

Introduction to Data Science

FIT1043

Monash University

The word "WELCOME" is written in a large, bold, and colorful font. Each letter features a unique pattern: 'W' has pink and purple swirls; 'E' has a rainbow plaid; 'L' has blue and yellow dots; 'C' has red and orange polka dots; 'O' has vertical stripes in yellow, green, and blue; 'M' has horizontal stripes in purple, pink, and yellow; and 'E' has a yellow sunburst pattern.

to FIT1043: Introduction to data science

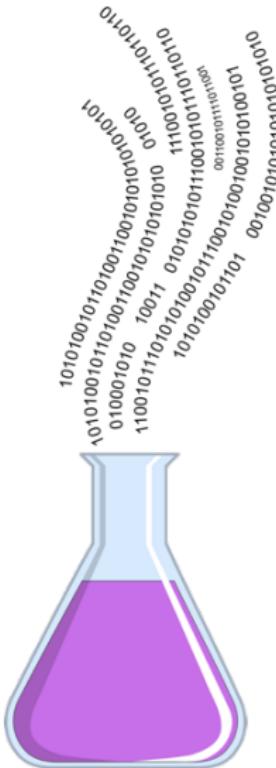


Image modified from: <https://pixabay.com/en/beaker-glass-science-lab-23417/>
and <https://openclipart.org/detail/202770/binary-code-wave>

Before we get started

Do the following:

1. get out your device
2. login to Moodle and download these slides from the FIT1043 page
3. complete [this survey](#)
4. download [this ePub](#) from Alexandria

About this Unit

Resources

1. review of Moodle

- ▶ contains Unit Orientation, Assessments and Discussion Forums
- ▶ as well as Lecture Notes, which contain active links to recommended videos & readings

2. review of Alexandria + ePubs

- ▶ contains LOTS of additional resources and exercises
- ▶ use as an online textbook

Resources

1. review of Moodle

- ▶ contains Unit Orientation, Assessments and Discussion Forums
- ▶ as well as Lecture Notes, which contain active links to recommended videos & readings

2. review of Alexandria + ePubs

- ▶ contains LOTS of additional resources and exercises
- ▶ use as an online textbook

3. additional textbook:

- ▶ No “perfect” *Introduction to Data Science* textbook available
- ▶ But a good introductory text available for purchase is:
The Art of Data Science by Peng & Matsui

Resources

1. review of Moodle

- ▶ contains Unit Orientation, Assessments and Discussion Forums
- ▶ as well as Lecture Notes, which contain active links to recommended videos & readings

2. review of Alexandria + ePubs

- ▶ contains LOTS of additional resources and exercises
- ▶ use as an online textbook

3. additional textbook:

- ▶ No “perfect” *Introduction to Data Science* textbook available
- ▶ But a good introductory text available for purchase is:
The Art of Data Science by Peng & Matsui

4. be aware also of the:

- ▶ library services available
- ▶ special consideration policies
- ▶ disability support available

Getting Started

1. no tute this 1st week
2. but there are activities to do, so check them on Moodle
 - ▶ **Module 1: Data Science and Data in Society**

Getting Started

1. no tute this 1st week
 2. but there are activities to do, so check them on Moodle
 - ▶ Module 1: Data Science and Data in Society
 3. how these classes are run
 - ▶ watch videos and read background material between classes
 - ▶ prepare for labs/tutes
 - ▶ bring a device to lectures to participate
 - ▶ see Module 7 in ePub for Data Science Resources

Contacts

Need help?

1. ask questions during tutorials and lectures
 - ▶ please interrupt me with questions!
2. check for relevant **Discussions** on Moodle
 - ▶ note in particular the “Assessments” discussion threads
3. attend the consultation hour of the tutors or the lecturer
 - ▶ consultation hours in Moodle
4. send email to tutor or lecturer

Motivation for the Unit

Data Science is in its **growth phase**:

- ▶ every academic & industry community wants to claim credit
- ▶ huge community of “leading international experts,” “highly sought-after consultants,” and “thought leaders” to confuse you with advice, blogs, guidelines, ...
- ▶ huge growth in software and services

We try and cover **the full extent of what makes Data Science**:

- ▶ background and context
- ▶ leading review articles, lectures, introductions
- ▶ academic surveys and national programmes

Prerequisites

You will need:

- ▶ high school level of mathematics and statistics
- ▶ a “critical mindset”:
 - ▶ you will read/view a variety of material
 - ▶ different levels of quality and standards
 - ▶ some sales, some educational, some journalistic
- ▶ basic exposure to information technology and internet businesses:
 - ▶ software, science or business computing
 - ▶ MS Excel
 - ▶ laptop or tablet, ...
 - ▶ Amazon, Google, Twitter, ...
- ▶ basic exposure to “business thinking”

Warning

Alexandria links to a LOT of content:

- ▶ videos, blogs, articles, ...
- ▶ there is **way too much** for you to read it all in detail!

Strategy:

- ▶ limit your time per week
- ▶ get the big picture from articles/videos
- ▶ find out what is out there
- ▶ focus in on the details when you need something for assessment (or on stuff you want for your own development)

Unit Schedule: Modules

| Module | Week | Content |
|---------------|-------------|------------------------------------|
| 1. | 1 | overview and look at projects |
| | 2 | (job) roles, and the impact |
| 2. | 3 | data business models |
| | 4 | application areas and case studies |
| 3. | 5 | characterising data and "big" data |
| | 6 | data sources and case studies |
| 4. | 7 | resources and standards |
| | 8 | resources case studies |
| 5. | 9 | data analysis theory |
| | 10 | data analysis process |
| 6. | 11 | issues in data management |
| | 12 | data management frameworks |

Any Other Questions?

... before we get started?

... otherwise, use the discussion board!

Your Background

Describe a little bit about yourself: (link only
works for the lecturer ;-)

FIT1043: Introduction to Data Science

Module 1: Data Science and Data in Society

Lecture 1: Overview and a Look at Projects

Monash University

Unit Schedule: Modules

| Module | Content |
|--------|--|
| 1. | overview and look at projects (job) roles, and the impact |
| 2. | data business models application areas and case studies |
| 3. | characterising data and "big" data data sources and case studies |
| 4. | resources and standards resources case studies |
| 5. | data analysis theory data analysis process |
| 6. | issues in data management data management frameworks |

Overview of Data Science (ePub section 1.1+1.3)

a quick overview of the context

Who are the Data Scientists?

