

FIT1043 Introduction to Data Science

Week 2: Data Scientist Roles and Skills,
Impact of Data Science & Business Models
with Data

Dr. Sicily Ting

School of Information Technology
Monash University Malaysia

Week 1 Coverage

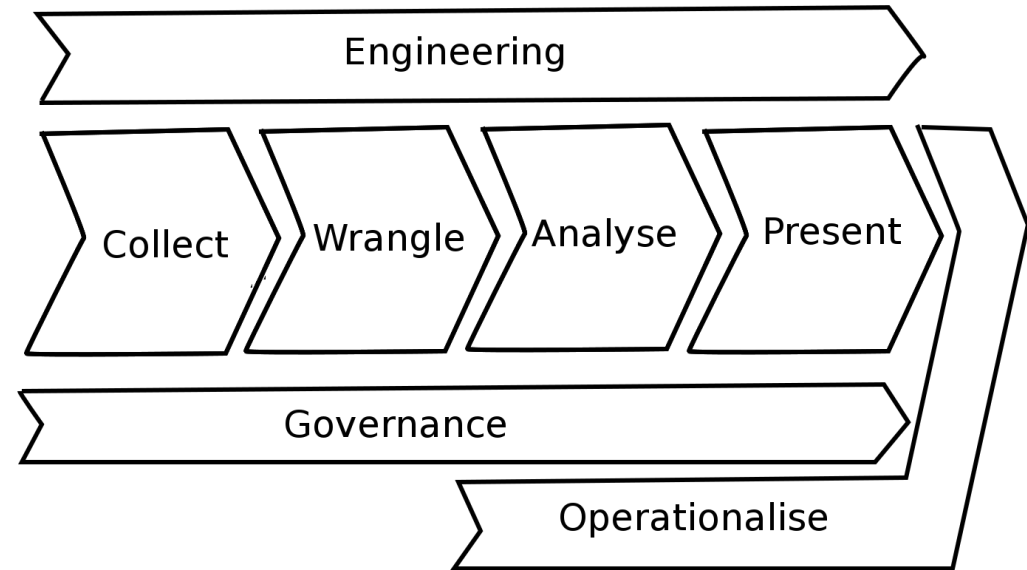
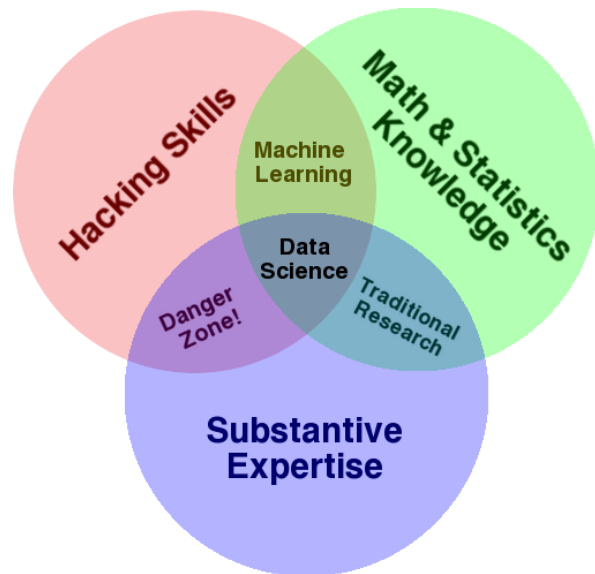
Python for Data Science

Overview of Data Science



Week 1 Coverage

- Why study data science?
 - We had a look at “Data Scientist” as a job last week. How about searching for:
 - [Data Analyst](#)
 - [Machine Learning](#)
- Drew Conway’s Venn Diagram
- Usefulness of Machine Learning
- Data Science Process and Our Standard Value Chain



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

We will refer to this throughout the semester!

Collection

- Getting the data

Engineering

- Storage and computational resources across full lifecycle

Governance

- Overall management of data across full lifecycle

Wrangling

- Data pre-processing, cleaning

Analysis

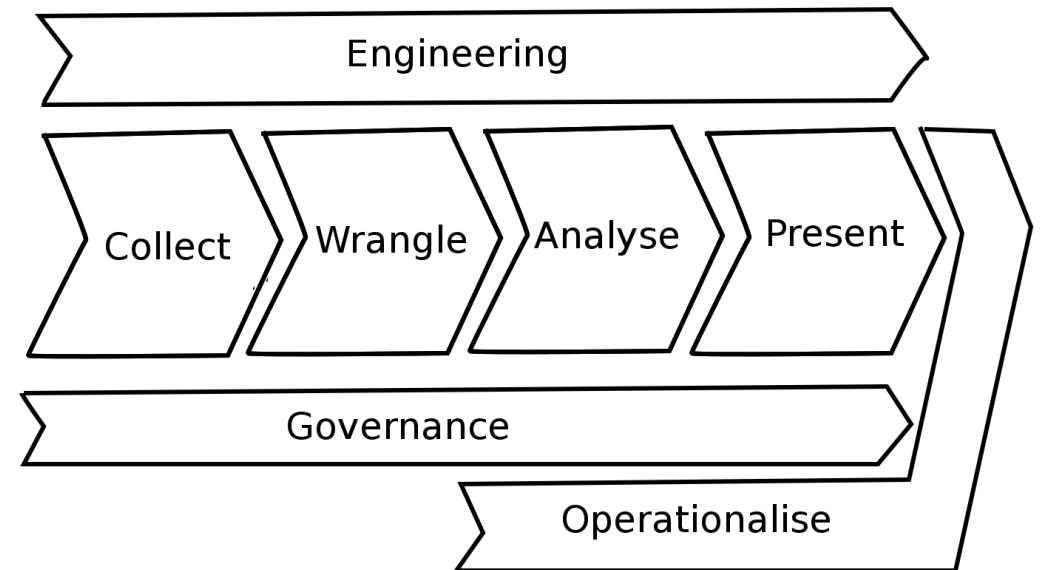
- Discovery (learning, visualisation, etc.)

