet

- Sarah Elhammadi



- Dylan Green



- Nam Hee Kim (Head TA)



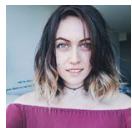
- Frederik Kunstner



- Ke (Mark) Ma



- Lironne Kurzman



- Benjamin Paul-Dubois-Taine



- Michael Przystupa



- Shahriar Shayesteh



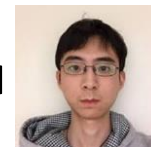
- Betty Shea



- Karl Slakov



- Yihan (Joey) Zhou



Assignments

- There will be 6 Assignments worth 30% of final grade (for 340):
 - Usually a combination of math, programming, and very-short answer.
- Assignment 1 will be on webpage soon, and is due next Friday.
 - Submission instructions will posted on webpage/Piazza.
 - The assignment should give you an idea of expected background.
 - Make sure to submit before the deadline and check your submission.
- Start early, there is a lot there.
 - Don't wait to see you if get off the waiting list to start.

Working in Teams for Assignments

- Assignment 1 must be done individually.
- Assignments 2-6 can optionally be done in pairs.
 - We haven't quite figured out how partners will work this term.
 - Will depend on whether submission is through “handin” or Canvas.
 - We expect you won't need to have the same partner for all assignments.
- All the various permutations of partners are allowed:
 - Partners can be from different sections of 340.
 - One of you can be in 340 and one can be in 532M.
 - Partnering with an auditor is ok.

Late “Class” Policy for Assignments

- Assignments will be due at midnight on the due date.
- If you can't make it, you can use “late classes”:
 - For example, if assignment is due on a Friday:
 - Handing it in Monday is 1 late class.
 - Handing it in Wednesday is 2 late classes.
 - There is no penalty for using “late classes”, but you will get a mark of 0 on an assignment if you:
 - Use more than 2 late classes on the assignment.
 - Use more than 4 late classes across all assignments.
- We'll release solutions to assignments after 2 “late classes”.
 - We'll try to put grades up within 10 days of this.

Assignment Issues

- **No extensions will be considered** beyond the late days.
 - Also, since you can submit more than once, you have no excuse not to submit something preliminary by the deadline.
- Further, due to grouchiness, these issues are a 50% penalty:
 - Missing names or student IDs on assignments.
 - Corrupted .zip submission files or not using a .zip file.
 - Submitting the wrong assignment (year or number).
 - Incorrect assignment names in submission files.
 - Not including answers in the correct location in the .pdf file.

Programming Language: Julia

- 3 most-used languages in these areas: Python, Matlab, and R.
- We will be using Julia which is a free and fast high-level language.
 - See the list of common Julia commands on the course webpage.
- No, you cannot use Matlab/R/TensorFlow/Python/etc.
 - Assignments have prepared code: we won't translate to many languages.
 - TAs shouldn't have to know many languages to grade.

Waiting List and Auditing

- Right now only CS students can register directly.
 - All other students need to **sign up for the waiting list to enroll**.
- We're going to start registering people from the waiting list.
 - Being on the **waiting list is the only way to get registered**:
 - <https://www.cs.ubc.ca/students/undergrad/courses/waitlists>
 - You might be registered without being notified, be sure to check!
 - They might also ask to submit a prereq form, let me know if you have issues.
- Because the room is full, we **may not have seats for auditors**.
 - If there is space, I'll describe (light) auditing requirements then.

Getting Help

- Many students find the assignments long and difficult.
- But there are many **sources of help**:
 - **TA office hours** and **instructor office hours**.
 - Starting in the second week of class.
 - Times will be posted on the course webpage.
 - **Piazza** (for general questions).
 - **Weekly tutorials** (optional).
 - Starting in second week of class.
 - Will go through provided code, review background material, review big concepts, and/or do exercises.
 - **Other students** (ask your neighbor for their e-mail).
 - **The web** (almost all topics are covered in many places).