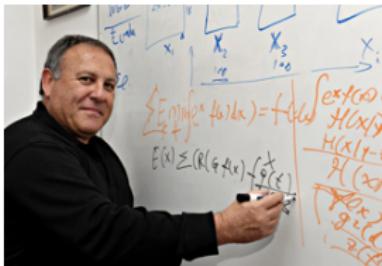
Vote by clicking on the [survey form](#)



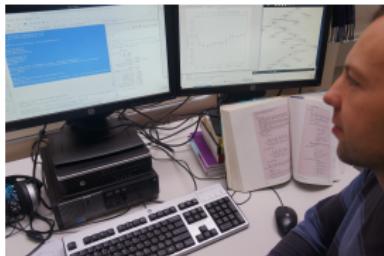
person A



person B



person C



person D

What is Data Science?

“name contains the word ‘science’, so it can’t be one”

- ▶ *Note: this is an old joke ...*

What is Data Science?

“name contains the word ‘science’, so it can’t be one”

- Note: this is an old joke ...

“data science is what a data scientist does”

- ### ► *a circular definition!*

What is Data Science?

“name contains the word ‘science’, so it can’t be one”

- Note: this is an old joke ...

“data science is what a data scientist does”

- ### ► *a circular definition!*

“data science is the technology of handling and extracting value from data”

- less circular and a bit more useful

What is Data Science?

“name contains the word ‘science’, so it can’t be one”

- Note: this is an old joke ...

“data science is what a data scientist does”

- ### ► *a circular definition!*

“data science is the technology of handling and extracting value from data”

- less circular and a bit more useful

“machine learning on big data”

- *useful, but too narrow!*

Machine Learning Definition

(well understood and agreed on)

Machine Learning (ML) is concerned with the development of algorithms and techniques that allow computers to *learn*.

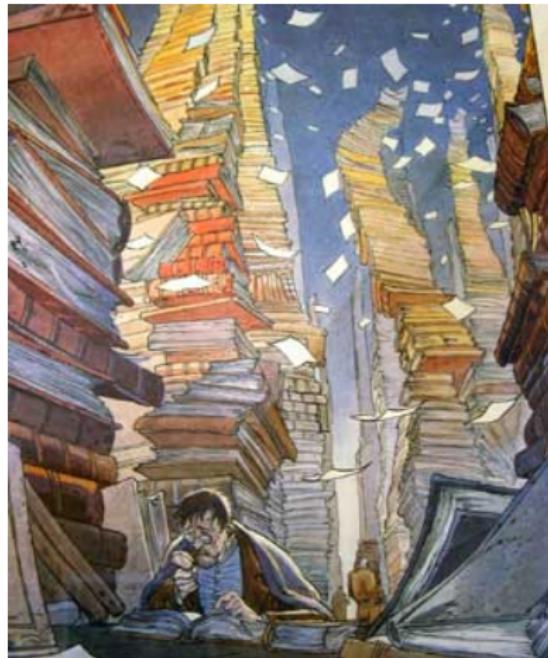
- ▶ concerned with building computational artifacts
 - ▶ but the underlying theory is statistics

Why use Machine Learning?

Machine learning is useful when:

- ▶ Human expertise is not available
e.g. Martian exploration
 - ▶ Humans cannot explain their expertise (as a set of rules), or their explanation is incomplete and needs tuning
e.g. speech recognition
 - ▶ Many solutions need to be adapted automatically
e.g. user personalisation
 - ▶ Situation changes over time
e.g. junk email
 - ▶ There are large amounts of data
e.g. discover astronomical objects
 - ▶ Humans are expensive to use for the work
e.g. handwritten zipcode recognition

Why use Machine Learning?



because you do not want
to be this poor guy!

- ▶ sifting through all
the data by hand

Why use Machine Learning?

Other reasons for needing Machine Learning:

- ▶ the information society
 - ▶ information warfare
 - ▶ information overload
 - ▶ information access

Exercise: Google these to find out about them!

Data Science Examples

Some famous data science projects and investigations:

1. Google's spell checker and translation engine
 - ▶ we'll learn about these in Module 5
 2. Amazon.com's recommendation engine
 3. Public health: "saturated fat is not bad for you after all"
 - ▶ many more of this type of investigation will be coming ...
 4. Microsoft's predictive analytics for traffic

Another Example: Melbourne Datathon 2016

Background for this is Section 1.2 of the ePub

- ▶ Seek.com is an online jobs website, which provided data and tasks for a data science competition.
- ▶ The tasks were:
 - ▶ **label prediction:** predict if a job is in the “Hotel and Tourism” category
 - ▶ **data exploration:** what useful information can be discovered from the data that Seek can use.
- ▶ See their own description of the *business context and dataset.*

Datathon Questions

- ▶ how did Seek come up with their prediction task?
 - ▶ why is it important to them?
 - ▶ did a data scientist come up with the task?
 - ▶ all Datathon participants had to destroy their copies of the data at the end of the Datathon: why?
 - ▶ how would you present results of exploratory analysis to Seek.com management? see
one such presentation by the 4Quarters team

Datathon Questions, cont.

- ▶ how much data is there?
 - ▶ what software/systems could you use to do the prediction task?
 - ▶ could you introduce/find auxiliary data to do the prediction better? is that “cheating”?
 - ▶ how would you estimate how well your predictions are going?
 - ▶ how would Seek.com “fairly” evaluate participants in the Datathon?

Historical Context

Links to resources providing historical background to data science:

- ▶ [Wolfram Alpha: computable knowledge history](#)
- ▶ [Cloud Infographic: Evolution Of Big Data](#)
- ▶ [The Web Technology timeline](#)
- ▶ [A brief history of Data Science](#)

The Rise of Big Data

in [Foreign Affairs](#), by Cukier and Mayer-Schoenberger

Data Science interest is related to the arrival of “Big Data”

- ▶ **data collection** has changed:
 - ▶ lots of data, but more messy
 - ▶ don't look for perfect models – settle for finding patterns
 - ▶ examples: Google's *language translation* and *flu trends*
- ▶ **datafication**:
 - ▶ taking all aspects of life and turning them into data
 - ▶ e.g. NYC using big data to improve public services and lower costs
- ▶ the “information society” has come of age
 - ▶ and data brokers have started amassing huge data about individuals: *big data could become Big Brother*

Homework

From Section 1.1:

- ▶ watch Cukier's TED talk on "Big Data"
- ▶ watch the video on "Big Data" by Tim Smith

The Data Science Process

(ePub section 1.2)

what happens in a Data Science project?

- ▶ illustrating the process
 - ▶ a quick walkthrough illustrating the steps
 - ▶ the standard value chain
 - ▶ our model of the process

The Data Science Process: Illustrating the Process

a quick walkthrough illustrating the steps

The Data Science Process

Many different tasks come together to complete a Data Science project

- ▶ a data scientist should be familiar with most, but doesn't need to be an expert in all

Not all are labelled as Data Science

- ▶ some from other field such as computer engineering, business, ...



Pitching ideas for data science projects to
investors/managers.

"Young Business Man Holding a Tablet" by Pic Basement, CC-BY 2.0

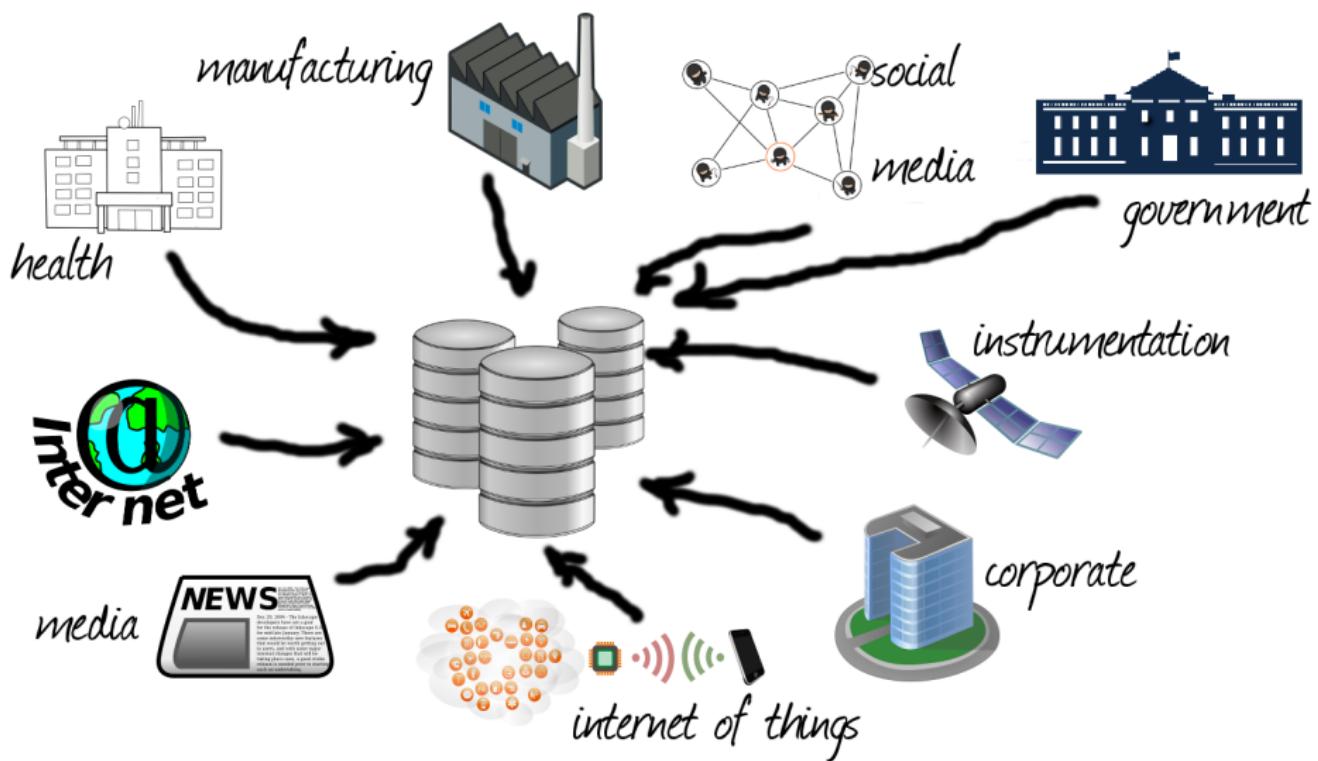


Collecting data: researchers preparing to x-ray a patient.

by Stephen Ausmus acquired from USDA ARS, public domain.

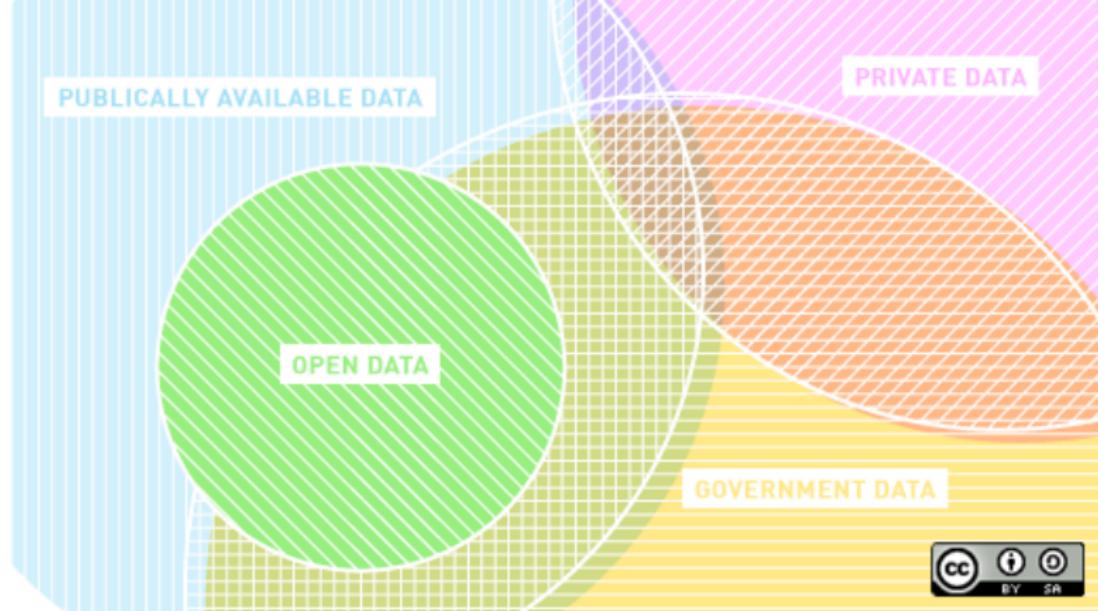


Monitoring: Scientists watch over data collected by the gravimeter & magnetometer instruments.



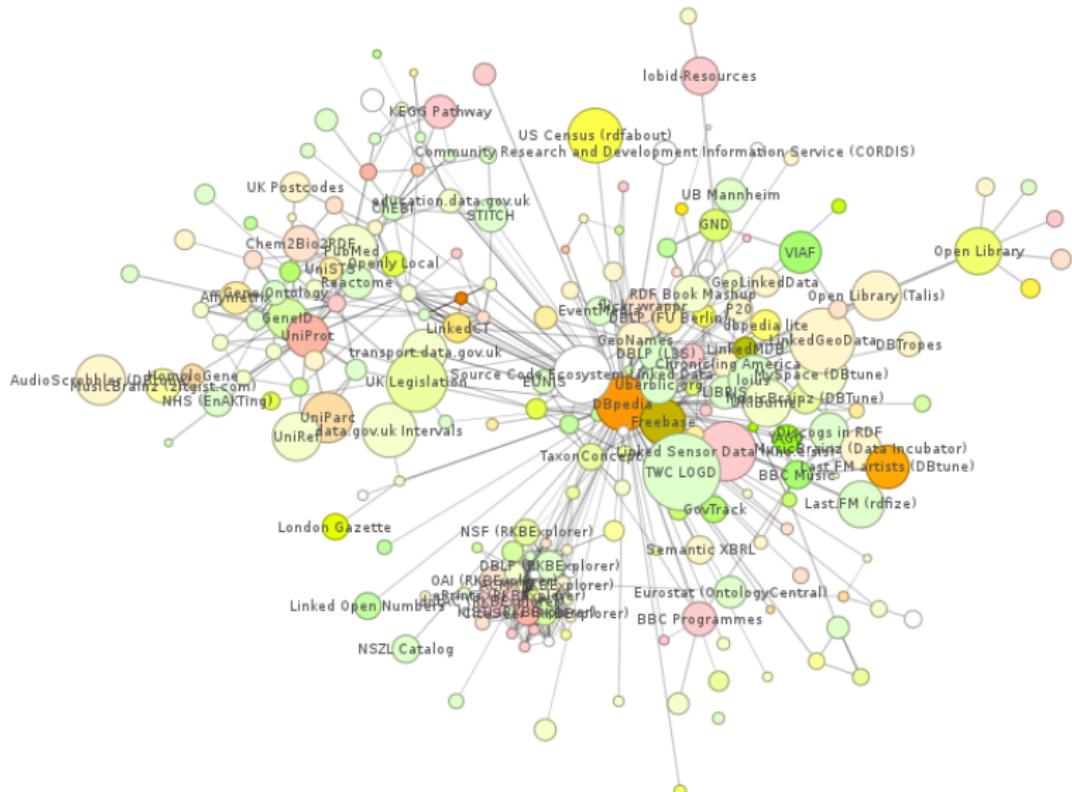
Integration: Data can come from many different sources.

icons from [by Openclipart.org](http://Openclipart.org), public domain



Some of the best data is Open (publicly available) Data.

by Libby Levi for opensource.com, CC-BY-SA 2.0



Interpretation: The Linked Open Data (LOD) graph can ascribe meaning (semantics) to data.

by Open Knowledge, CC-BY-SA 2.0





Navigation: of data standards and formats

“The Web is Agreement” cropped, by Paul Downey, CC-BY 2.0

Understanding a database schema.



archiving



storage



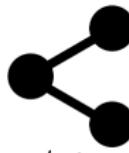
privacy



legal & compliance



safety



sharing



metadata



management



ethics

Governance cares for the data and its subjects.

icons from [by Openclipart.org](http://Openclipart.org), public domain;

Good and Evil by AJC ajcann.wordpress.com, CC-BY-SA 2.0



Slide 41 / 66



Data engineers make the back-end work

File Edit Code View Plots Session Project Build Tools Help

RStudio

Project: (None)

Files Plots Packages Help

R packages available saps mararie notes.R

Source on Save Run Source

genes5

In selection Match case Whole word Regex Wrap

```

96 ~ for(i in 1:nrow(p_pure)){ ## voor elke gene set
97   conGenes<-intersect(genes,unique(as.character(geneSets[,i])))
98   if(length(conGenes)<2) {## als er geen overlappen, doe dan iets
99     p_pure[,i,"Size"]<-length(conGenes) # stop het aantal overlappende genen in de matrix
100   }
101 
102   # Global
103   data<-scale(data[,x,i.element(genes,conGenes)]) # data genlist voor alle patienten
104   lab<-kmeans(data[,2])$cluster # clusteren patienten in groep 1 of 2
105   survtest<-survdiff(Surv(time[event==1]-lab)
106   p_pure[i,"Global"]<- 1 - pchisq(survtest$chisq, 1)
107 
108   # For ovary
109   if(anType=="Ov"){
110     data<-scale(data[,x=="Angiogenic",is.element(genes,conGenes)])
111     lab<-kmeans(data[,2])$cluster
112     survtest<-survdiff(Surv(time[st=="Angiogenic"],event[st=="Angiogenic"])-lab)
113     p_pure[i,"Angiogenic"]<- 1 - pchisq(survtest$chisq, 1)
114 
115     data<-scale(data[,x!="Non-angiogenic",is.element(genes,conGenes)])
116     lab<-kmeans(data[,2])$cluster
117     survtest<-survdiff(Surv(time[st=="Non-angiogenic"],event[st=="Non-angiogenic"])-lab)
118     p_pure[i,"Non-angiogenic"]<- 1 - pchisq(survtest$chisq, 1)
119   }
120 
121   # For Breast
122   if(anType=="Br"){
123     data<-scale(data[,x=="ER+/HER2- High Prolif",is.element(genes,conGenes)])
124     lab<-kmeans(data[,2])$cluster
125     survtest<-survdiff(Surv(time[st=="ER+/HER2- High Prolif"],event[st=="ER+/HER2- High Prolif"])-lab)
126     p_pure[i,"ER_H"]<- 1 - pchisq(survtest$chisq, 1)
127 
128     data<-scale(data[,x=="ER+/HER2- Low Prolif",is.element(genes,conGenes)])
129     lab<-kmeans(data[,2])$cluster
130     survtest<-survdiff(Surv(time[st=="ER+/HER2- Low Prolif"],event[st=="ER+/HER2- Low Prolif"])-lab)
131     p_pure[i,"ER_L"]<- 1 - pchisq(survtest$chisq, 1)
132 
133     data<-scale(data[,x=="HER2+",is.element(genes,conGenes)])
134   }