Visualization

- Arguing the case that the results are significant and useful

Operationalize

- Putting the results to work, so as to gain benefits or value



Week 1

Overview of data science

Engineering

Weeks 9-10

Week 4

Collect

Wrangle

Analyse

Present

Week 3

Weeks 5-7

Week 11

Governance

Operationalise

Weeks 2 & 8

Tools for data science

| Week | Activities | Assignments |
|------|---|--------------|
| 1 | Overview of data science | |
| 2 | Introduction to Python for data science | |
| 3 | Data visualisation and descriptive statistics | |
| 4 | Data sources and data wrangling | |
| 5 | Data analysis theory | Assignment 1 |
| 6 | Regression analysis | |
| 7 | Classification and clustering | |
| 8 | Introduction to R for data science | |
| 9 | Characterising data and "big" data | Assignment 2 |
| 10 | Big data processing | |
| 11 | Issues in data management | |
| 12 | Industry guest lecture (tentative) | Assignment 3 |

Weekly Quiz From Week 2-11
 - The quiz will open for 48 hours and you are allowed to have 2 attempts

Week 2 Outline

Introduction to Python for Data Science

- Motivation to studying Python
- Python data types (Video Lectures)
- Essential libraries

Overview of data science (con't)

- Data science roles and skills
- Impact of data science
- Business models with data

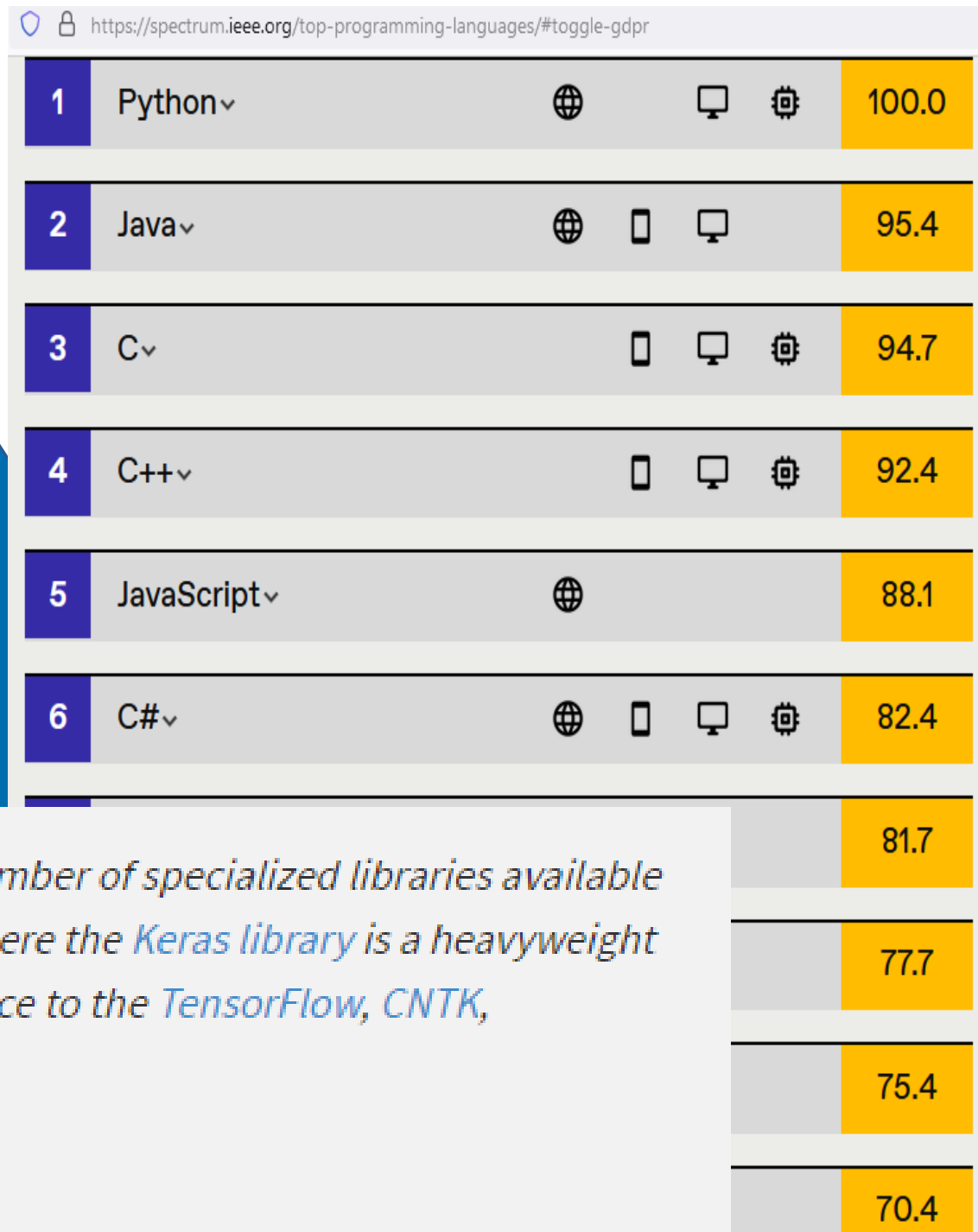
Learning Outcomes

Week 2

By the end of this week you should be able to:

- Comprehend essentials for coding in Python for data science
- Explain and interpret given **Python** codes
- Explain **different data science roles** and skills and comprehend the differences between them
- Explain **Impact** of data science
- Explain the **data business models** for organizations

IEEE Top Programming Languages in 2019








| | | | |
|---|------------|-----------|-------|
| 1 | Python | 🌐 🖥️ ⚙️ | 100.0 |
| 2 | Java | 🌐 📱 🖥️ | 95.4 |
| 3 | C | 📱 🖥️ ⚙️ | 94.7 |
| 4 | C++ | 📱 🖥️ ⚙️ | 92.4 |
| 5 | JavaScript | 🌐 | 88.1 |
| 6 | C# | 🌐 📱 🖥️ ⚙️ | 82.4 |
| | | | 81.7 |
| | | | 77.7 |
| | | | 75.4 |
| | | | 70.4 |


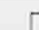









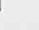




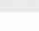

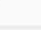
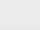


Python's popularity is driven in no small part by the vast number of specialized libraries available for it, particularly in the domain of artificial intelligence, where the *Keras library* is a heavyweight among deep-learning developers: Keras provides an interface to the *TensorFlow*, *CNTK*, and *Theano* deep-learning frameworks and tool kits.

IEEE Spectrum

5 Years Running

| Language Rank | Types | Spectrum Ranking |
|---------------|---|------------------|
| 1. Python |   | 100.0 |
| 2. C |    | 99.7 |
| 3. Java |    | 99.5 |
| 4. C++ |    | 97.1 |
| 5. C# |    | 87.7 |
| 6. R |  | 87.7 |
| 7. JavaScript |   | 85.6 |
| 8. PHP |  | 81.2 |
| 9. Go |   | 75.1 |
| 10. Swift |   | 73.7 |

2017

| Language Rank | Types | Spectrum Ranking |
|---------------|--|------------------|
| 1. Python |    | 100.0 |
| 2. C++ |    | 99.7 |
| 3. Java |    | 97.5 |
| 4. C |    | 96.7 |
| 5. C# |    | 89.4 |
| 6. PHP |  | 84.9 |
| 7. R |  | 82.9 |
| 8. JavaScript |   | 82.6 |
| 9. Go |   | 76.4 |
| 10. Assembly |  | 74.1 |

2018

| Rank | Language | Type | Score |
|------|------------|---|-------|
| 1 | Python |    | 100.0 |
| 2 | Java |    | 96.3 |
| 3 | C | | |
| 4 | C++ | | |
| 5 | R | | |
| 6 | JavaScript | | |
| 7 | C# | | |
| 8 | Matlab | | |
| 9 | Swift | | |
| 10 | Go | | |

2019

| Rank | Language | Type | Score |
|------|------------|---|-------|
| 1 | Python |    | 100.0 |
| 2 | Java |    | 95.3 |
| 3 | C |    | 94.6 |
| 4 | C++ |    | 87.0 |
| 5 | JavaScript |  | 79.5 |
| 6 | R |  | 78.6 |
| 7 | Arduino |  | 73.2 |
| 8 | Go |   | 73.1 |
| 9 | Swift |   | 70.5 |
| 10 | Matlab |  | 68.4 |
| 11 | Ruby |   | 66.8 |

2020

Python's Role in Data Science

Many tools out there for data science.

