

# **FIT1043 Introduction to Data Science**

Week 4

Ian K T Tan

School of Information Technology  
Monash University Malaysia

*With materials from Wray Buntine, Mahsa Salehi*

# Learning Outcomes

## Week 4

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

- *Explain open data and linked open data*
- *Explain how to access to new data sources through APIs*
- *Identify how different APIs work*
- **Inspect data quality problems in datasets and recommend solutions to fix them**
- *Use data wrangling operations in Python*

# Data Wrangling

Manipulating data to make it directly usable for analysis

## Week 4

- Common types of data that are at our disposal.
- Data wrangling and the goals of wrangling
- Data quality
- Data auditing
- Examples

# Data

```
access.log
108840 ss/ektlab/Locale/Makstekt.pod HTTP/1.1" 404 538 "-" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108841 ss/ektlab/LMP/media.types HTTP/1.1" 404 535 "-" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108842 ss/mt-static/images/ban-bq.gif HTTP/1.1" 404 538 "-" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108843 cas/ HTTP/1.1" 404 518 "http://192.168.203.11/drupal/modules/acquia/fivestars/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108844 ss/mt-static/images/decrease.gif HTTP/1.1" 404 538 "-" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108845 ss/pub/TWiki/ATasteOfTWiki/BT_logo.gif HTTP/1.1" 404 544 "-" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108846 cas/ HTTP/1.1" 404 518 "http://192.168.203.11/drupal/modules/acquia/fivestars/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108847 filefield_token.module HTTP/1.1" 200 1930 "http://192.168.203.11/drupal/modules/acquia/filefield/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108848 ivestor.info?naxxx=NSFTW HTTP/1.1" 200 614 "http://192.168.203.11/drupal/modules/acquia/fivestars/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108849 ivestor.info?naxxx="22--%3E%3C/style%3E%3C/script%3E%3Cscript%3Enetsparker(0x0037E2)%3C/script%3E HTTP/1.1" 200 61
108850 ivestor.info?naxxx="2BNSFTW%2B" HTTP/1.1" 200 614 "http://192.168.203.11/drupal/modules/acquia/fivestars/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108851 ivestor.info HTTP/1.1" 200 614 "" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108852 ivestor.info HTTP/1.1" 200 614 "" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108853 cas/ HTTP/1.1" 404 518 "http://192.168.203.11/drupal/modules/acquia/fivestars/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108854 ivestor.info HTTP/1.1" 200 614 "http://netsparker.com/p" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108855 ivestor.info?22--%3E%3C/style%3E%3C/script%3E%3Cscript%3Enetsparker(0x0037E4)%3C/script%3E HTTP/1.1" 404 570 "http://192.168.203.11/drupal/modules/acquia/fivestars/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108856 ivestor.info?http://www.netsparker.com? HTTP/1.1" 200 614 "http://192.168.203.11/drupal/modules/acquia/fivestars/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108857 ivestor.info?http://netsparker.com/p HTTP/1.1" 200 614 "http://192.168.203.11/drupal/modules/acquia/fivestars/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108858 ivestor.info?naxxx="0d%0ans%3anetsparker056650%3dvuln HTTP/1.1" 200 614 "http://192.168.203.11/drupal/modules/acquia/fivestars/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108859 ivestor.info?naxxx="0ans%3anetsparker056650%3dvuln HTTP/1.1" 200 614 "http://192.168.203.11/drupal/modules/acquia/fivestars/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108860 ivestor.info?22ns=22netsparker(0x0037E5) HTTP/1.1" 404 548 "http://192.168.203.11/drupal/modules/acquia/fivestars/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108861 ivestor.info HTTP/1.1" 200 614 "%0d%0ans%3anetsparker056650%3dvuln" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108862 ivestor.info?2522ns%253D%2522netsparker%25280x0037E6%2529 HTTP/1.1" 404 548 "http://192.168.203.11/drupal/modules/acquia/fivestars/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108863 ivestor.info HTTP/1.1" 200 614 "%0ans%3anetsparker056650%3dvuln" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108864 ivestor.info?0d%0ans%3anetsparker056650%3dvuln HTTP/1.1" 404 544 "http://192.168.203.11/drupal/modules/acquia/fivestars/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108865 ivestor.info?"" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108866 ivestor.info?0ans%3anetsparker056650%3dvuln HTTP/1.1" 404 542 "http://192.168.203.11/drupal/modules/acquia/fivestars/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108867 ivestor.install?naxxx=NSFTW HTTP/1.1" 200 5998 "http://192.168.203.11/drupal/modules/acquia/fivestars/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108868 filefield_widget.inc HTTP/1.1" 200 14909 "http://192.168.203.11/drupal/modules/acquia/filefield/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108869 ivestor.install?naxxx="2BNSFTW%2B" HTTP/1.1" 200 5998 "http://192.168.203.11/drupal/modules/acquia/fivestars/" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108870 /filefield_field.inc.bak HTTP/1.1" 404 262 "-" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108871 /filefield.theme.inc.bak HTTP/1.1" 404 262 "-" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
108872 /filefield_field.inc HTTP/1.1" 404 262 "-" "Mozilla/4.0 (compatible: MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4342.2546)"
```

