

About COMP9318 (2020 t1)

Wei Wang @ CSE, UNSW

February 16, 2020

Introduction

Lecturer-in-charge:

Prof. Wei Wang

School of Computer Science and Engineering

Office: K17 507

E-mail: weiw@cse

Ext: 9385 7162

<http://www.cse.unsw.edu.au/~weiw>

Research Interests:

- Knowledge graph / natural language processing
- AI security
- DB + AI
- High-dimensional data / Similarity query processing

Course Info

- Homepage: <http://www.cse.unsw.edu.au/~cs9318>
- Communications:
 - Main form: Piazza Forum:
<https://piazza.com/class/k6k9ru836577bj>
 - **Email:** `weiw AT cse.unsw.edu.au`:
 - Only for matters that cannot/should not be resolved via piazza.
- Lectures:
 - 1800 – 2000 MON, Rex Vowels Theatre
 - 1800 – 2000 THU, Physics Theatre
- Tutorials: several *online* tutorials + `ipython` notebooks
- Consultations:
 - Use Piazza
 - Weekly by tutors: 1300-1400, K17-508
 - LiC: in lectures or by appointment only.
 - We are considering adding online consultation too.

Assessment (Tentative)

Due to the uncertainty related to the travel ban and university's response, the following is the tentative one.

Overview

- 1 written assignments + 1 programming project + lab
- `lab = np.mean(sorted([lab1, lab2, lab3, lab4, lab5], reverse=True)[:3])`

No late submission allowed for labs. Read the spec of assignment/project to find out late penalty policies.

Default project

- up to 2 students per team
- TBD

Research project

- I have a few topics available. First-come-first-serve. Talk to me for details.

Finally ...

Exam

- If you are ill on the day of the exam, **do not attend** the exam — I will **not** accept medical special consideration claims from people who have already attempted the exam.

Final Mark

- Final mark (tentative)

$$final_mark = 0.15 \cdot ass1 + 0.20 \cdot proj1 + 0.10 \cdot lab + 0.55 \cdot exam$$

- Also requires $exam \geq 40$.

Special Arrangement for Remote Students

For students currently stranded in China

- Please fill out the form: <http://au.mikecrm.com/xatjby0>
- We are organizing additional resources to help.

Warning I

This course has

- Broad coverage
- Heavy workload
- High fail rate $\geq 20\%$
- Plagiarism is not allowed. Make sure you read all *types* of plagiarism, esp. **collusion** in <https://student.unsw.edu.au/plagiarism>.

Specially, we do not accept personal plea or excuses; if you have valid reasons that affect your performance, apply for a UNSW Special Consideration:

<https://student.unsw.edu.au/special-consideration>.

Warning II

Example excuse

- I spent so much time and effort on this course but still failed?
- I did the work by myself and may have shared it with my classmate for discussion.
- If I fail this course, I will [...]. Please.

Lecture Slides

- Contains many materials not found in the text/reference books.

Text Book

- Jensen *et al*, *Multidimensional Databases and Data Warehousing*. (Accessible from a UNSW IP)
- Han *et al*, *Data Mining: Concepts and Techniques*, 1st/2nd edition, Kaufmann Publishers.

Reference Books

- Charu Aggarwal, *Data Mining: The Textbook*, Springer, 2015.
- Tan *et al*, *Introduction to Data Mining*, Addison-Wesley, 2005.
- Leskovec *et al*, *Mining of Massive Datasets* (ver 2.1), Available at <http://infolab.stanford.edu/~ullman/mmds.html>

Software

- Anaconda
- Python 3
- Jupyter notebook
- Python libs such as `numpy`, `pandas`, `matplotlib`, `scikit-learn`, ...

Reading Materials

- Papers from machine learning/data mining conferences/journals, white papers, surveys, etc.
- All available from the course Web page.

Schedule (tentative)

Week	Contents	Assignments
1a	Course overview + Math review	
1b	Math review + Data warehousing and OLAP	
2a	Data warehousing and OLAP	lab1
2b	Data warehousing and OLAP	
3a	Data Preprocessing	
3b	Data Preprocessing	
4a	Classification	lab2
4b	Classification	
5a	Classification	assignment/project
5b	Classification	
6a	Classification	lab3
6b	Clustering	assignment/project
7a	Clustering	
7b	Clustering	
8a	Clustering	lab4
8b	Association Rule Mining	
9a	Holiday	
9b	Association Rule Mining	
10a	Review	lab5

Course Objective and Requirements

Objectives:

- Cover practically useful data mining/machine learning algorithms and concepts
- Foster **deeper** understanding of maths, models, and algorithms
- Gain hands-on experience with solving real problems

Requirements:

- You need to have a solid background in Maths (Linear Algebra, Calculus, Probability & Statistics) and programming (mainly python).
- **Understand** (not memorize) concepts/equations/algorithms.
 - Ask *why*.
 - Describe it in your own language to a layman.

Feedback welcome (throughout the course).

Example

Example

John got a positive result for the α test, and the probability that patients with the deadly β disease having a positive α test result is 99%. Should John be worried about having the β disease?

Example

Example

John got a positive result for the α test, and the probability that patients with the deadly β disease having a positive α test result is 99%. Should John be worried about having the β disease?

$$P(\beta \mid \alpha) = \frac{P(\alpha \mid \beta)P(\beta)}{P(\alpha)} = 0.99 \frac{P(\beta)}{P(\alpha)}$$

Example

Example

John got a positive result for the α test, and the probability that patients with the deadly β disease having a positive α test result is 99%. Should John be worried about having the β disease?

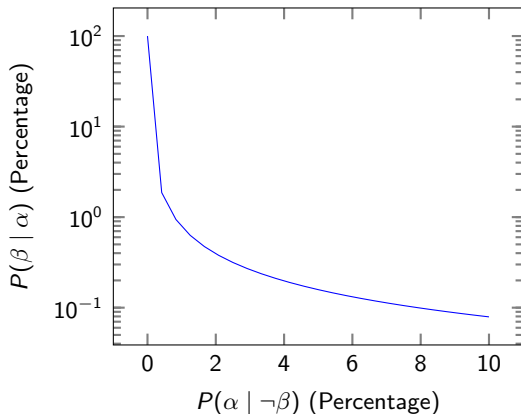
$$P(\beta \mid \alpha) = \frac{P(\alpha \mid \beta)P(\beta)}{P(\alpha)} = 0.99 \frac{P(\beta)}{P(\alpha)}$$

$$P(\beta \mid \alpha) = \frac{P(\alpha \mid \beta)P(\beta)}{P(\alpha \mid \beta)P(\beta) + P(\alpha \mid \neg\beta)P(\neg\beta)}$$

Example

Exercise

Exercise: plot the function $P(\beta \mid \alpha)$ with respect to $P(\alpha \mid \neg\beta)$ given $P(\beta) = \frac{8}{100,000}$.



Example

Example

John got a positive result for the α test.

All patients with the deadly β disease have a positive α test result.

Does John have the β disease?

Example

Example

John got a positive result for the α test.

All patients with the deadly β disease have a positive α test result.

Does John have the β disease?

$\beta \rightarrow \alpha$ is true does not imply that the converse, $\alpha \rightarrow \beta$ is true.

For those new to the computing environment at CSE, UNSW

- Use Linux/command line.
 - Project marked on linux servers
 - You need to be able to upload, run, and test your program under linux.
- Assignment/Project submission
 - Give to submit. **Watch out** for possible error messages.
 - Classrun. Check your submission, marks, etc. Read <https://wiki.cse.unsw.edu.au/give/Classrun>
 - Common errors:
 - File corrupt (during SFTP?), not in the correct format.
 - Submission not accepted by the system (wrong filename? too large? ...).
- Lab submission: our home-made Web submission system.

Other Specialised Courses

Other specialised courses in the **Database** or **Data Science** stream:

- COMP9319: Advanced algorithms on compression, text/XML databases, etc.
- COMP9313: Big data systems (hadoop, spark, etc)
- COMP6714: Information retrieval, Natural language processing, Search engines.

Other machine learning courses:

- COMP9417: Machine Learning and Data Mining
- COMP9444: Neural Networks and Deep Learning
- COMP9418: Advanced Machine Learning

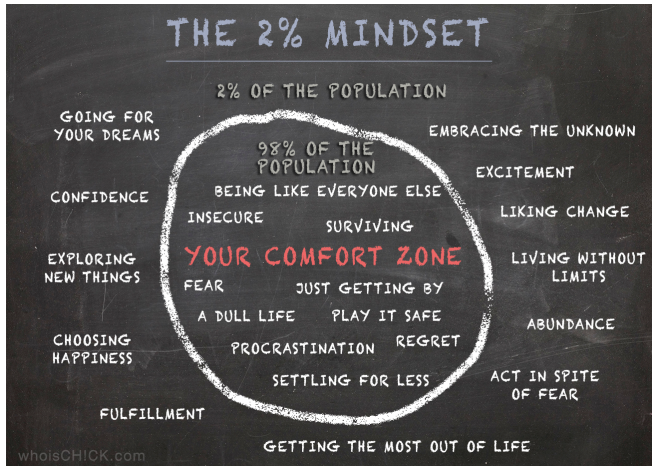
Things to ponder:

- The long-term impact of the latest development in AI/DS/Hardware.
- What do you want out of this course?

Requirement:

- **Plan ahead** for the course.
- Learning happens outside your **comfortable zone**.
- **Review** teaching materials after the lecture.
- Use the Jupyter **notebooks**.

Make Errors and Learning Sth. New



Source:

<http://combiboilersleeds.com/images/comfort-zone/comfort-zone-0.jpg>