```

102:15 (Untitled): R Script

Workspace History

Data

| | |
|--------|--------------------------|
| dat | 1678x11247 double matrix |
| dat.st | 1670x11247 double matrix |
| dat.x | 1670x11247 double matrix |
| dat1 | 1670x5 double matrix |
| dat5 | 1670x17 double Matrix |

Console ~/Documents/data/saps paper data/molsigdb.v3.0.enriched.Rf ↵

| | | | | | |
|--------------|--------------|--------------|--------------|--------------|--------------|
| TCGA-61-1984 | TCGA-61-1906 | TCGA-61-1987 | TCGA-61-1918 | TCGA-61-1911 | TCGA-61-1913 |
| 1 | 1 | 1 | 1 | 1 | 1 |
| TCGA-61-1914 | TCGA-61-1915 | TCGA-61-1917 | TCGA-61-1918 | TCGA-61-1919 | TCGA-61-1995 |
| 1 | 1 | 2 | 1 | 2 | 1 |
| TCGA-61-1998 | TCGA-61-2000 | TCGA-61-2002 | TCGA-61-2003 | TCGA-61-2008 | TCGA-61-2009 |
| 1 | 1 | 2 | 2 | 2 | 2 |
| TCGA-61-2012 | TCGA-61-2016 | TCGA-61-2017 | TCGA-61-2018 | TCGA-61-2087 | TCGA-61-2088 |
| 1 | 1 | 1 | 1 | 1 | 1 |
| TCGA-61-2092 | TCGA-61-2094 | TCGA-61-2095 | TCGA-61-2096 | TCGA-61-2097 | TCGA-61-2098 |
| 1 | 1 | 2 | 1 | 2 | 2 |
| TCGA-61-2101 | TCGA-61-2102 | TCGA-61-2104 | TCGA-61-2109 | TCGA-61-2110 | TCGA-61-2111 |
| 2 | 2 | 1 | 1 | 1 | 1 |
| TCGA-61-2113 | X1 | X101 | X109 | X11 | X112 |
| 2 | 2 | 2 | 2 | 1 | 2 |
| X113 | X114 | X120 | X126 | X127 | X128 |
| 1 | 2 | 1 | 1 | 2 | 2 |
| X118 | X14 | X146 | X143 | X146 | X147 |
| 1 | 1 | 2 | 2 | 2 | 1 |
| X157 | X159 | X16 | X163 | X164 | X165 |
| 2 | 1 | 2 | 2 | 2 | 1 |
| X167 | X168 | X182 | X2 | X216 | X217 |
| 2 | 1 | 1 | 1 | 1 | 2 |
| X234 | X240 | X252 | X3 | X30 | X314 |
| 1 | 2 | 2 | 1 | 2 | 1 |
| X317 | X336 | X34 | X345 | X346 | X347 |
| 2 | 2 | 2 | 1 | 2 | 1 |
| X35 | X352 | X355 | X358 | X36 | X362 |
| 1 | 2 | 2 | 2 | 1 | 2 |
| X363 | X37 | X41 | X43 | X46 | X65 |
| 2 | 1 | 1 | 2 | 1 | 1 |
| X89 | X9 | | | | |
| 1 | 2 | | | | |

> lab[1:4]

| | | | | | |
|------------|--------------|--------------|--------------|--|--|
| 1_Cy5_5258 | 101_Cy5_5379 | 103_Cy5_S117 | 105_Cy5_5457 | | |
| 2 | 1 | 2 | 2 | | |

> lab[1:28]

| | | | | | | | |
|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| 1_Cy5_5258 | 101_Cy5_5379 | 103_Cy5_S117 | 105_Cy5_5457 | 107_Cy5_5425 | 11_Cy5_5463 | 111_Cy5_5482 | 121_Cy5_5235 |
| 2 | 1 | 2 | 2 | 1 | 2 | 2 | 2 |
| 13_Cy5_S429 | 131_Cy5_5267 | 137_Cy5_S423 | 147_Cy5_S355 | 149_Cy5_S111 | 151_Cy5_S293 | 155_Cy5_S431 | 157_Cy5_S341 |
| 2 | 1 | 2 | 2 | 1 | 1 | 2 | 2 |
| 159_Cy5_S482 | 163_Cy5_5232 | 165_Cy5_5444 | 17_Cy5_S413 | | | | |
| 2 | 2 | 2 | 2 | | | | |

> ?kmeans

> |

Inspecting and cleaning the data.

"rstudio" by mararie, CC-BY-SA 2.0



$$\zeta(s) = \sum_{n=1}^{\infty} \frac{\mu(n)}{n^s},$$

$$\eta(s) = \sum_{k=0}^{\infty} \frac{(-1)^k}{k^s} = 1 - \frac{1}{2^s} + \frac{1}{3^s} - \frac{1}{4^s} + \dots$$

$$\zeta(n) = \sum_{k=1}^{\infty} \frac{1}{k^n} = \frac{1}{1^n} + \frac{1}{2^n} + \frac{1}{3^n} + \dots = \sum_{k=1}^{\infty} \frac{1}{k^n}$$

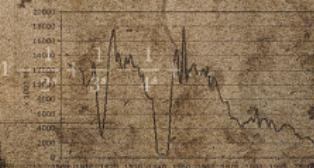
$$(n) = 1 + \frac{1}{2^n} + \int_3^{\infty} s(t) dt = \sum_{k=0}^{\infty} \frac{1}{k^n}$$

$$= \sum_{k=0}^{\infty} \frac{(-1)^k}{k^s} = 1 - \frac{1}{2^s} + \frac{1}{3^s} - \frac{1}{4^s} + \dots$$

$A=0$

$$\sum_{n=1}^{\infty} \frac{1}{n^s} = \prod_{p \text{ prime}} \frac{1}{1 - p^{-s}} \quad \text{for } s > 1$$

$$= \sum_{k=0}^{\infty} \frac{(z-1)^k}{k!} = 1 + \frac{z-1}{1!} + \frac{(z-1)^2}{2!} + \frac{(z-1)^3}{3!} + \dots$$



$$\sum_{n=1}^{\infty} \frac{1}{n^2} = \frac{\pi^2}{6}$$

1 1 1 1

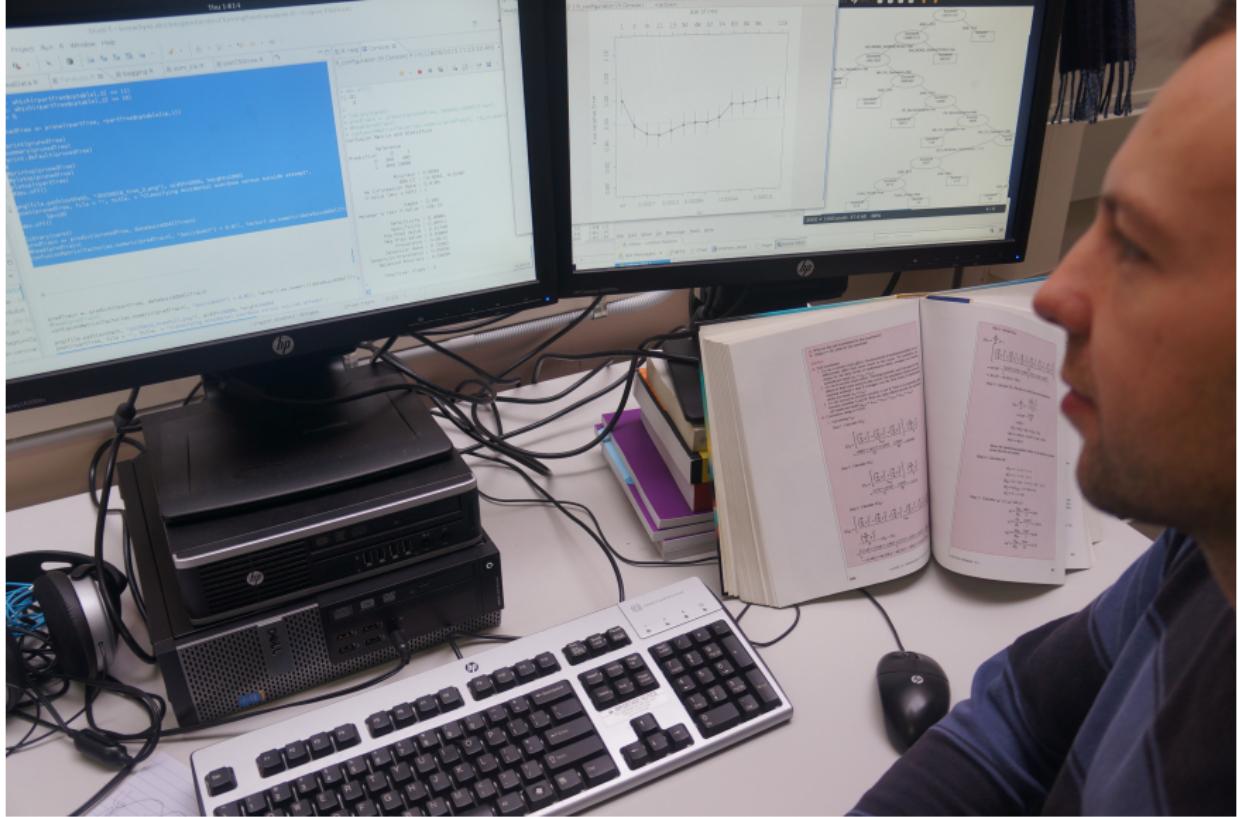
I - J S C T J A L

1 33

卷之三

卷之三

Proposing a conceptual/mathematical/functional model.



Analyst building models with his favourite tool.

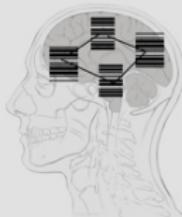
Data



Information



Knowledge



Understanding



Wisdom



Facts

No relations, patterns
or principles

Who, What,
When, Where
Gives Meaning

How-to
Inside our heads
Application of Information

Answers the question
Why?

What is best?