# Databases

## Relational Databases

- Oracle
- MySQL
- MariaDB
- MS-SQL

## Containing data for

- Customer relationship management
- Mortgage, hire purchase, business loan system
- Salesforce automation
- Credit card system, ATM transactions, retail banking
- Human resource and payroll





# Files

## Examples

- ***System log files***
- Spreadsheets
- PDF files
- Image files
- Raw text files
- Formatted text files



# System Log Files

```
access.log
108840 22/exllib/Locale/Maketext.pod HTTP/1.1" 404 538 "-" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV
108841 22/exllib/LWP/media.types HTTP/1.1" 404 535 "-" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1; .3
108842 22/mt-static/images/ban-bg.gif HTTP/1.1" 404 538 "-" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV
108843 22/ HTTP/1.1" 404 518 "http://192.168.203.11/drupal/modules/acquia/fivestar/" "Mozilla/4.0 (compatible; MS
108844 22/mt-static/images/decrease.gif HTTP/1.1" 404 539 "-" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; S
108845 22/pub/TWiki/ATasteOfTWiki/BY_logo.gif HTTP/1.1" 404 544 "-" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT
108846 22/ HTTP/1.1" 404 518 "http://192.168.203.11/drupal/modules/acquia/fivestar/" "Mozilla/4.0 (compatible; MSIE
108847 filefield_token.module HTTP/1.1" 200 1930 "http://192.168.203.11/drupal/modules/acquia/filefield/" "Mozilla/4
108848 ivestar.info?nsxxt=NSFTW HTTP/1.1" 200 614 "http://192.168.203.11/drupal/modules/acquia/fivestar/" "Mozilla/
108849 ivestar.info?nsxxt=%22--%3E%3C/style%3E%3C/script%3E%3Cscript%3Enetsparker(0x0037E2)%3C/script%3E HTTP/1.1"
108850 ivestar.info?nsxxt=%2BNSFTW%2B' HTTP/1.1" 200 614 "http://192.168.203.11/drupal/modules/acquia/fivestar/" "M
108851 ivestar.info HTTP/1.1" 200 614 "'\`--></style></script><script>netsparker(0x0037E3)</script>" "Mozilla/4.0 (co
108852 ivestar.info HTTP/1.1" 200 614 "'%2BNSFTW%2B'" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1; .NET CL
108853 22/ HTTP/1.1" 404 518 "http://192.168.203.11/drupal/modules/acquia/fivestar/" "Mozilla/4.0 (compatible; MSIE 6
108854 ivestar.info HTTP/1.1" 200 614 "http://netsparker.com/n" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV
108855 ivestar.info?%22--%3E%3C/style%3E%3C/script%3E%3Cscript%3Enetsparker(0x0037E4)%3C/script%3E HTTP/1.1" 404 570 "h
108856 ivestar.info?http://www.netsparker.com? HTTP/1.1" 200 614 "http://192.168.203.11/drupal/modules/acquia/fivestar/
108857 ivestar.info?http://netsparker.com/n HTTP/1.1" 200 614 "http://192.168.203.11/drupal/modules/acquia/fivestar/"
108858 ivestar.info?nsxxt=%0d%0ans%3anetsparker056650%3dvuln HTTP/1.1" 200 614 "http://192.168.203.11/drupal/modules/ac
108859 ivestar.info?nsxxt=%0ans%3anetsparker056650%3dvuln HTTP/1.1" 200 614 "http://192.168.203.11/drupal/modules/acqui
108860 ivestar.info/%22ns=%22netsparker(0x0037E5) HTTP/1.1" 404 548 "http://192.168.203.11/drupal/modules/acquia/fivestar
108861 ivestar.info HTTP/1.1" 200 614 "%0d%0ans%3anetsparker056650%3dvuln" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT
108862 ivestar.info/%2522ns%253D%2522netsparker%25280x0037E6%2529 HTTP/1.1" 404 548 "http://192.168.203.11/drupal/modules
108863 ivestar.info HTTP/1.1" 200 614 "%0ans%3anetsparker056650%3dvuln" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1
108864 ivestar.info%0d%0ans%3anetsparker056650%3dvuln HTTP/1.1" 404 544 "http://192.168.203.11/drupal/modules/acquia/fives
108865 ivestar.info?'\`--></style></script><script>netsparker(0x0037E7)</script> HTTP/1.1" 200 614 "http://192.168.203.11/c
108866 ivestar.info%0ans%3anetsparker056650%3dvuln HTTP/1.1" 404 542 "http://192.168.203.11/drupal/modules/acquia/fivestar/
108867 ivestar.install?nsxxt=NSFTW HTTP/1.1" 200 5598 "http://192.168.203.11/drupal/modules/acquia/fivestar/" "Mozilla/4.0
108868 filefield_widget.inc HTTP/1.1" 200 14509 "http://192.168.203.11/drupal/modules/acquia/filefield/" "Mozilla/4.0 (comp
108869 ivestar.install?nsxxt=%2BNSFTW%2B' HTTP/1.1" 200 5598 "http://192.168.203.11/drupal/modules/acquia/fivestar/" "Mos
108870 /filefield_field.inc.bak HTTP/1.1" 404 262 "-" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1
108871 /filefield.theme.inc.bak HTTP/1.1" 404 262 "-" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1
108872 /filefield field.inc= HTTP/1.1" 404 262 "-" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.43
4
```



# Web & Crowd Sources

Open Data, REST API, Web Scrapping, etc.

## Huge!

- What can we do with Web and other crowd sourced data?
- What type of web data?

## Data on the Web

- News: Sports, Finance, World, Nation
- Blogs
- Corporates: Products & Services
- Government



## Watching Our World Unfold

### A Global Database of Society

Supported by [Google Jigsaw](#), the GDELT Project monitors the world's broadcast, print, and web news from nearly every corner of every country in over 100 languages and identifies the people, locations, organizations, counts, themes, sources, emotions, counts, quotes, images and events driving our global society every second of every day, creating a free open platform for computing on the entire

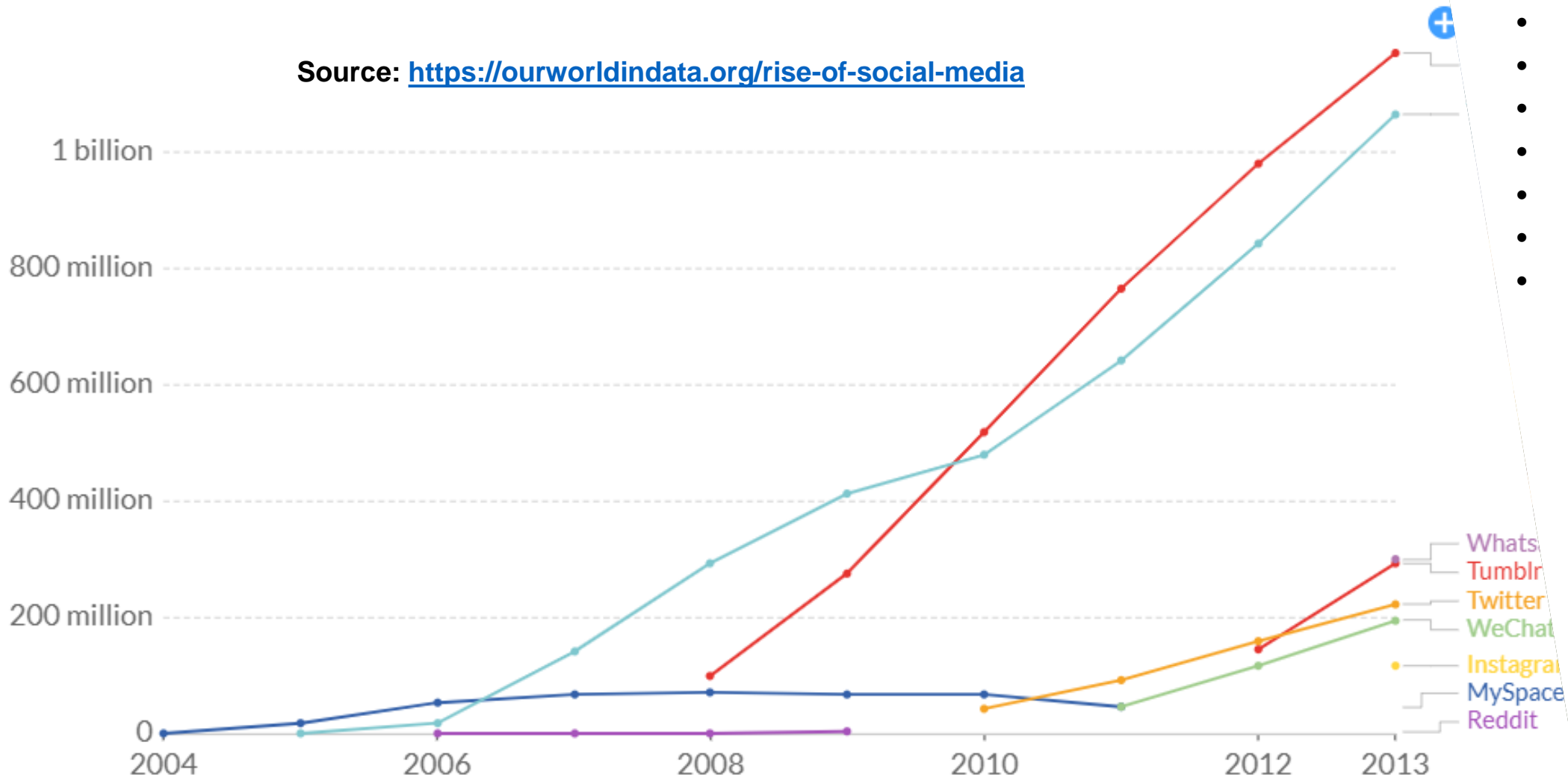
Apr 20, 2014



# Number of people using social media platforms

Estimates correspond to monthly active users (MAUs). Facebook, for example, measures MAUs as users that have logged in during the past 30 days. See source for more details.

Source: <https://ourworldindata.org/rise-of-social-media>

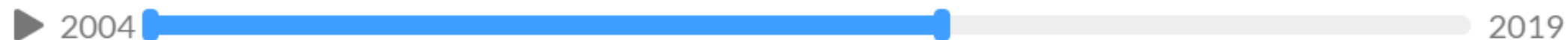


## Social Media

- *Twitter*
- *Facebook*
- *LinkedIn*
- *Instagram*
- *Pinterest*
- *TikTok*
- *YouTube*

Source: Statista and TNW (2019)

CC BY

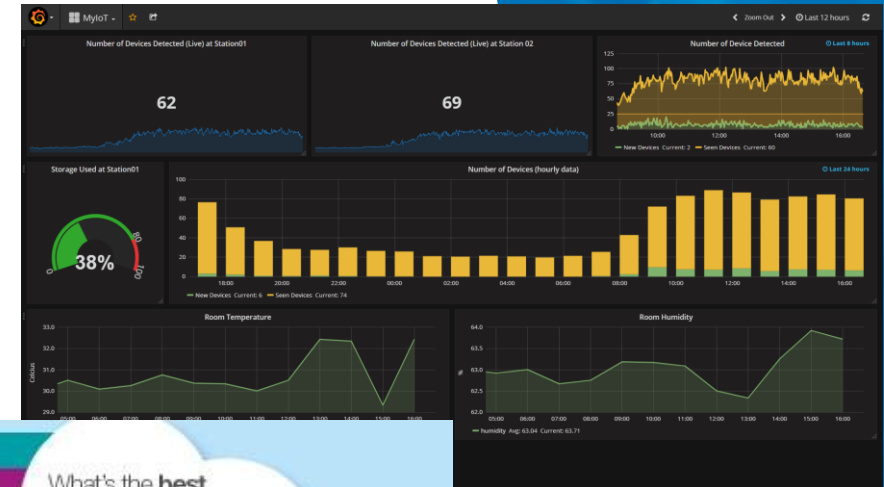


# Internet of Things (IoT)

## The Future

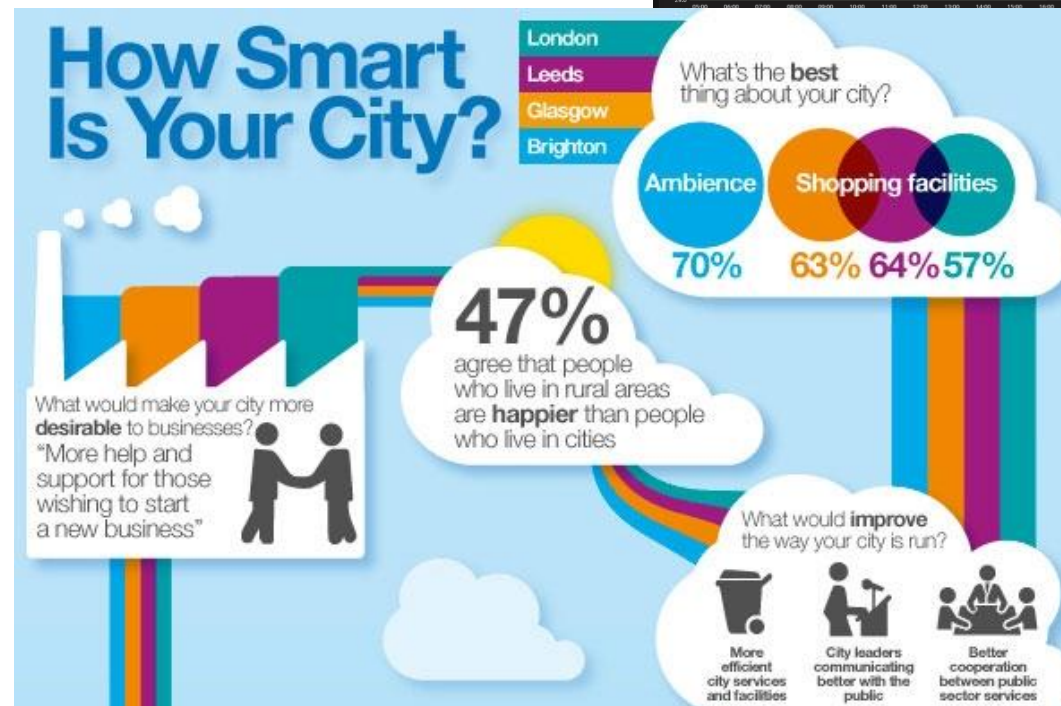
### Where?

- Utilities – Water, Power, Street Lighting, Traffic Lights, Waste Disposal
- Vehicle to Vehicle, Autonomous Driving
- Monitoring Systems



### Mobile Phone Data

- Location
- Browsing history
- Usage history
- Personal information



# Data Wrangling





	C	D	E	F	G	H	I	J	K	L
	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-famil	White	Male	2174	
no	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0
	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-famil	White	Male	0	0
	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0	0
	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0	0
	284582	Masters	14	Married-civ-spouse	Exec-managerial	Wife	White	Female	0	0
	160187	9th	5	Married-spouse-absen	Other-service	Not-in-famil	Black	Female	0	0
no	209642	HS-grad	9	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0
	45781	Masters	14	Never-married	Prof-specialty	Not-in-famil	White	Female	14084	0
	159449	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	5178	0
	280464	Some-colleg	10	Married-civ-spouse	Exec-managerial	Husband	Black	Male	0	0
	141297	Bachelors	13	Married-civ-spouse	Prof-specialty	Husband	Asian-Pac-Is	Male	0	0
	122272	Bachelors	13	Never-married	Adm-clerical	Own-child	White	Female	0	0
	205019	Assoc-acdm	12	Never-married	Sales	Not-in-famil	Black	Male	0	0
	121772	Assoc-voc	11	Married-civ-spouse	Craft-repair	Husband	Asian-Pac-Is	Male	0	0
	245487	7th-8th	4	Married-civ-spouse	Transport-moving	Husband	Amer-Indian	Male	0	0

Ideal Data



## REPORT:

Clinical note: transformed AML. Ongoing fevers.? Source. ? fungal infection

## Report:

Axial 1.25 mm slices at 10 mm intervals taken in inspiration with selected images in the prone position.

No mediastinal or hilar lymphadenopathy. Heart size is normal. Borderline enlargement of the main pulmonary outflow tract. There is smooth interlobular septal thickening throughout both lungs, which may be secondary to fluid overload. There is a background of emphysematous changes, predominantly in the upper lobes. A 5 x 8 mm nodule is identified in the right upper lobe (image 10). It is well-circumscribed with no evidence of surrounding ground-glass opacity. No calcification or cavitation of this lesion. The visualised portions of the liver and spleen appear normal, allowing for lack of intravenous contrast.

## Conclusion:

Single nodule in right upper lobe has a non-specific appearance but given the clinical history, this could represent a focus of fungal infection.

Reported by: Dr. Christopher J. Lee

PJL/PJL

A1.2f

Result type: CT Chest Hi Resolution

Result date: 11 January 2005 12:21

Result status: Auth (Verified)

Result title: CTCHEHI

Data Reality

# Why Wrangle?

## Working with Raw Data is challenging

- Data comes in all shapes and sizes
- Different files have different formatting
- Mistakes in data entries

***We need techniques to cleanse and prepare the data!***

# Goals of Data Wrangling

Transform Data into Usable Data

Raw data  $\Rightarrow$  **Data Wrangling**  $\Rightarrow$  Tidy data  $\Rightarrow$  Data Analysis  $\Rightarrow$  Knowledge

*Data + Wrangling + Analysis = Data Product*



# What is Data Wrangling?

*Process of **transforming** “raw” data into data that can be **analysed** to generate valid **actionable results** and **insights***

## Steps

- Data pre-processing
- Data preparation
- Data cleansing
- Data transformation
- Etc.

**Data Quality**



# Source of Data Quality Issues

## Causes

- Interpretability issue
- Data format issue
- Inconsistent and faulty data
- Missing and incomplete data
- Outliers
- Duplicates

# Interpretability Issue

## Data quality problems

- Is there a proper documentation about the data?
  - Without proper documentation (i.e., a [data dictionary](#)), it is not possible for us to use the synthetic data.
- We might be able to guess the meaning of each column,
  - But in general, we need a data dictionary to explain the fields

*Is this data interpretable?*

```
32,1,1,95,0,?,0,127,0,.7,1,?,?,1
34,1,4,115,0,?,?,154,0,.2,1,?,?,1
35,1,4,?,0,?,0,130,1,?,?,?,7,3
36,1,4,110,0,?,0,125,1,1,2,?,6,1
38,0,4,105,0,?,0,166,0,2.8,1,?,?,2
38,0,4,110,0,0,0,156,0,0,2,?,3,1
38,1,3,100,0,?,0,179,0,-1.1,1,?,?,0
38,1,3,115,0,0,0,128,1,0,2,?,7,1
38,1,4,135,0,?,0,150,0,0,?,?,3,2
38,1,4,150,0,?,0,120,1,?,?,?,3,1
40,1,4,95,0,?,1,144,0,0,1,?,?,2
```



# Data format issue

## Data quality problems

- Data from different sources often have different data formats and are generated from different processes.
- It is challenging to integrate and manipulate data in different formats.

```
1 {  
2   "meta" : {  
3     "view" : {  
4       "id" : "tdvh-n9dv",  
5       "name" : "Melbourne bike share",  
6       "attribution" : "City of Melbourne, Australia",  
7       "averageRating" : 0,  
8       "category" : "Transport & Movement",  
9       "createdAt" : 1428898164,  
10      "description" : "Melbourne Bike Share is a joint RACV/Victoria  
11      "displayType" : "table",
```

JavaScript Object Notation (JSON)

```
<response>  
  <row>  
    <row _id="155" _uuid="7C09387D-9E6C-4B42-9041-9A98B88F54  
      <id>2</id>  
      <featurename>Harbour Town - Docklands Dve - Dockland  
      <terminalname>60000</terminalname>  
      <nbbikes>9</nbbikes>  
      <nemptydoc>14</nemptydoc>  
      <uploaddate>1453986006</uploaddate>  
      <coordinates human_address="{&quot;address&quot;:&qu
```

Extensible Mark-up Language (XML)

# Inconsistent and faulty data

## Data quality problems

- Mistyped data
- Inconsistent entry
- Extraneous data
- Etc.

*Mark Johnson, 31, 21/Aug/1985, 180, M, 0433010010, Melbourne VIC*

*Mr. Christian, Peter, 34, 21-09-1982,, M, 0433010118, Sydney NSW*

*Ethan Steedman, 32, 01/01/1982, 170, M, 0433210019, Sydney NSW*

# Missing values

## Data quality problems

- Data values that should be presented in a dataset but that are absent for many reasons.

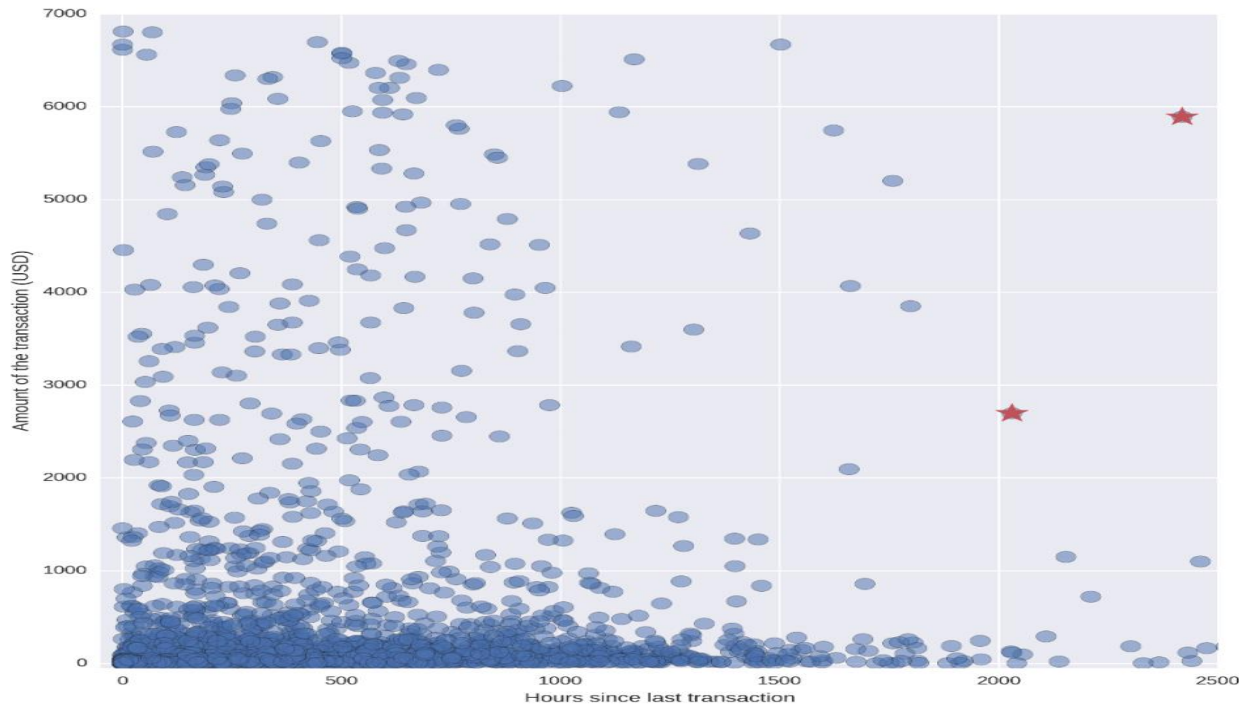
*Missing values in the Switzerland heart disease data set are indicated by “?”.*

```
32,1,1,95,0,?,0,127,0,.7,1,?,?,1
34,1,4,115,0,?,?,154,0,.2,1,?,?,1
35,1,4,?,0,?,0,130,1,?,?,?,7,3
36,1,4,110,0,?,0,125,1,1,2,?,6,1
38,0,4,105,0,?,0,166,0,2.8,1,?,?,2
38,0,4,110,0,0,0,156,0,0,2,?,3,1
38,1,3,100,0,?,0,179,0,-1.1,1,?,?,0
38,1,3,115,0,0,0,128,1,0,2,?,7,1
38,1,4,135,0,?,0,150,0,0,?,?,3,2
38,1,4,150,0,?,0,120,1,?,?,?,3,1
40,1,4,95,0,?,1,144,0,0,1,?,?,2
```

# Outliers

## Data quality problems

- An observation that lies in an abnormal distance from the majority of the other observations in the dataset.