Python has gained popularity over the last few years

- easy to learn
- flexible and multi-purpose
- great libraries
- well designed computer language
- good visualization for basic analysis

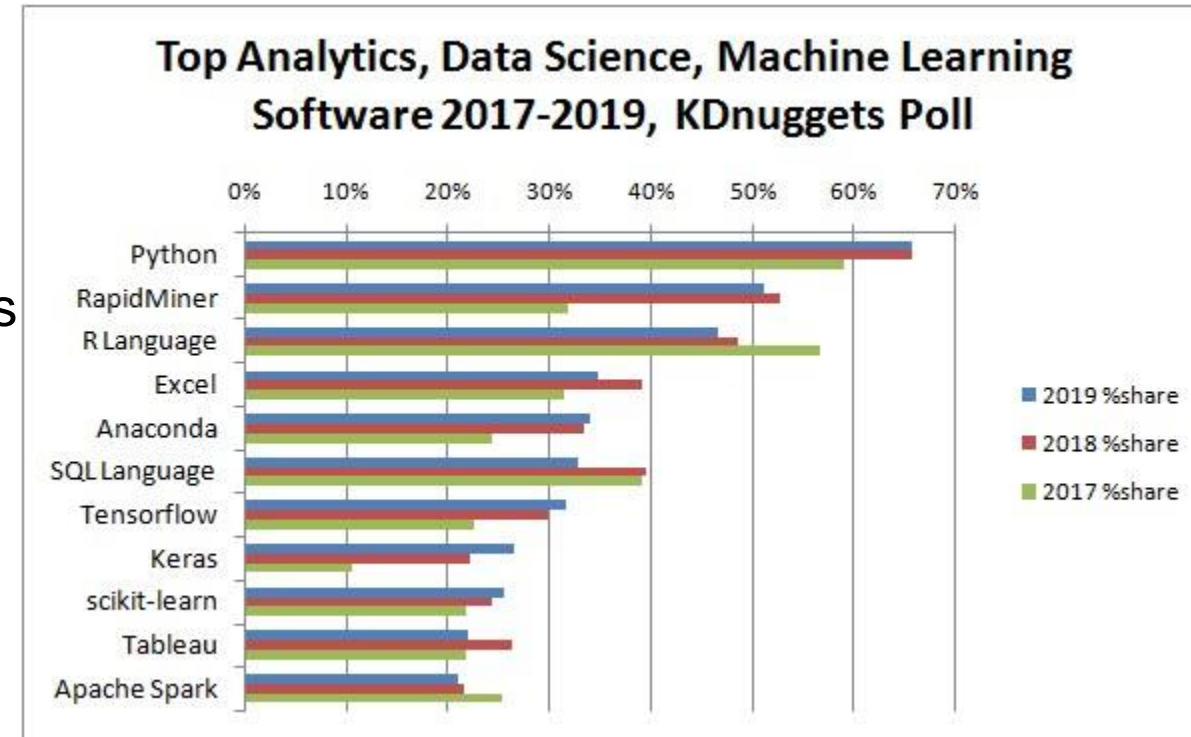









image source: kdnuggets.com

Applications on Channels

| | | | |
|---|---|---|--|
|  jupyter notebook 5.0.0 Web-based, interactive computing notebook environment. Edit and run human-readable docs while describing the data analysis. Launch |  IPy qtconsole 4.3.0 PyQt GUI that supports inline figures, proper multiline editing with syntax highlighting, graphical calltips, and more. Launch |  spyder 3.1.4 Scientific Python Development Environment. Powerful Python IDE with advanced editing, interactive testing, debugging and introspection features. Launch |  glueviz 0.10.4 Multidimensional data visualization. Explore relationships between related data. Install |
|  orange3 3.4.1 Component based data mining framework. Data visualization and data analysis for novice and expert. Interactive workflows with a large toolbox. Install |  R rstudio 1.0.136 A set of integrated tools designed to help you be more productive with R. Includes R essentials and notebooks. Install |  ANACONDA. | |

Anaconda:
Environment
Manager.



Data Scientist Roles

ePub Section 1.4



For better understanding the different kinds of data scientists:

- Reviewing:
Analyzing the Analyzers from Harris, Murphy and Vaisman
- Interviews:
From *Data Analytics Handbook*

Roles of a Data Scientist

Analyzing the Analyzers <http://www.oreilly.com/data/free/files/analyzing-the-analyzers.pdf>

A [quote from Quora from Jason Widjaja](#):

- Data analysts are primarily people who develop insights with data,
- Data scientists are primarily people who develop data models and products, that in turn produce insights, and
- Data engineers are primarily people who manage data infrastructure, automate data processing and deploy models at scale.

(Note the use of the word “primarily”!)

see also [Job Comparison](#) – Data Scientist vs Data Engineer vs Statistician

Skills of Data Scientists

Analyzing the Analyzers, Harris, Murphy and Vaisman, 2013

Business:

Product development, business

Machine learning/Big data:

Unstructured data, structured data, machine learning, big and distributed data

Mathematics/Operations research:

Optimisation, mathematics, graphical models, algorithms

Programming:

Systems administration, back end programming, front end programming

Statistics:

Visualisation, temporal statistics, surveys and marketing, spatial statistics, science, data manipulation

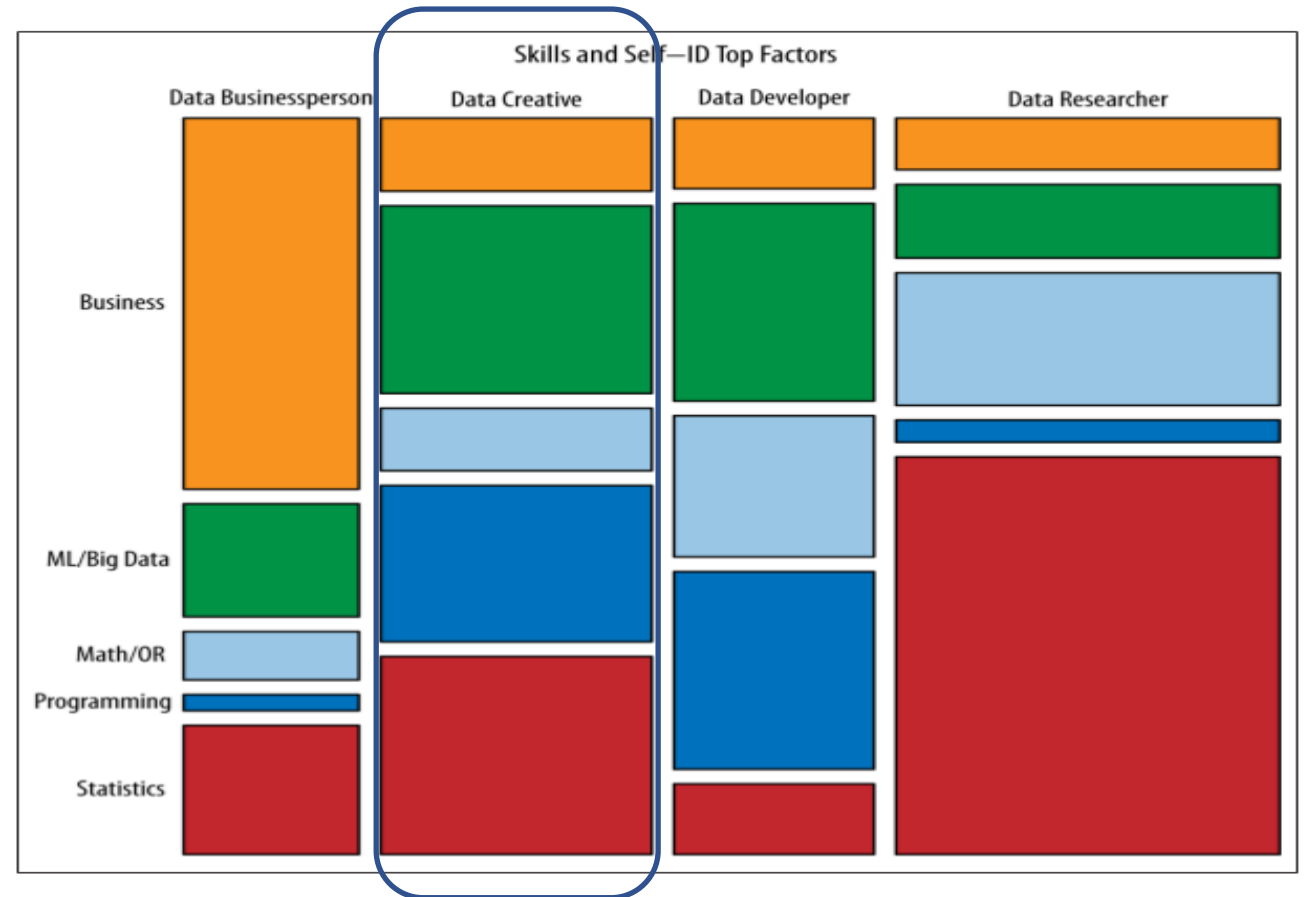
NB. typical data scientist doesn't have to know all of these!

Mapping Styles to Skills

Analyzing the Analyzers, Harris, Murphy and Vaisman, 2013

The Variety of Data Scientist (pages 14 – 16)

- Data Businesspeople
- Data Creatives
- Data Developers
- Data Researchers



Roles of a Data Scientist 2

Interviews from Data Analytics Handbook (<https://www.teamleada.com/handbook>)

From Data Analytics Handbook

The Data Analytics Handbook is a four volume set of long interviews from industry and academic professionals in the field.

Volume 1 deals with practitioners:

- What exactly do the sexy “data scientists” do?
- What other professions are there in big data?
- What tools do they use to accomplish their tasks?
- How can I enter the industry if I don’t have a Ph.D. in Statistics?

Lessons from the DA Handbook

Summary (important bits)

- **Communication** skills are underrated.
- The biggest challenge for a data analyst is the **Collection** and **Wrangling** steps.
- A data scientist is better at statistics than a software engineer and better at software engineering than a statistician.
- The data industry is still nascent (growing) and the roles less well defined so you get to interact with many parts of the company from engineering to business intelligence to product managers.
- Keep a **curiosity** about working with data, a quality as important as your technical abilities.