Midterm and Final

- Midterm worth 20% and a (cumulative) final worth 50%
 - Closed-book.
 - One doubled-sided ‘cheat sheet’ for midterm, two doubled-sided pages for final.
 - No need to pass the final to pass the course (but recommended).
- Midterm is tentatively schedule for 6:30pm October 17th.
 - Let us know if you have a conflict that cannot be resolved.
- I don’t control when the final is, don’t make travel plans before December 18th.
 - If it’s scheduled early, we may restrict the number “late classes” for the last assignment.
- There will be two types of questions:
 - ‘Technical’ questions requiring things like pseudo-code or derivations.
 - Similar to assignment questions, and will only be related topics covered in assignments.
 - ‘Conceptual’ questions testing understanding of key concepts.
 - All lecture slide material except “bonus slides” is fair game here.

Lectures

- All slides will be posted online (before lecture, and final version after).
 - I'll also post Julia versions of Mike's demos on the webpage.
- Please ask questions: you probably have similar questions to others.
 - I may deflect to the next lecture or Piazza for certain questions.
- Be warned that the **course we will move fast** and **cover a lot of topics**:
 - Big ideas will be covered slowly and carefully.
 - But a bunch of other topics won't be covered in a lot of detail.
- Isn't it wrong to have only have shallow knowledge?
 - In this field, it's **better to know many methods** than to know 5 in detail.
 - This is called the “no free lunch” theorem: different problems need different solutions.

Videos from Previous Offering

- Videos of Mike's January 2018 offering of the course:
 - https://www.youtube.com/playlist?list=PLWmXHcz_53Q02ZLeAxigki1JZFfCO6M-b
- You may find these useful:
 - Material is almost identical, but now you can rewind (or fast-forward).
 - Mike is a more experienced teacher than I am.

Bonus Slides

- I will include a lot of “bonus slides”.
 - May mention advanced variations of methods from lecture.
 - May overview big topics that we don’t have time for.
 - May go over technical details that would derail class.
- You are **not expected to learn** the material on these slides.
 - But they’re useful if you want to take 540 or work in this area.
- I’ll use this colour of background on bonus slides.

Code of Conduct

- Do not post offensive or disrespectful content on Piazza.
- If you have a problem or complaint, let me know (maybe we can fix it).
- Do not distribute any course materials without permission.
- Do not record lectures without permission.
- Think about **how/when to ask for help**:
 - Don't ask for help after being stuck for 10 seconds. Make a reasonable effort to solve your problem (check instructions, Piazza, and Google).
 - But **don't wait until the 10th hour of debugging before asking for help**.
 - If you do, the assignments could take all of your time.
- There will be no post-course grade changes based on grade thresholds:
 - 48% will not be rounded to 50%, and 70% will not be rounded to 72%, and so on.

Cheating and Plagiarism

- Read about UBC's policy on "academic misconduct" (cheating):
 - <http://www.calendar.ubc.ca/Vancouver/index.cfm?tree=3,54,111,959>
- When submitting assignments, **acknowledge all sources**:
 - Put "I had help from Sally on this question" on your submission.
 - Put "I got this from another course's answer key" on your submission.
 - Put "I copied this from the Coursera website" on your submission.
 - Otherwise, this is **plagiarism** (course material/textbooks are ok with me).
- **At Canadian schools, this is taken very seriously.**
 - Automatic grade of zero on the assignment.
 - Could receive 0 in course, be expelled from UBC, or have degree revoked.

Course Outline

- Next class discusses “exploratory data analysis”.
- After that, the remaining lectures focus on five topics:
 - 1) Supervised Learning.
 - 2) Unsupervised learning.
 - 3) Linear prediction.
 - 4) Latent-factor models.
 - 5) Deep learning.
- “What is Machine Learning?” (overview of many class topics)

Photo I took in the UK on the way home from the “Optimization and Big Data” workshop:

