

Data Mining

INFS 4203/7203

Miao Xu

miao.xu@uq.edu.au

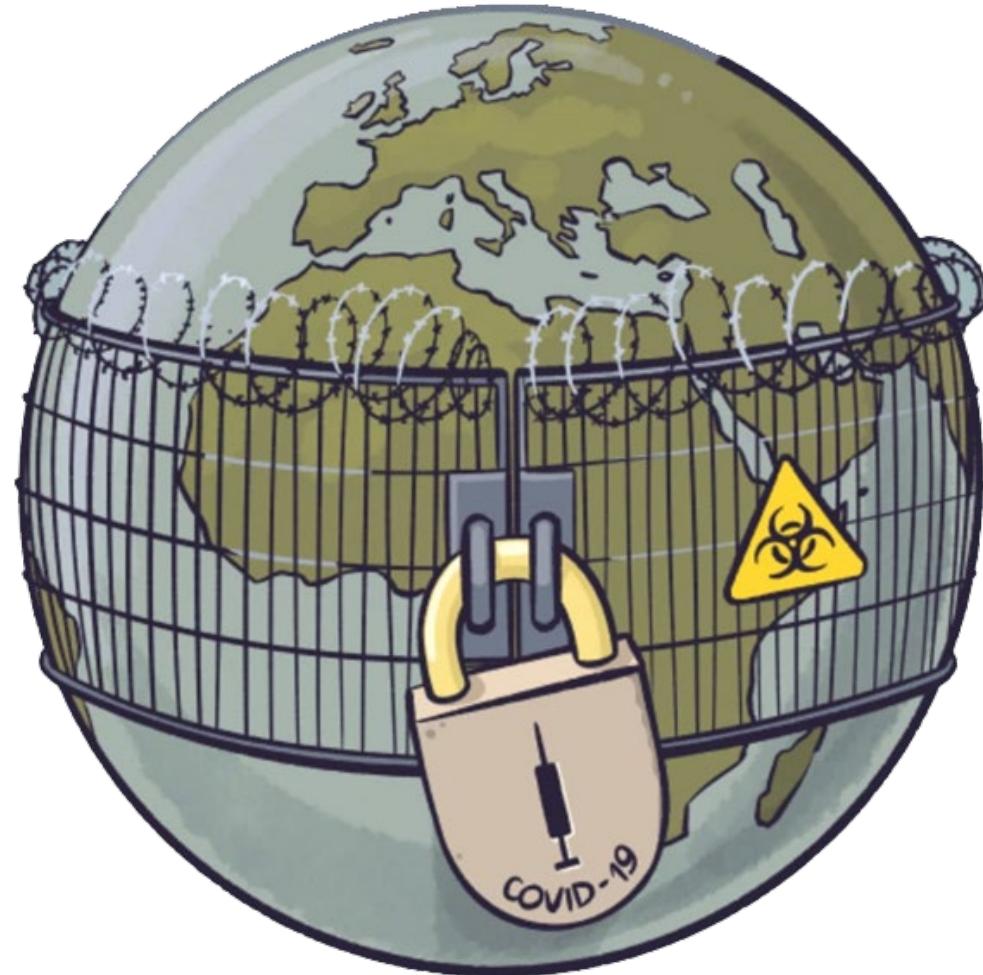
The University of Queensland, 2020 Semester 2



**WELCOME
BACK TO
the
Classroom**

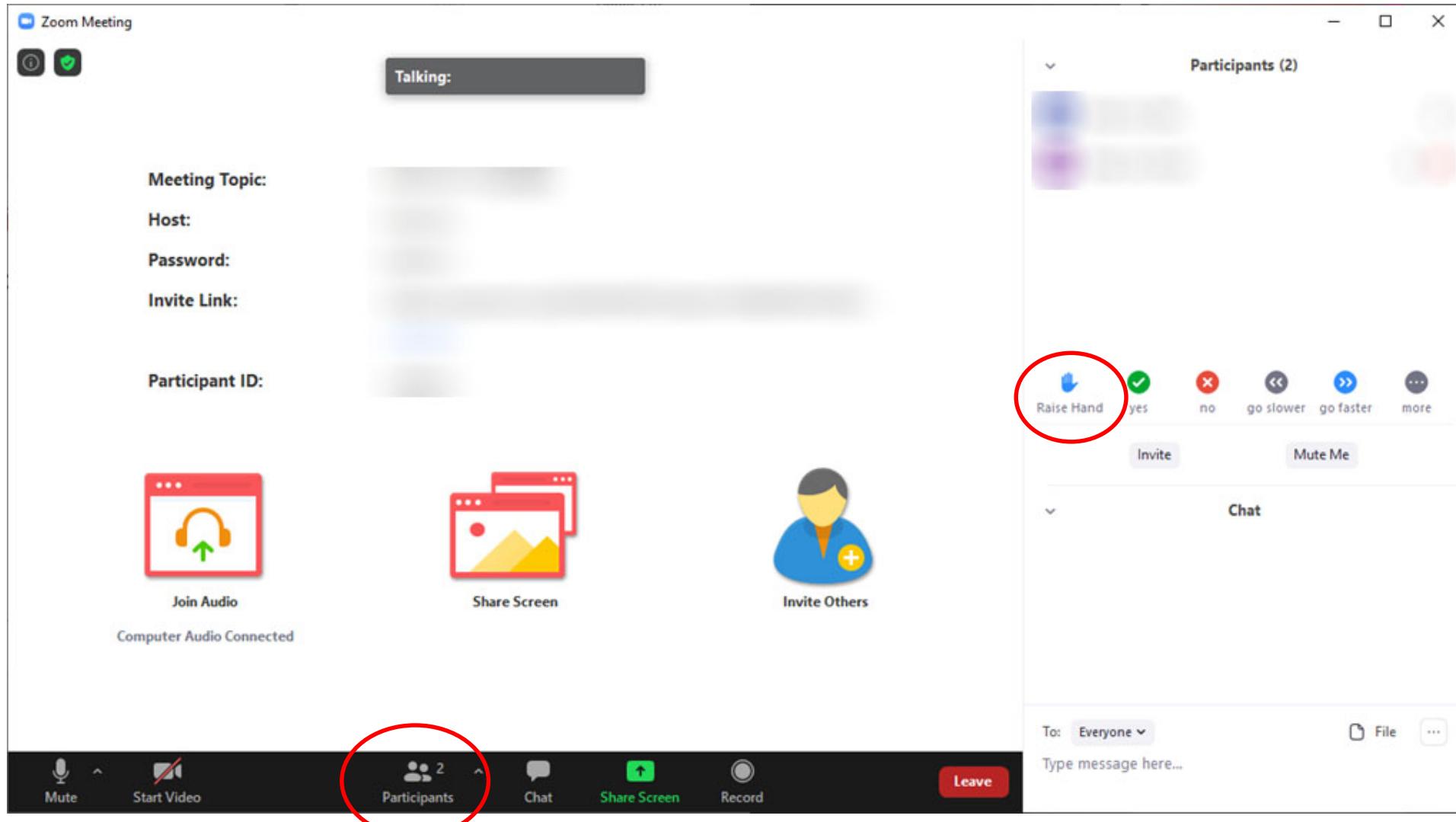
A mixture of external mode and flexible mode

- Lectures:
 - Interactively online
- Contact:
 - Interactively online
- Tutorials:
 - Online and on campus