Career as a Data Scientist

Your CV

To become a specialist you need:

- Solid **machine learning** and **statistics**
- Related mathematics (1st+2nd year in many degrees)
- Solid prototyping (proof of concepts)
 - R, **Python**, Java
 - **Github**
 - Competitions, e.g. **Kaggle**
- **Unix experience** (Linux, Mac OSX). This unit provides an introduction and background only

Python for Data Science

Basic Python in
Pre-Class Activity





MONASH
University

ePub Section 1.6



Impact of Data Science

Some examples of how data science is impacting others:

- Your life in the cloud
 - Datafication of you
- Social good
 - Numerous examples and very rewarding
- Futurology
 - Healthcare and automobiles

Your Life on the Cloud

From Year Zero: Our life timelines begin

Our personal information is increasingly stored in the cloud:

- Social life (Facebook, etc.),
- Career (LinkedIn),
- Search history (Google, etc.),
- Health and medical (Fitbit, etc.),
- Music (Apple, Spotify, etc.).

This provides many, **many advantages**:

- Personal agents, computerised support for health.

But also **some disadvantages**:

- Security and privacy breaches.

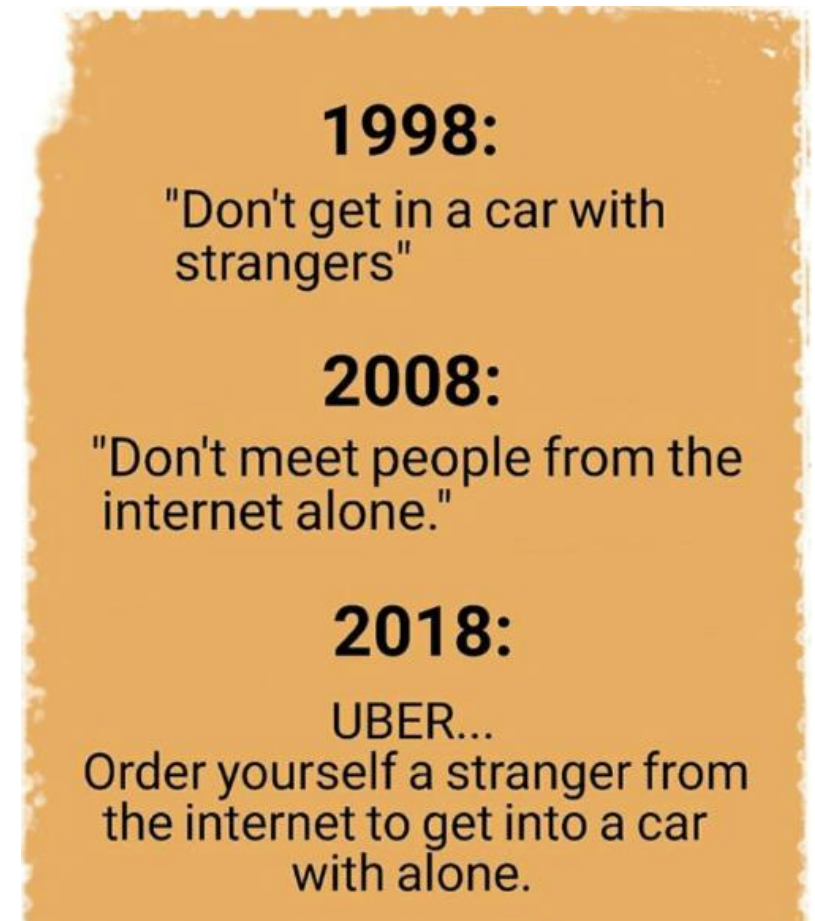


Image source: me.me

Your Life on the Cloud

But

- Corporate leakage to government (security, tax, etc.)
- What if you don't have rights to access/delete our own data?
- Security and privacy breaches
- What if we've changed our ways?
- The department of pre-crime
- Corporate mergers

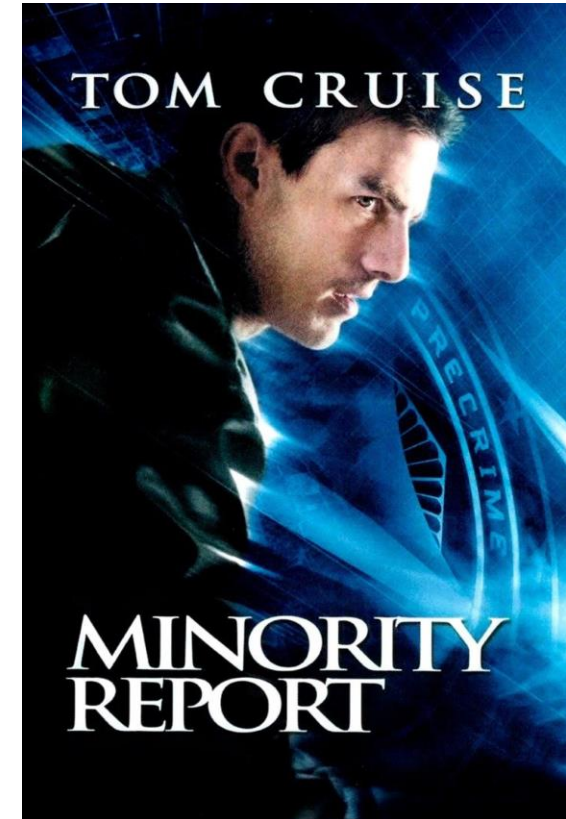


Image source: imdb

Your Life on the Cloud

Social Scoring (<https://www.youtube.com/watch?v=xuqbx8tyW1Y>)