Doing the right things
What should be done

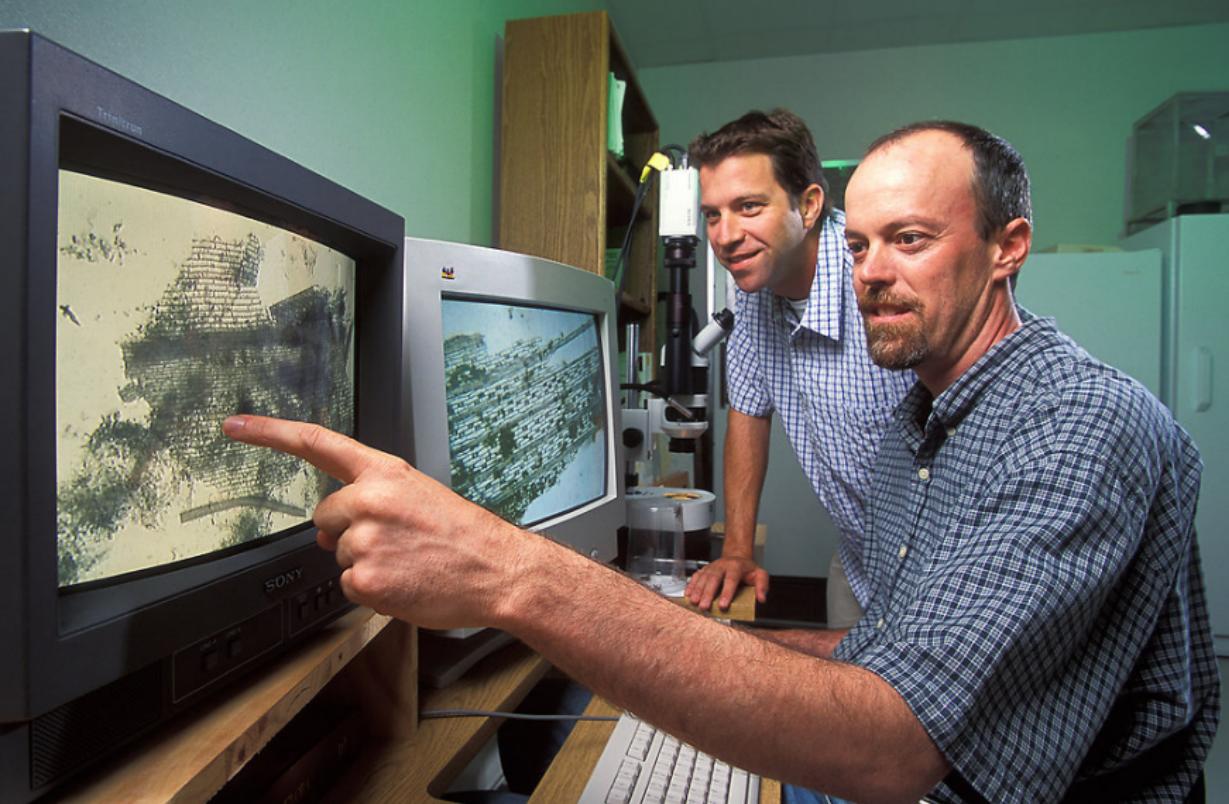


Analysis, statistics and/or machine learning works on the data.



Choosing appropriate visualizations for the data. Many different options exist!

"Visualization Matrix" cropped, by Lauren Manning, CC-BY 2.0

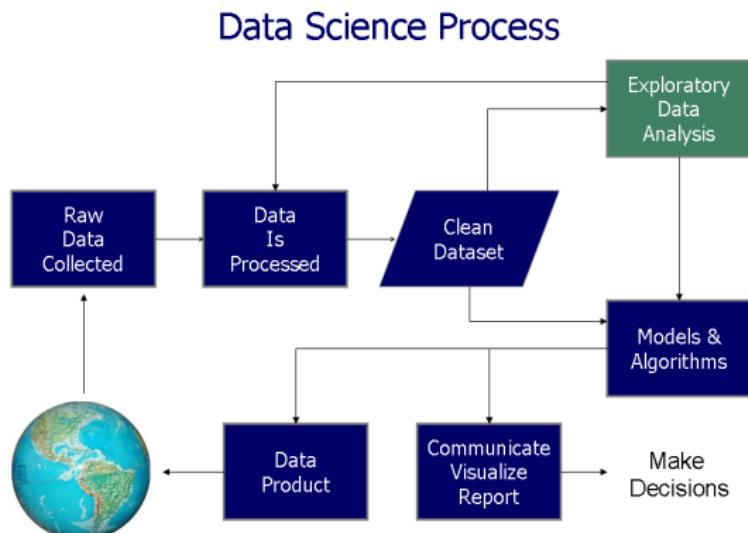


Visualising data to interpret it and present results.

by Stephen Ausmus acquired from USDA ARS, public domain.

The Data Science Process

- ▶ A flowchart showing the overall data science process





Operationalization: putting the results to work.

The Data Science Process: Our Standard Value Chain

our model of the process

Parts of a Data Science Project

Collection: gathering the data

Engineering: storage and computational resources across full lifecycle

Governance: overall management of data across full lifecycle

Wrangling: data preprocessing, cleaning

Analysis: discovery (learning, visualisation, etc.)

Presentation: arguing the case that the results are significant and useful

Operationalisation: putting the results to work, so as to gain benefits or value

We call this the **Standard Value Chain**.

Interpretations of the Parts

Following [Jeff Hammerbacher's](#) UC Berkeley 2012 course notes, we will interpret these four entities:

- ▶ business analyst
- ▶ programmer
- ▶ enterprise
- ▶ web company

Interpretations: the Business Analyst

Collection: copy and paste into Excel

Engineering: use Excel to store and retrieve

Wrangling: use Excel functions, VBA

Analysis: charts

Interpretations: the Programmer

Collection: web APIs, scraping, database queries

Engineering: flat files

Wrangling: Python and Perl, etc.

Analysis: Matplotlib in Python, R

Interpretations: the Enterprise

Collection: application databases, intranet files, server logs

Engineering: Teradata, Oracle, MS SQL Server

Wrangling: Talend, Informatica

Analysis: Cognos, Business Objects, SAS, SPSS

Interpretations: the Web Company

Collection: application databases, server logs, crawl data

Engineering: Hadoop/Hive, Flume, HBase

Wrangling: Pig, Oozie

Analysis: dashboards, R

What is Data Science? (ePub section 1.3)

how can we define or circumscribe data science?

Wikipedia Definitions

Data Science is the extraction of knowledge from data, which is a continuation of the field of data mining and predictive analytics.

Big data is a broad term for data sets so large or complex that traditional data processing applications are inadequate.

Pivotal Definition

Data Science: The use of statistical and machine learning techniques on big multi-structured data in a distributed computing environment to identify correlations and causal relationships, classify and predict events, identify patterns and anomalies and infer probabilities, interest and sentiment.

NIST Big Data Working Group Definition

Data Science is the empirical synthesis of actionable knowledge from raw data through the complete data lifecycle process.

A data scientist is a practitioner who has sufficient knowledge in the overlapping regimes of business needs, domain knowledge, analytical skills, and software and systems engineering to manage the end-to-end data processes through each stage in the data lifecycle.

Journal of Data Science

Definition

Data Science is almost everything that has something to do with data: collecting, analyzing, modeling..... yet the most important part is its applications — all sorts of applications.