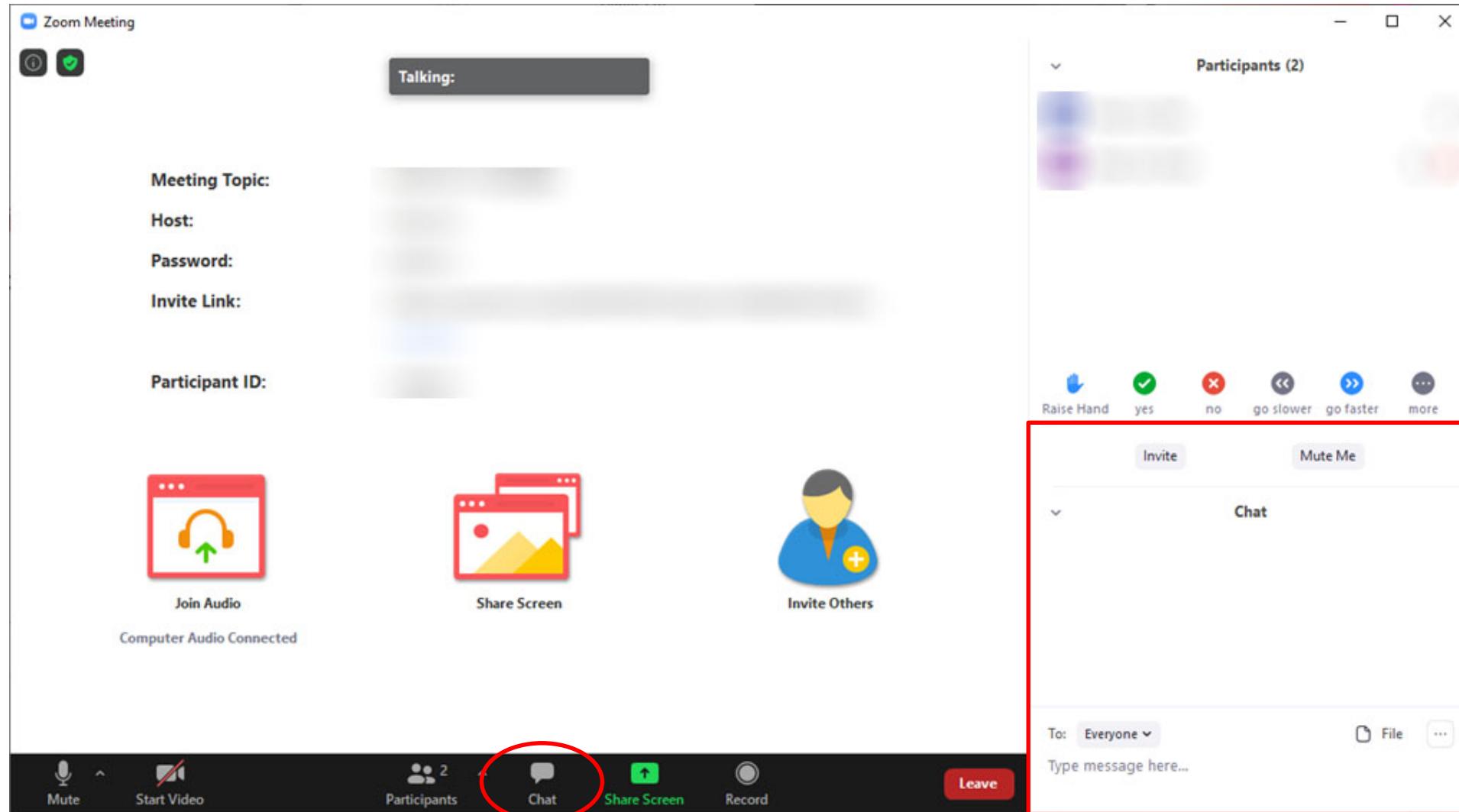


Stay healthy and enjoy the course!

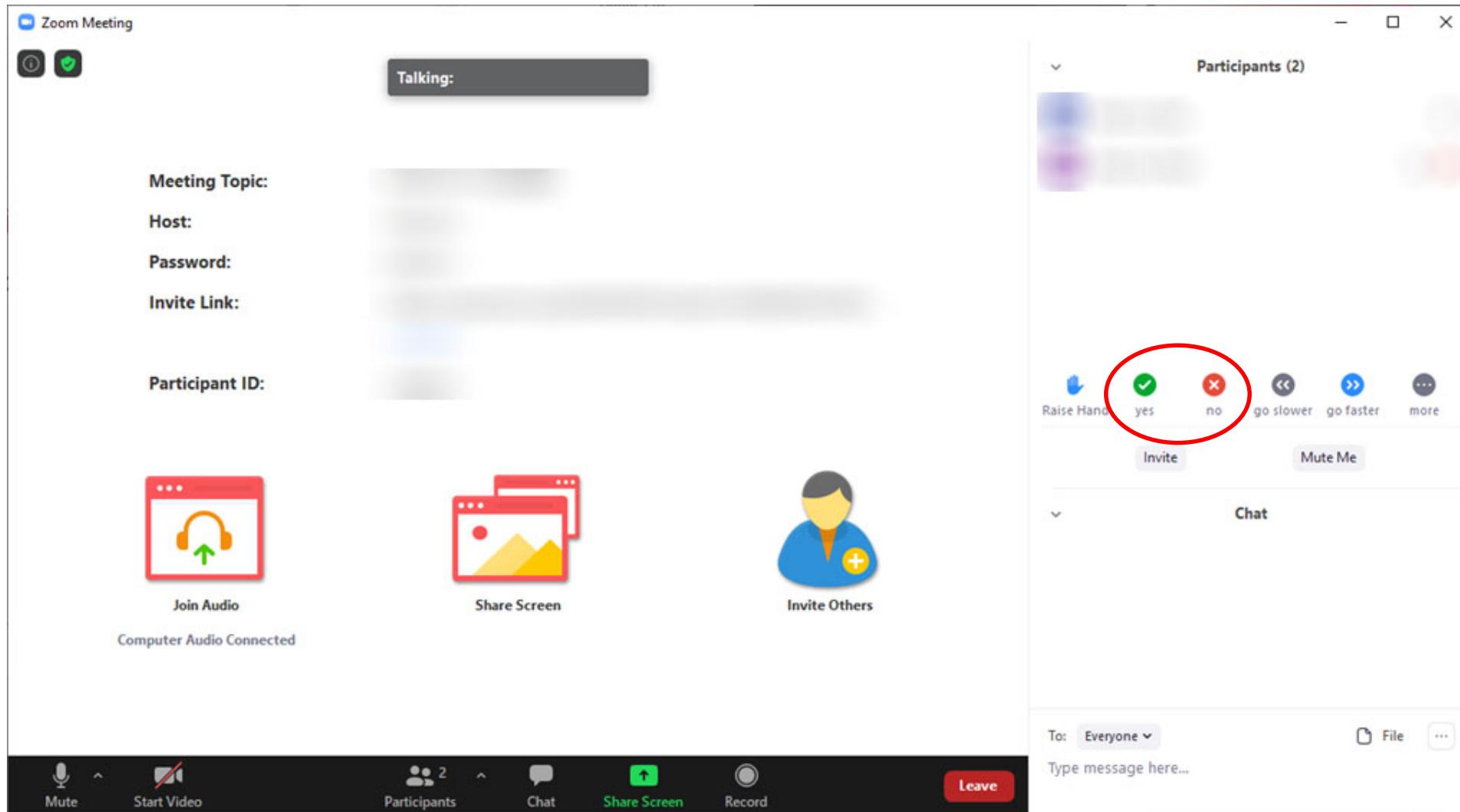
Zoom tips: asking a question by “Raise Hand”



Zoom tips: asking a question by “Chat”



Zoom tips: answering a yes/no question



This lecture:

- Basic information for INFS4203/7203
- Introduction to Data Mining

Basic Information for INFS4203/7203

Outline

- Who are your fellows?
- What you are supposed to learn?
- What you are supposed to do?
- Where you can ask for help?

Outline

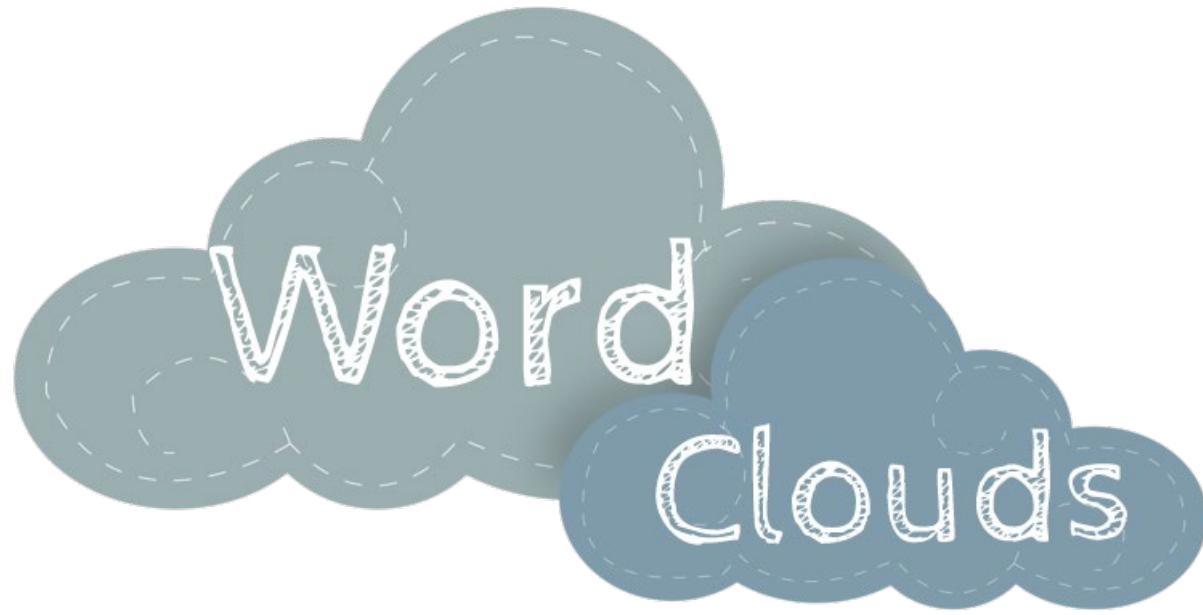
- Who are your fellows?
- What you are supposed to learn?
- What you are supposed to do?
- Where you can ask for help?

Quick facts

- Bachelor/Master/Grad Cert./Grad Dip./Study Abroad in
 - Data Science/ Computer Science/ Information Technology/ Bioinformatics/ Business/ Commerce/ Engineering Science/ Engineering/Geographic Information Science/ Applied Econometrics/ Science/ Mathematics...

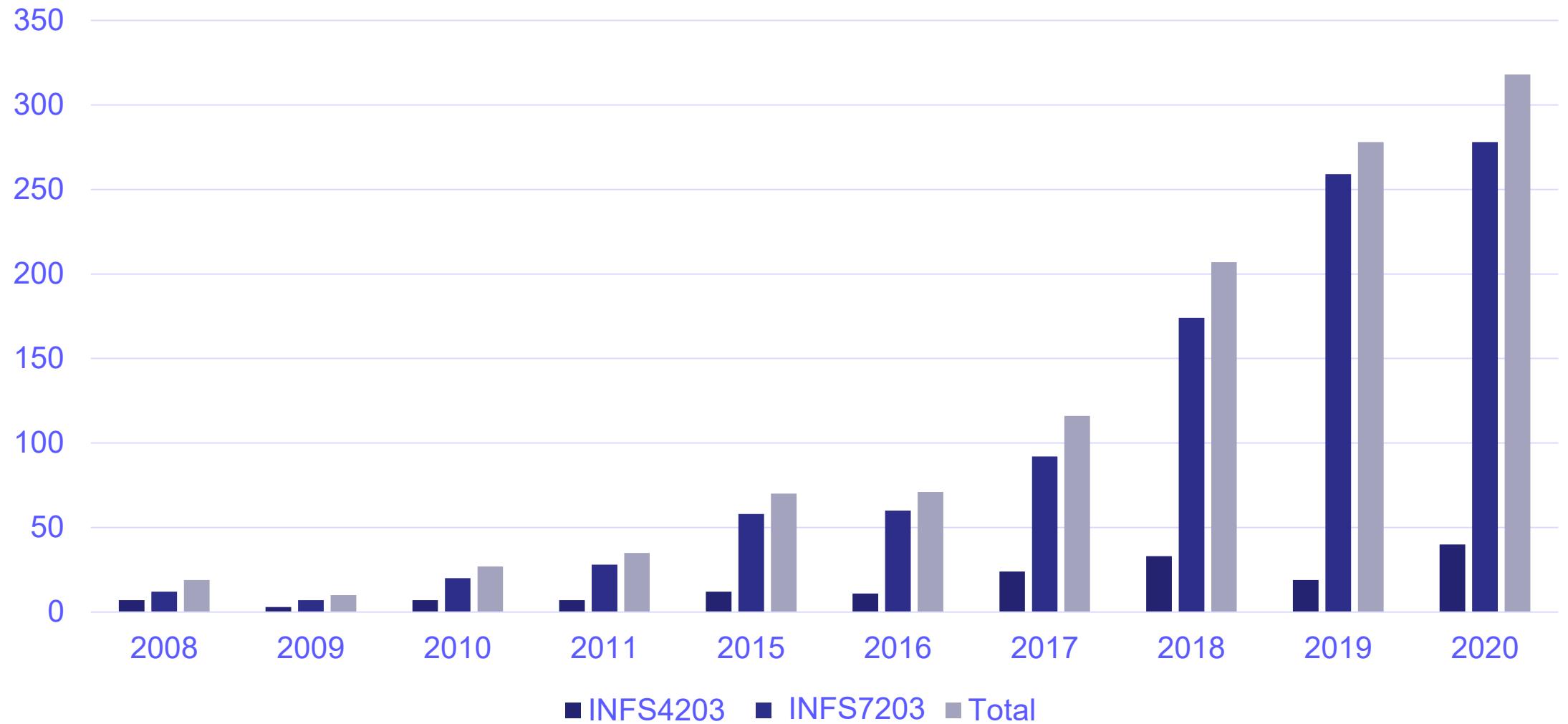
Course	Mode	Enrolment
INFS4203 (Undergraduate)	Flexible	29
	External	13
INFS7203 (Postgraduate)	Flexible	193
	External	119

Activity: what is your major



apps.elearning.uq.edu.au/wordcloud/75753

Enrolment



Activity: group discussion

- Task 1: introduce yourself
- Task 2: discuss why you take the course
- Task 3: what do you expect from this course
- Task 4: one representative records Task 2 and 3 in
<https://padletuq.padlet.org/miaoxu/m7lfg4880nf26jis>

Outline

- Who are your fellows?
- What you are supposed to learn?
- What you are supposed to do?
- Where you can ask for help?

INFS courses

- INFS1200/7900 Intro to Information Systems
- **INFS2200/7903 Relational Database Systems**
- INFS3200/7907 Advanced Database Systems
- INFS3202/7202 Web Information Systems
- INFS3208/7208 Cloud Computing
- **INFS4203/7203 Data Mining**
- INFS4205 Advanced Technology for High Dimensional Data
- INFS7410 Information Retrieval
- INFS7450 Social Media Analytics
- INFS7901 Database Principles

Pre-Requisites

INFS1200/7900 vs. INFS2200/7903

- Recall what you have learned?
- What are the difference between these two courses?

INFS1200/7900

Intro to Information Systems

INFS2200/7903

Relational Database Systems

INFS1200/7900 vs. INFS2200/7903



INFS1200/7900

Intro to Information Systems

INFS2200/7903

Relational Database Systems

INFS4200/7203 Data mining

“Data mining is the analysis of (often large) observational data sets to find unsuspected relationships and to summarize the data in novel ways that are both understandable and useful to the data owner.”

[D. Hand et al., Principles of Data Mining]



Data mining vs. machine learning

Data mining vs. machine learning



- Database originated
- Data-centred



- AI originated
- Learning-centred

Yes/No questions

- Is data mining the same as machine learning?
- Is data mining totally different from machine learning?
- Is data mining the result of the evolution of database technologies?
- Is data mining the result of the evolution of machine learning research?

Outline

- Who are your fellows?
- What you are supposed to learn?
- **What you are supposed to do?**
- Where you can ask for help?

Learning activities

	Lectures	Tutorials	Contact
When	Fri. 14:00-16:00	8 sessions on campus 45 sessions online	Mon. 17:00-20:00
Where	Zoom	14-115/Zoom	Zoom

- Tutorials:
 - Start in Week 3
 - Problem Solving
 - Python Language
- <https://www.learnpython.org/>
- Contact:
 - Start in Week 2
 - Consultation
 - Not Compulsory
 - Join/Leave Freely

Staying Healthy @ UQ



Stay home if you are unwell



Cover your mouth and nose when you sneeze or cough



Avoid touching your face



Wash your hands thoroughly



Don't share personal items



Clean surfaces



Maintain space between each other



Put used tissues in the bin



Call your General Practitioner (doctor) or UQ Health Care and explain your symptoms

CRICOS 00025B



Need the facts? about.uq.edu.au/coronavirus-advice-uq-community

Course website

- <https://learn.uq.edu.au/>

INFS4203/7203

Data Mining

[INFS4203/7203] Data
Mining (St Lucia &
external). Semester 2,
2020, Flexible Delivery

(INFS4203/7203)

Announcements

Course Profile (eCP)

Course Staff

Course Help

Learning Resources

Assessment

Discussion Board

My Grades

Library Links

Announcements

Learning Resources for Orientation Week

Posted on: Thursday, 30 July 2020 13:24:05 PM AEST

- Please check the website at least once a week

Please check the learning resources for orientation week. Although lectures don't start until Friday Week 1, there are still things

Course Link/Learning Resources/Orientation Week: Welcome

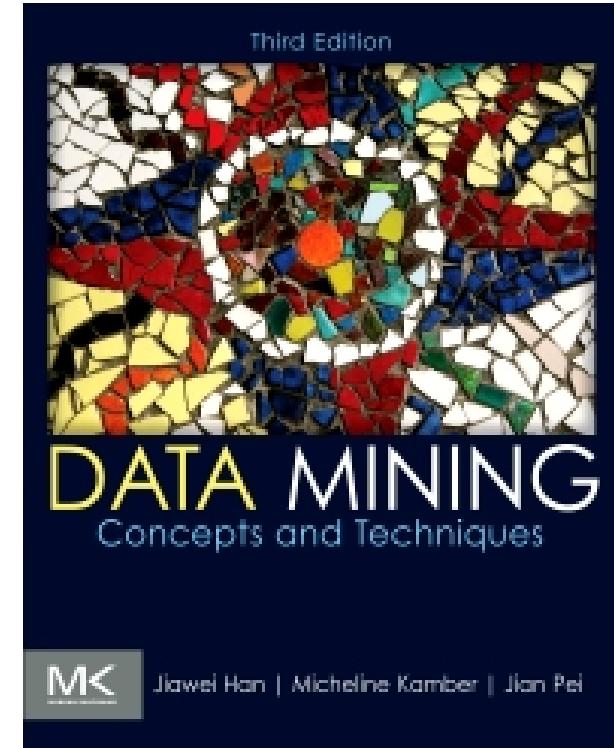
A Warm Welcome to the INFS4203/7203 Data Mining Blackboard Site

Learning materials

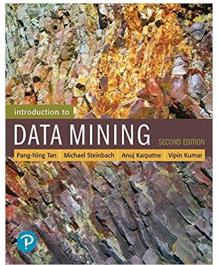
- Lecture notes/records
- Tutorial notes/records
- Weekly summarized frequently asked questions

Recommended reference book

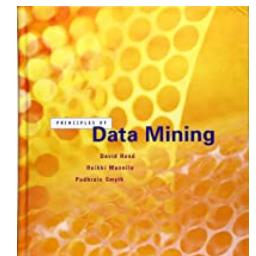
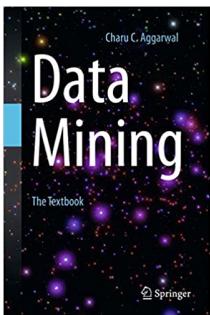
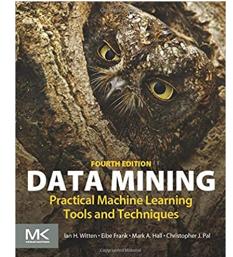
- Jiawei Han, Micheline Kamber, Jian Pei. *Data mining concepts and techniques 3rd ed.* Morgan Kaufmann, 2012.
- Additional chapters from 2nd ed
 - Mining stream and time-series data
 - Mining graph data
 - Mining text and web data



Other data mining books (for your interest)



- Pang-Ning Tan, Michael Steinbach, Anuj Karpatne and Vipin Kumar. *Introduction to Data Mining*, 2nd ed. Pearson, 2018. (the library link is for a previous version)
- Ian H. Witten, Eibe Frank, Mark A. Hall and Christopher J. Pal. *Data Mining: Practical Machine Learning Tools and Techniques*. Morgan Kaufmann, 4th edition, 2016.
- Charu C. Aggarwal. *Data Mining: The Textbook*. Springer, 2015.
- David Hand, Heikki Mannila and Padhraic Smyth. *Principles of Data Mining*. The MIT Press, 2001.



Individual assessment

Assessment	When	Percentage
Assignment1	Due by Aug. 27 th	4%
Assignment2	Due by Sep. 10 th	6%
Assignment3	Due by Oct. 8 th	6%
Assignment4	Due by Oct. 22 nd	6%
Assignment5	Due by Oct. 29 th	3%
Mid-term	Sep. 18 th in class	25% Online Close-Book
Final-term		50% Invigilated Online Close-Book



Uh, Mom?
Remember that paper
you wrote for me about nuclear
reactors? What does it mean
when the big, red "meltdown"
light is flashing?

AI Awareness Poster, Brigham Young University
https://live-academicintegrity.pantheonsite.io/wp-content/uploads/2017/12/462c19_1dc7602ba4ec4f4586d803e7c43904cc_mv2-1.gif

Outline

- Who are your fellows?
- What you are supposed to learn?
- What you are supposed to do?
- Where you can ask for help?

Course staff

- Course coordinator: Dr. Miao Xu
 - Email: miao.xu@uq.edu.au
 - Office: General Purpose South Building (78) 633, St Lucia
- Tutors (see Blackboard for their contacts)
 - Mr. Rocky Chen
 - Ms. Sivangi Mund
 - Ms. Skye Sun
 - Mr. Zijian Wang
 - Mr. Ziwei Wang

Ask for help

- Ask in class or immediately before/after class
- Post it on Blackboard discussion board Piazza
- Bring it to “Contact” time (Mon. 5pm-8pm)
 - Email tutors
 - Email course coordinator
 - Visit the office: by appointment

You are encouraged to visit the Piazza to share your thoughts and help your fellows!



Lecture 1: Introduction to Data Mining

Example 1: Mining transaction data

Bread
Coke
Milk



Beer
Bread
Milk



Beer
Coke
Diapers
Milk



Beer
Bread
Diapers
Milk



Diapers
Milk



Anything interesting?

- How many items are there?
- Do they appear frequently?
- Do they appear together?
- Any patterns?
- What can be the implicit reasons?

Example 1: Mining transaction data

Bread
Coke
Milk



Beer
Bread
Milk



Beer
Coke
Diapers
Milk



Beer
Bread
Diapers
Milk



Diapers
Milk



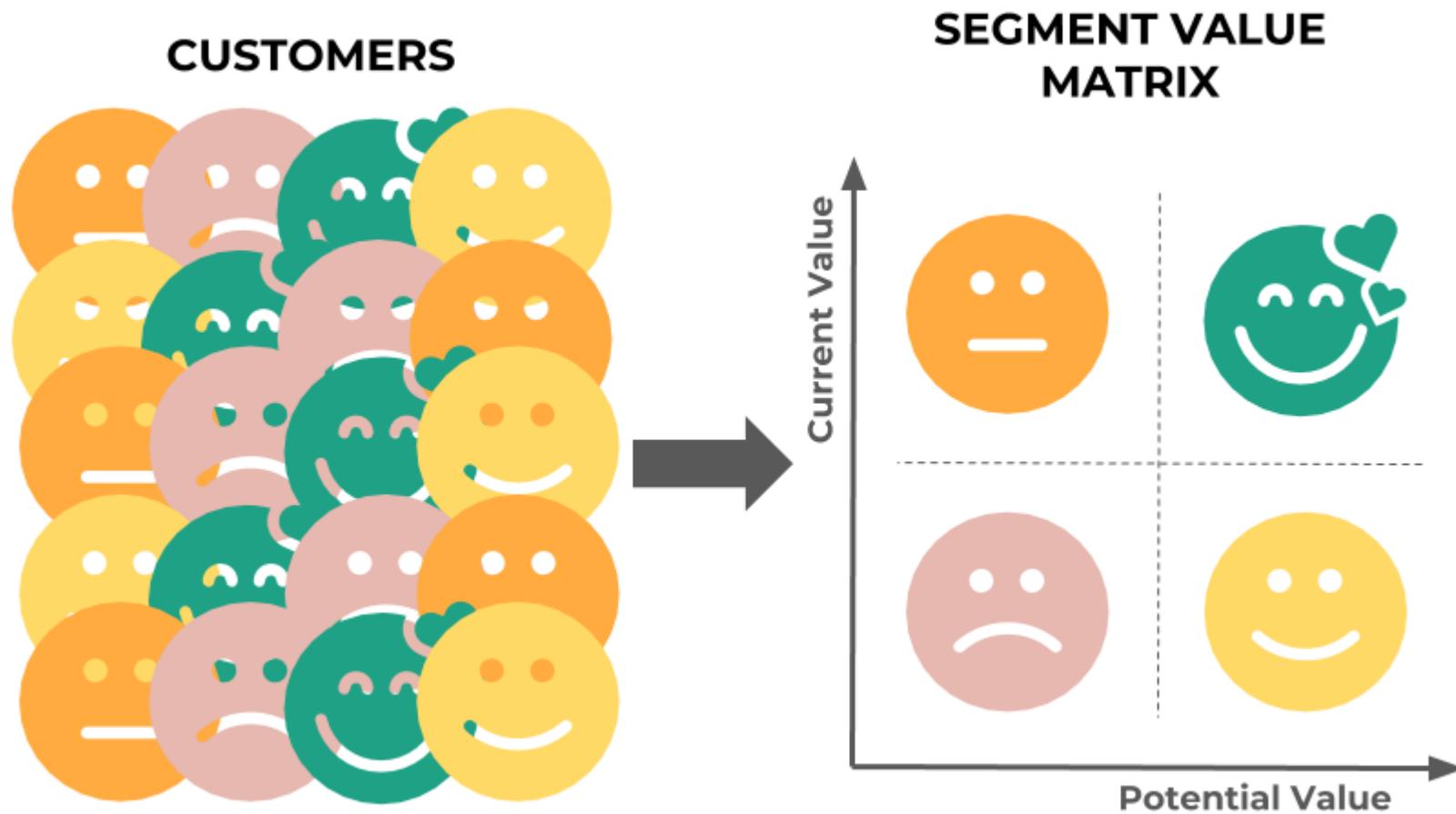
Anything interesting?

- ✓ Bread → Milk 100%
- ✓ Diapers → Milk 100%
- ✓ Diapers → Beer 66%

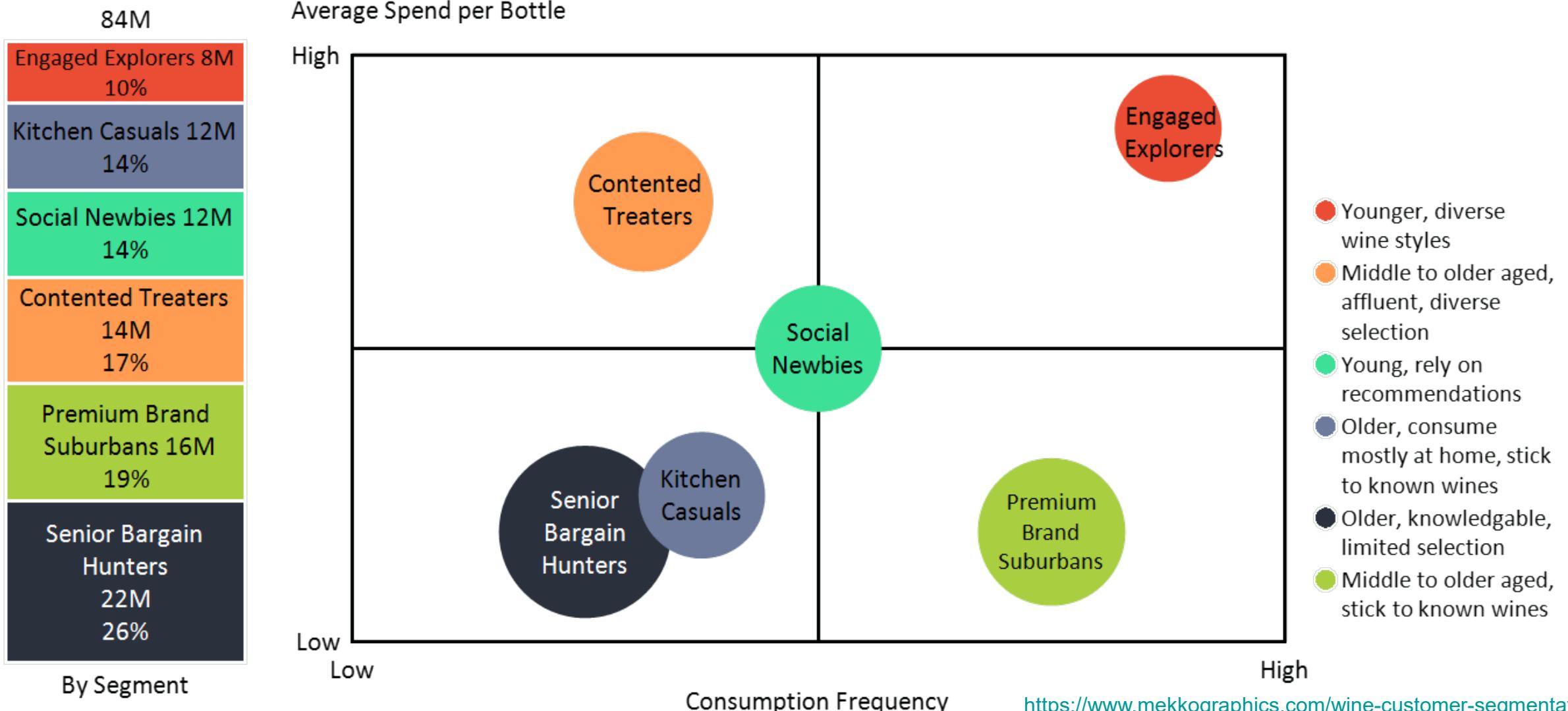




Example 2: Mining customer data



Example 2: Mining customer data



Example 3.1: Mining video surveillance data



[https://cdn.vox-cdn.com/thumbor/9Hirq_5MsKwoDhz1biJWMJ1BbQY=/0x0:601x371/620x413/filters:focal\(253x138:349x234\):gifv\(\):no_upscale\(\)/cdn.vox-cdn.com/uploads/chorus_image/image/60183483/stop_thief_ai_cctv_automated_surveillance.0.gif](https://cdn.vox-cdn.com/thumbor/9Hirq_5MsKwoDhz1biJWMJ1BbQY=/0x0:601x371/620x413/filters:focal(253x138:349x234):gifv():no_upscale()/cdn.vox-cdn.com/uploads/chorus_image/image/60183483/stop_thief_ai_cctv_automated_surveillance.0.gif)

Example 3.2: Credit card fraud detection

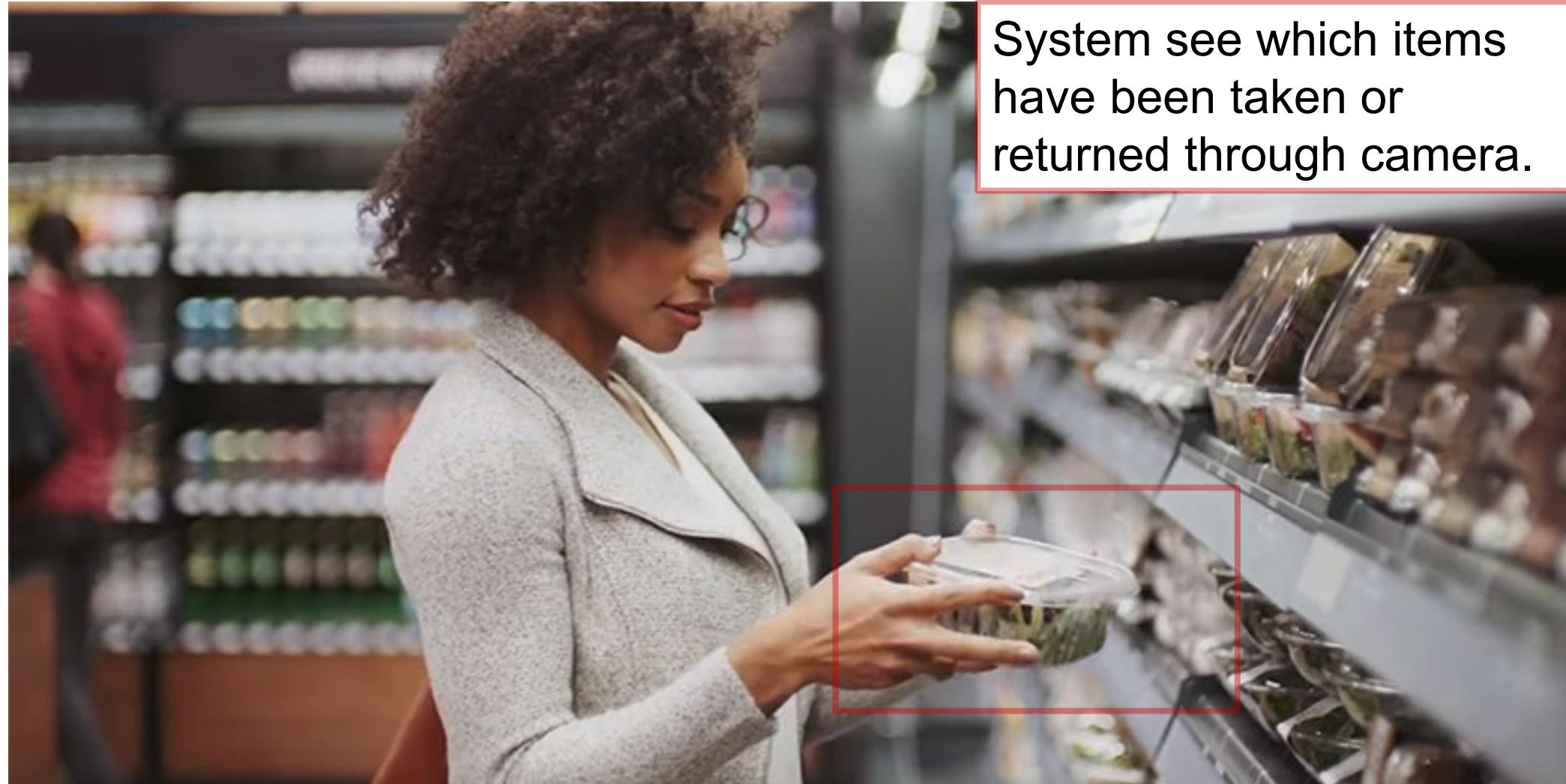


<https://logentries.com/product/anomaly-detection/>

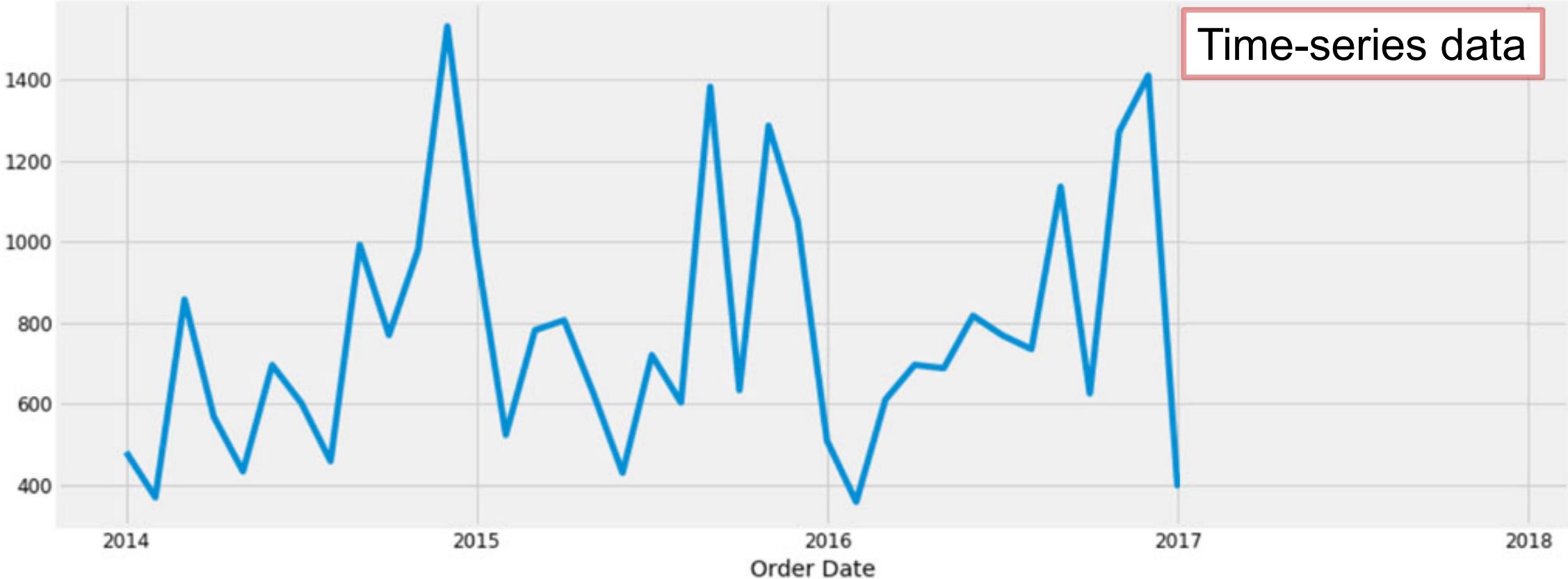
Example 4: Amazon Go



Example 4: Amazon Go



Example 5.1: sales prediction



Example 5.1: sales prediction



Example 5.2: an online store

Review Snapshot by PowerReviews



4.1

18 reviews Write A Review

Ratings Distribution

5 Stars		9
4 Stars		4
3 Stars		3
2 Stars		2
1 Star		0

Pros

- 12 Smells/Tastes Great
- 11 Soothing
- 10 Effective
- 7 Healing
- 1 Long Lasting

Cons

- 9 Not Long-Lasting
- 2 Bad Taste
- 2 Greasy
- 1 Bad Smell
- 1 Ineffective

Describe Yourself

9 Brand Buyer

6 Budget Buyer

Best Uses

14 Treat Chapped Lips

9 Daily Use

2 Sun Protection

1 Prevent Wind Burn

Most Liked Positive Review



Pleasantly surprised

I have used the same brand of lip balm for over 5 years and I'm glad I was able to try something new! I loved the smell and it helped my chapped lips. This tube is a little bigger than your typical lip balm, but I found that it was much easier to find in my backpack. I hope they add more flavors (ma...)

[Read complete review](#)

Most Liked Negative Review



Not the best lip balm

I was looking forward to trying this lip balm, but I was a little disappointed. The color of the tube and the tube itself are cute and different, but the actual lip balm was less than ideal. It didn't last very long on my lips and I didn't really see a change after using it for a couple days. I also...

[Read complete review](#)

Text/web data

<https://www.powerreviews.com/blog/how-to-get-more-product-reviews/>

After learning the data mining course, you will have a fundamental knowledge of how these are achieved.

Outline

- Why data mining?
- What is data mining?
- What can data mining do?
- Who are using data mining?

Outline

- Why data mining?
- What is data mining?
- What can data mining do?
- Who are using data mining?

Why data mining

- How will you describe the haystack?
- How will you describe the golden needle?
- What is the meaning of fire on the top?
- Is this task hard?
- Why is it hard?



**“Search for a golden needle in a haystack.
So much hay and so little time”**

Why data mining

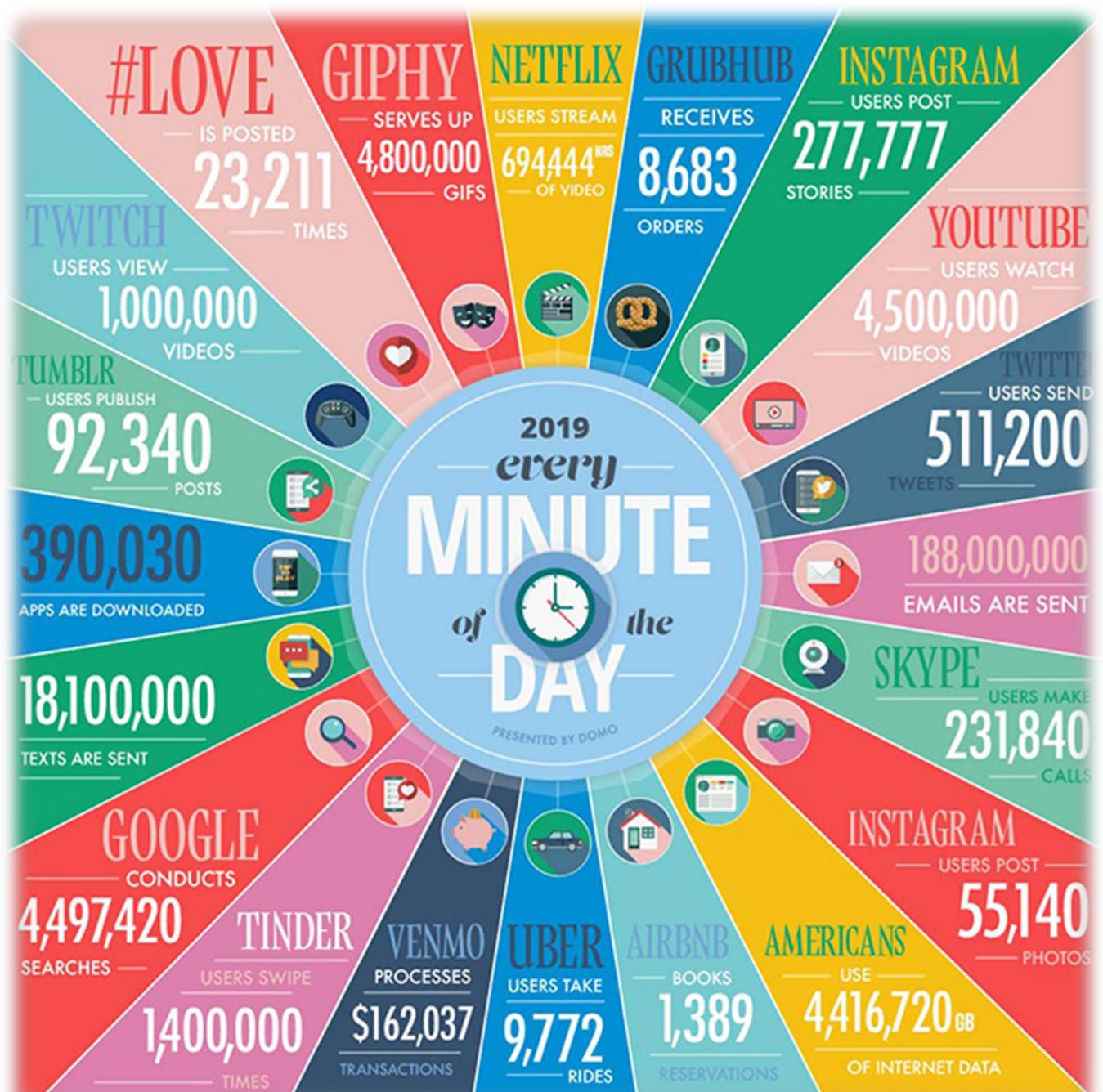
- ✓ Volume
- ✓ Value
- ✓ Velocity
- ✓ Variety



**“Mining golden needle in a haystack.
So much hay and so little time”**

Data on the web

- ✓ Volume
- ✓ Value
- ✓ Velocity
- ✓ Variety



Data in healthcare

Data gives us the potential to improve quality of health care meanwhile reducing cost.

1PB = 1024 TB

The Healthcare Data Explosion

2012
500 petabytes

Worldwide healthcare data is expected to grow to **50 times** the current total

2020
25,000 petabytes

<https://www.scoop.it/topic/mobile-health-how-mobile-phones-support-health-care>

Cool Facts About Big Data in Healthcare



In healthcare, Big Data holds enormous promise, and we're already seeing its impact in areas such as precision medicine, biopharmaceutical R&D productivity, mobile health, telemedicine and more.



Big Data analytics could reduce pharmaceutical R&D costs as much as **\$70 billion**

Potential \$300 to \$450 billion in reduced healthcare spending



...through widespread application of Big Data



10 Years

The amount of time it took to decode the human genome

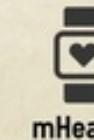


It can now be done in...

..about 1 week



The biggest opportunity for Big Data in Healthcare



mHealth



EMRs



Internet of things

Outline

- Why data mining?
- **What is data mining?**
- What can data mining do?
- Who are using data mining?

What is data mining

“Data mining is the analysis of (often large) observational data sets to find **unsuspected** relationships and to summarize the data in **novel** ways that are both **understandable** and **useful** to the data owner.”

[D. Hand et al. , Principles of Data Mining]

Associates the words with the objects

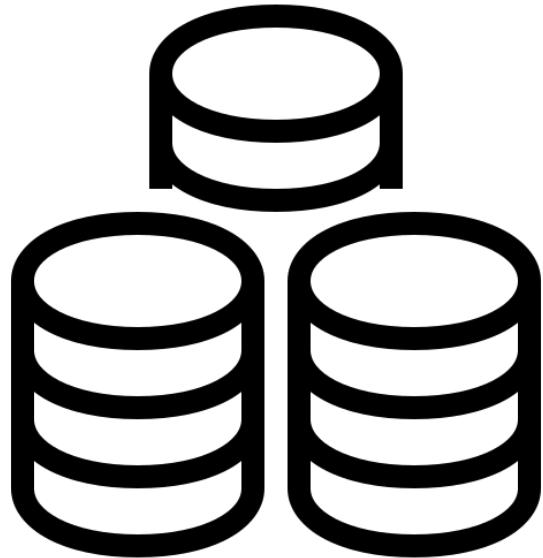
Large

Unsuspected

Novel

Understandable

Useful



Data



Mining
Technique



Knowledge

Associates the words with the objects



Mining
Technique

**Efficient in both
time and space**

Novel
Useful
Understandable
Unsuspected



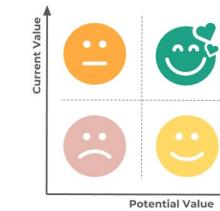
Knowledge

Outline

- Why data mining?
- What is data mining?
- What can data mining do?
- Who are using data mining?

What can data mining do

- Association rule mining (Week 3)
 - How to find highly correlated items out
- Clustering (Week 4, 5)
 - How to segment data with similar traits
- Anomaly detection (Week 6)
 - How to detect outliers in a dataset
- Classification (Week 8, 9, 10)
 - How to construct the mapping from the input space to the nominal output space



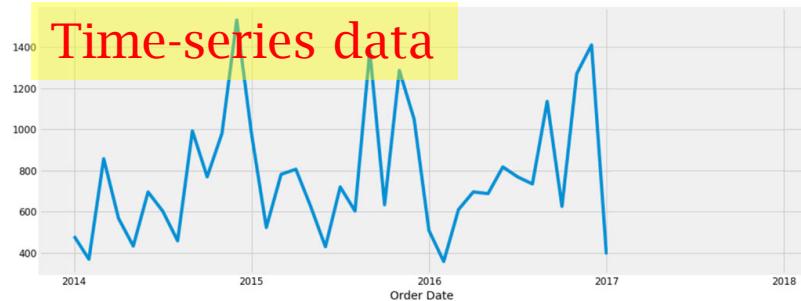
Flat data

“Flat” data: vectors and matrix

A WEEK-BY-WEEK LOOK AT COVID-19 IN ENGLAND							
	MAY 21	JUNE 4	JUNE 18	JUNE 25	JULY 2	JULY 9	JULY 17
% OF POPULATION INFECTED	0.25%	0.10%	0.06%	0.09%	0.04%	0.03%	0.04%
TOTAL CURRENTLY INFECTED	137,000	53,000	33,000	51,000	25,000	14,000	24,000
NEW CASES PER WEEK	61,000	39,000	26,900	22,000	25,000	11,900	11,900
NEW CASES PER DAY	8,714	5,500	3,800	3,142	3,571	1,700	1,700
R RATE	0.7 - 1.0	0.7 - 0.9	0.7 - 0.9	0.7 - 0.9	0.7 - 0.9	0.7 - 0.9	0.7 - 0.9
HOSPITAL ADMISSIONS	697	505	387	318	394	319	179
DEATHS ANNOUNCED	338	176	184	154	176	126	66
NEW POSITIVE TESTS	2,615	1,805	1,115	653	829	630	642