Social Good

[Data Science for Social Good](#) movement training data scientists to support community and charity.

Fight accidents



Fight disease



Teradata University Network

Trending Topics

Hands-On

Must See

Preparing Students and Educators for Success in a Data-Driven World

Big Data / AI

Analytics

Big Data / AI

Data Science

Data Visualization

Database

Marketing

BENEFITS FOR:



Faculty Members



College Students



High School Students



Working Professionals

We believe in change for good
through the actions of many

Simply Giving connects people with causes to help them
make a bigger difference!



simplygiving.com

Helping people make a difference

OUR STORY

At SimplyGiving, we're inspired by the power of the Internet to do good. We provide an incredible way for anyone to fundraise for their favourite charity or a

Health Care Futurology

Some areas where significant impact is to be made in

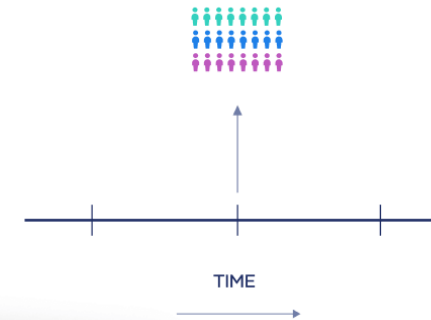
- Your stomach can be instrumented to assess contents, nutrients, etc.
- Your bloodstream can be instrumented too assess insulin levels, etc.
- Your “health” dashboard can be online and shared by your GP
- Health management organisations (HMO) tying funding levels to patient care performance
- GP/HMO will know about your ice cream/beer binge last night and you missing your morning run

Longitudinal studies feasible

- Longitudinal studies is a method in which data is gathered for the same subjects repeatedly over a period of time

Cross-sectional study

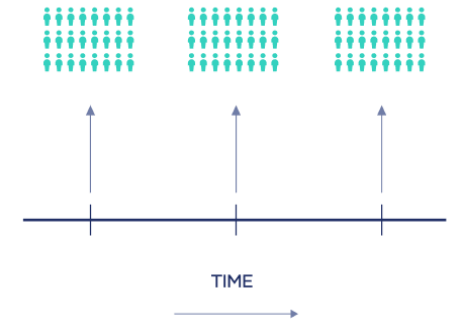
Data collected at one point in time



 Scribbr

Longitudinal study

Data collected repeatedly over time



"Big Data - 2020 Vision"

John Schitka in Strata + Hadoop 2014



<https://www.youtube.com/watch?v=kque3iHbkcl>

Early Innovation (1760s-1900s)= European Inventions

1768 = First Self-Propelled Road Vehicle (Cugnot, France)



1876 = First 4-stroke engine (Otto, Germany)



1886 = First gas-powered, 'production' vehicle (Benz, Germany)



1888 = First four-wheeled electric car (Flocken, Germany)



Streamlining (1910s-1970s)= American Leadership

1910s = Model T / Assembly Line (Ford)



1920s-1930s = Car as Status Symbol...
Roaring '20s / First Motels

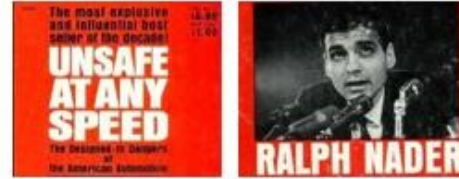


1950s = Golden Age...
Interstate Highway Act (1956)...
8 of Top 10 in Fortune 500 in Cars or Oil (1960)



Modernization (1970s-2010s)= Going Global / Mass Market

1960s = Ralph Nader / Auto Safety



1970s = Oil Crisis / Emissions Focus



1980s = Japanese Auto Takeover Begins...



1990s - 2000s = Industry Consolidation;
Asia Rising;
USA Hybrid Fail (Prius Rise)
DAIMLERCHRYSLER



Late 2000s = Recession / Bankruptcies / Auto Bailouts

Re-Imagination (Today) USA Rising Again

DARPA Challenge (2007, 2012, 2015)
Autonomy Inflection



Today=



+



+



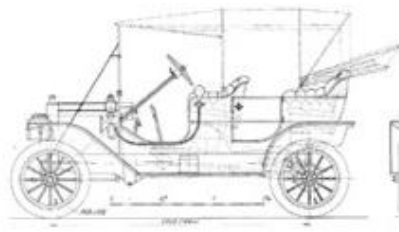
Car Industry Evolution (1760s – Today):

Driven by Innovation & Globalization

Pre-1980s

Analog / Mechanical

Used switches / wiring to route feature controls to driver



Today = Complex Computing

Up to 100 Electronic Control Units / car...
Multiple bus networks per car (CAN / LIN / FlexRay / MOST)...
Drive by Wire...