Definitions: Summary

narrow: “machine learning on big data”

broad: extraction of knowledge/value from data “through the complete data lifecycle process”

- ▶ broad concern with the different stages
- ▶ focus on the learning/knowledge discovery

Data Science Relationships



Related: Data Engineering

building scalable systems for storage, processing data

- ▶ e.g. Amazon Web Services, Teradata, Hadoop, ...
- ▶ databases, distributed processing,datalakes, cloud computing, GPUs, wrangling, ...
- ▶ huge, continuous improvement

Related: Data Analysis

performing analysis and understanding results

- ▶ e.g. R, Tableau, Weka, Microsoft Azure Machine Learning, ...
 - ▶ machine learning, computational statistics, visualisation, ...
 - ▶ huge, continuous improvement

Related: Data Management

managing data through its lifecycle

- ▶ e.g. ANDS, Talend, Master Data Management, ...
- ▶ ethics, privacy, providence, curation, backup, governance, ...
- ▶ huge, continuous improvement

Fits and Starts

- ▶ Data Analysis (John Tukey) in 1962
- ▶ Expert Systems in the 1980's
- ▶ Machine Learning in the 1980's
- ▶ Data Mining in the 1990's
- ▶ see *Business Week's "Database Marketing"* cover story September 1994

Data Science Emerges ~2000

- ▶ data analysis came of age 1990's
- ▶ William Cleveland publishes in 2001
Data Science: An Action Plan for ... the field of Statistics
- ▶ data engineering came of age 2000's (Dot.Com boom)
- ▶ (digital) data management came of age 2000's (Dot.Com boom)
- ▶ the data/information society
- ▶ business pressure on decision making
- ▶ "data" as a valuable asset
- ▶ Dot.Com companies show the way

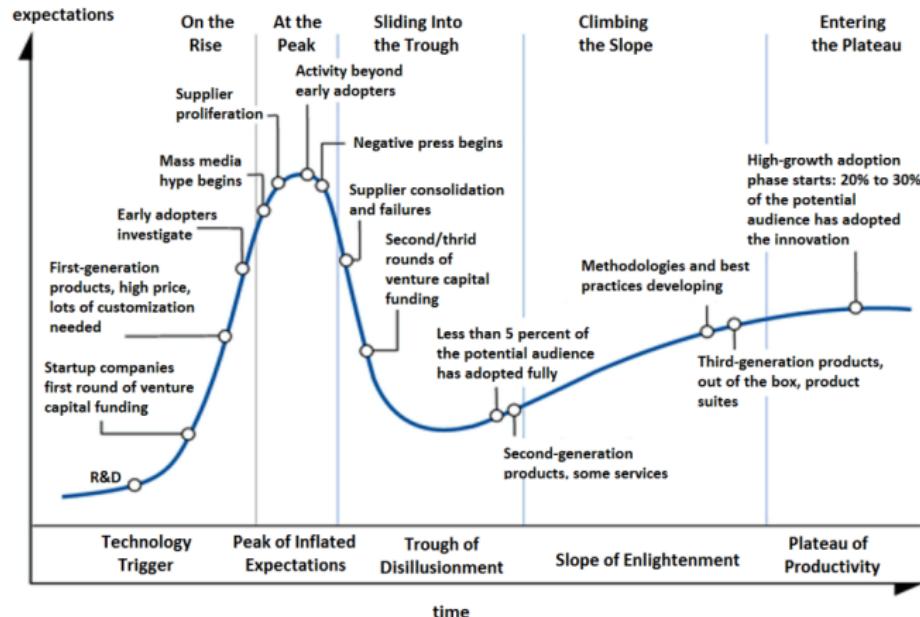
see also David Donoho's *50 years of Data Science* (PDF paper)

Homework

- ▶ look up Dot.Com boom in Wikipedia

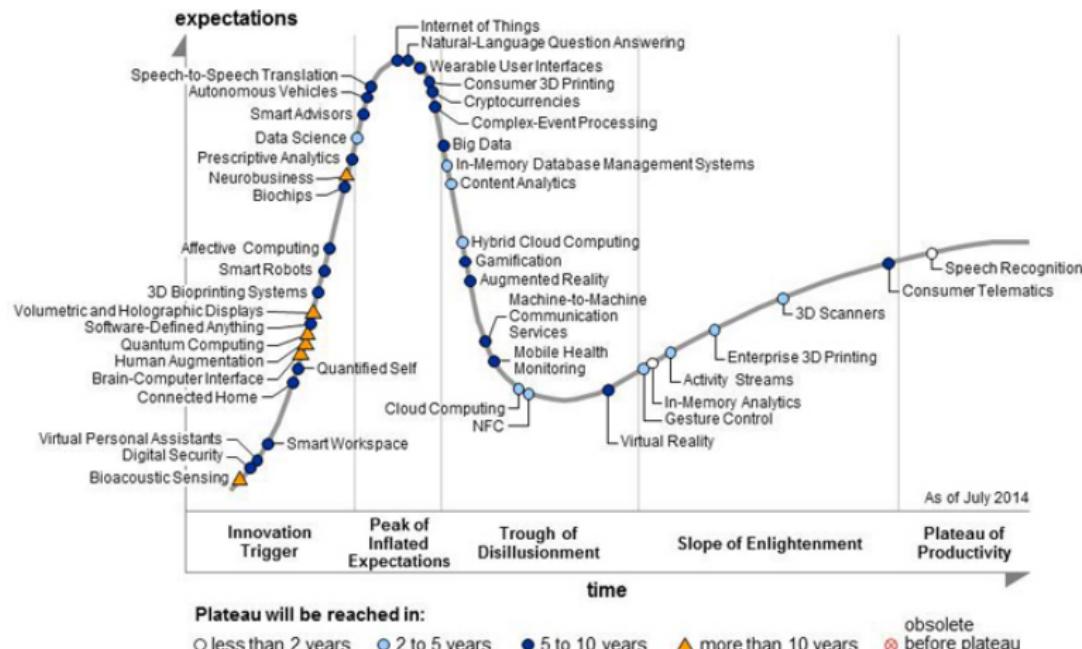
The Hype Cycle

- ▶ Gartner's Hype Cycle® attempts to quantify the level of maturity of various technologies:



Hype Cycle for 2014

- ▶ Can you spot Data Science?



Career as Data Scientist

Data science currently provides a strong career path.

To become a specialist you need:

- ▶ solid machine learning and statistics
- ▶ related mathematics (1st+2nd year in many degrees)
- ▶ solid programming/prototyping (Python, Java)
- ▶ Unix experience (Linux, Mac OSX)

See also:

- ▶ the infamous Metromap: [becoming a data scientist](#)
- ▶ and [Modern Data Scientist](#) on Pinterest

This unit provides an introduction and background (but not the core analyst skills) to become a data scientist

Data Science Research

Data Science is seeing major growth at universities internationally

Many research programs exist, including:

- ▶ US National Institute of Standards' Big Data Working Group (2013-2015)
- ▶ US National Academy of Sciences' Committee on the Analysis of Massive Data (2013)
- ▶ Alan Turing Institute for Data Science at London's new Knowledge Quarter (near National Library, 2016-)

End of Week 1