What can data mining do - continue

- Mining complex data (Week 11, 12)



Web data

Review Snapshot by PowerReviews

4.1 18 reviews Write A Review

Ratings Distribution	Pros	Cons
5 Stars	9	12 Smells/Tastes Great
4 Stars	4	11 Soothing
3 Stars	3	10 Effective
2 Stars	2	7 Healing
1 Star	0	1 Long Lasting

Describe Yourself: 9 Brand Buyer 6 Budget Buyer

Best Uses: 14 Treat Chapped Lips 9 Daily Use 2 Sun Protection 1 Prevent Wind Burn

Most Liked Positive Review: 5 Pleasantly surprised

I have used the same brand of lip balm for over 5 years and I'm glad I was able to try something new! I loved the smell and it helped my chapped lips. This tube is a little bigger than your typical lip balm, but I found that it was much easier to find in my backpack. I hope they add more flavors (m...)

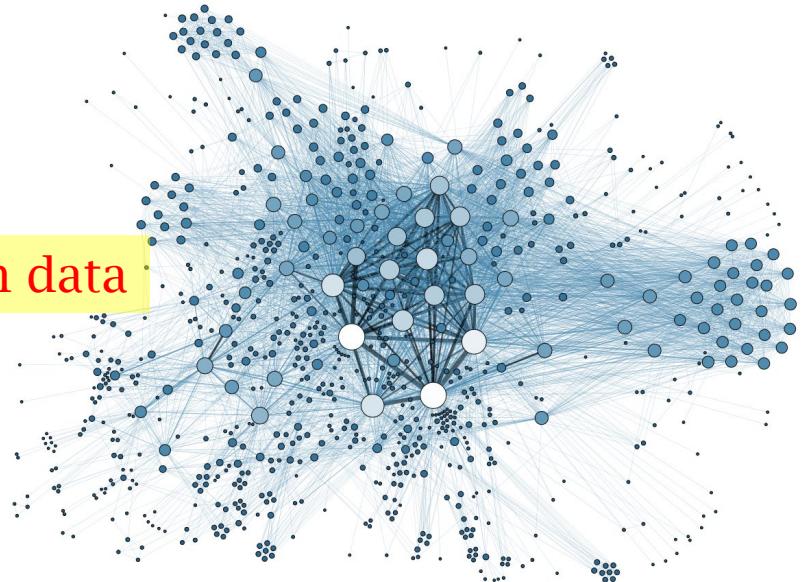
Read complete review

Most Liked Negative Review: 3 Not the best lip balm

I was looking forward to trying this lip balm, but I was a little disappointed. The color of the tube and the tube itself are cute and different, but the actual lip balm was less than ideal. It didn't last very long on my lips and I didn't really see a change after using it for a couple days. I also...

Read complete review

Graph data



Outline

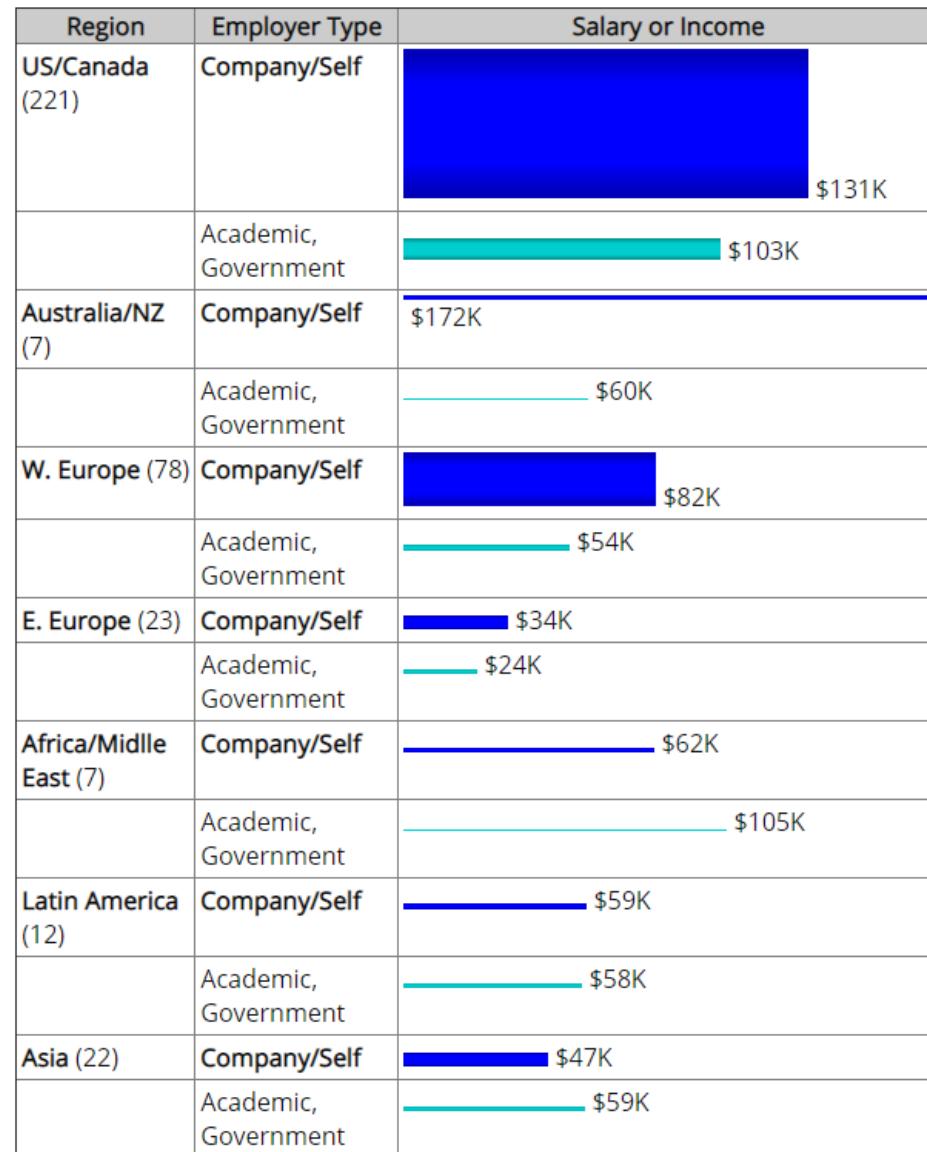
- Why data mining?
- What is data mining?
- What can data mining do?
- Who are using data mining?

Who need data mining

- Previous Sponsors and Supporters for KDD



Annual salary of data miners



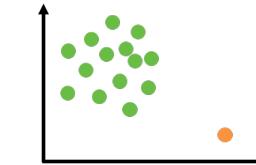
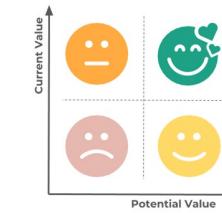
<https://www.kdnuggets.com/polls/2015/salary-analytics-data-science-data-mining.html>

Outline

- Why data mining?
- What is data mining?
- What can data mining do?
- Who are using data mining?

Next lecture

- Association rule mining (Week 3)
 - How to find highly correlated items out
- Clustering (Week 4, 5)
 - How to segment data with similar traits
- Anomaly detection (Week 6)
 - How to detect outliers in a dataset
- Classification (Week 8, 9, 10)
 - How to construct the mapping from the input space to the nominal output space



Q&A Time

- In Week 2:

No lecture/tutorial due to public holiday.

Contact is open for consultation.