Car Computing Evolution Since Pre-1980s:

1980s (to Present) **CAN Bus (Integrated Network)**

New regulatory standards drove need to monitor emissions in real time, hence central computer



Today = Smart / Connected Cars

Embedded / tethered connectivity...
Big Tech = New Tier 1 auto supplier
(CarPlay / Android Auto)...



Mechanical / Electrical → Simple Processors → Computers

1990s (to Present) **OBD (On-Board Diagnostics) II**

Monitor / report engine performance; Required in all USA cars post-1996



1990s-2010s **Feature-Built Computing + Early Connectivity**

Automatic cruise control...
Infotainment... Telematics... GPS / Mapping...



Tomorrow = Computers Go Mobile?...

Central hub / decentralized systems?
LIDAR...
Vehicle-to-Vehicle (V2V) /
Vehicle-to-Infrastructure (V2I) /
5G...
Security software...

"The Box" (Brooks & Bone)



Automobile Futurology

“Big data – 2020 vision” talk by SAP manager John Schitka

Self driving cars:

- How does the city replace traffic fine revenue?
- Can you drink and drive if the car is automatic?
- What happens to the taxi industry?
- What happens to the auto insurance industry?
- What happens to people still “self” driving, and their insurance?
- For the Ultimate Driving Machine, how will self-driving cars impact it?

Business Models with Data

ePub Section 2.3



What kinds of businesses do we have operating in the Data Science world?

Business Models

From Wikipedia:

- A [business model](#) describes the rationale of how an organization creates, delivers, and captures value, in economic, social, cultural or other contexts.

Examples of general classes:

- Retailer versus wholesaler
- Luxury consumer products
- Software vendor
- Service provider

What kinds of businesses
do we have operating in
the Data Science world?

Business Models for Data Science

Many Data Science companies fit into traditional IT business models.

- Software as a service (SaaS)
- Consulting
- Customer relationship management

*What are some
business models
specific to data
science?*

For example:

- SAS is both a software vendor and a consultancy, both traditional IT business models
- But there are business models somewhat unique to data-based businesses like data science.







Amazon is providing online infrastructure for online retailers.



Gaming accessories




Headsets Keyboards







Computer mice Chairs


"Alexa, turn on the lights."



Shop by Category

| | |
|--|--|
|  Computers & Accessories |  Video Games |
|  Baby |  Toys & Games |


AmazonBasics



Sign in for the best experience

[Sign in securely](#)

We ship over 45 million products around the world





**Amazon's
infrastructure,
which includes
packing.**

Amazon's smart warehouse



**And shipping
(logistics)**



**And shipping
(prime air)**



amazon

Amazon.com



- An assembly line for the retail industry, with support for embedded online retailers.
- Huge stock of books, DVDs, CDs, etc. *easily searchable*
- Extensive customer *reviews*

Amazon.com

Information-based differentiation:

Satisfies customers by providing a differentiated service:

- Superior information including **reviews** about products
- Superior **range**

Information-based delivery network:

They deliver information for others; retailers in the Amazon marketplace get:

- Customers **directed** to them
- Other retailers' support

Data Business Models

- Information **brokering** service:
 - Buys and sells data/information for others.
- **Information-based differentiation**:
 - Satisfies customers by providing a differentiated service built on the data/information. (www.amazon.com)
- Information-based delivery network:
 - Deliver data/information for others. (www.reuters.com)
(www.plentisoft.com)
- Information provider:
 - Business selling the data/information it collects. (www.Nielsen.com)

[“What a Big-Data Business Model Looks Like”](#) by Ray Wang in the Harvard Business Review claims these are unique in the data world.

Home Activities

Suggested Activities for the week

Videos

Watch [John Schitka “Big Data – 2020 Vision”](#)

From [Data Analytics Handbook \(Pt1\)](#) read the interviews of

- Abraham Cabangbang (2 pp) (pp 5 - 7)
- Ben Bregman (2 pp) (pp 13 - 15)
- Leon Rudyak (3 pp) (pp 16 - 19)



Recap: Learning Outcomes

Week 2

By the end of this week you should be able to:

- Comprehend essentials for coding in Python for data science
- Explain and interpret given **Python** codes
- Explain **different data science roles** and skills and comprehend the differences between them. (Read Analyzing the Analyzers, Harris, Murphy and Vaisman, 2013)
- Explain **Impact** of data science. (Data in cloud, Social good, Futurology)
- Explain the **data business models** for organizations, e.g. how amazon.com uses data to be more competitive.