# Duplicates

## Data quality problems

- Multiple data entries that correspond to the same piece of information

*Christoph Cleveland, 20, 10-10-1996, 50, M, 0433550210, Hobart TAS*

*Chris. Cleveland, 20, 10-10-1996, 176, M, 0433550210, Hobart TAS*

# Data Auditing



# Initial Data Auditing

Given a data set, what are the common initial auditing steps you would conduct?

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	who	adult_male	deck	embark_town	alive	alone	name
2	0	3	male	22	1	0	7.25	S	Third	man	TRUE		Southampton	no	FALSE	Braund, Mr. Owen
3	1	1	female	38	1	0	71.2833	C	First	woman	FALSE	C	Cherbourg	yes	FALSE	Cumings, Mrs. J.
4	1	3	female	26	0	0	7.925	S	Third	woman	FALSE		Southampton	yes	TRUE	Heikkinen, Miss.
5	1	1	female	35	1	0	53.1	S	First	woman	FALSE	C	Southampton	yes	FALSE	Futrelle, Mrs.
6	0	3	male	35	0	0	8.05	S	Third	man	TRUE		Southampton	no	TRUE	Allen, Mr. William
7	0	3	male		0	0	8.4583	Q	Third	man	TRUE		Queenstown	no	TRUE	Moran, Mr. James
8	0	1	male	54	0	0	51.8625	S	First	man	TRUE	E	Southampton	no	TRUE	McCarthy, Mr. Thomas
9	0	3	male	2	3	1	21.075	S	Third	child	FALSE		Southampton	no	FALSE	Palsson, Master G.
10	1	3	female	27	0	2	11.1333	S	Third	woman	FALSE		Southampton	yes	FALSE	Johnson, Mrs. Oscar
11	1	2	female	14	1	0	30.0708	C	Second	child	FALSE		Cherbourg	yes	FALSE	Nasser, Mrs. Nicholas
12	1	3	female	4	1	1	16.7	S	Third	child	FALSE	G	Southampton	yes	FALSE	Sandstrom, Miss. Margareta
13	1	1	female	58	0	0	26.55	S	First	woman	FALSE	C	Southampton	yes	TRUE	Bonnell, Miss.
14	0	3	male	20	0	0	8.05	S	Third	man	TRUE		Southampton	no	TRUE	Saunderscock, Mr. John
15	0	3	male	39	1	5	31.275	S	Third	man	TRUE		Southampton	no	FALSE	Andersson, Mr. Olof
16	0	3	female	14	0	0	7.8542	S	Third	child	FALSE		Southampton	no	TRUE	Vestrom, Miss.
17	1	2	female	55	0	0	16	S	Second	woman	FALSE		Southampton	yes	TRUE	Hewlett, Mrs. Mary
18	0	3	male	2	4	1	29.125	Q	Third	child	FALSE		Queenstown	no	FALSE	Rice, Master. Eugene
19	1	2	male		0	0	13	S	Second	man	TRUE		Southampton	yes	TRUE	Williams, Mr. Charles
20	0	3	female	31	1	0	18	S	Third	woman	FALSE		Southampton	no	FALSE	Vander Planke, Mrs. James
21	1	3	female		0	0	7.225	C	Third	woman	FALSE		Cherbourg	yes	TRUE	Masselmani, Mrs. Ines
22	0	3	male	25	0	0	26	S	Second	man	TRUE		Southampton	no	TRUE	Runyan, Mr. Jacob

# Data Auditing

Some general steps to perform data auditing (assuming you're given or have read the data into a DataFrame `df`)

- Dimension of the data: `df.shape()`
  - Number of rows
  - Number of columns

`(8500, 10)`

# Data Auditing

Head and tail rows: `df.head()` ; `df.tail()`

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	who	adult_male	deck	embark_town	alive	alone	name
0	0	3	male	22.0	1	0	7.2500	S	Third	man	True	NaN	Southampton	no	False	Braund, Mr. Owen Harris
1	1	1	female	38.0	1	0	71.2833	C	First	woman	False	C	Cherbourg	yes	False	Cumings, Mrs. John Bradley (Florence Briggs Th...
2	1	3	female	26.0	0	0	7.9250	S	Third	woman	False	NaN	Southampton	yes	True	Heikkinen, Miss. Laina
3	1	1	female	35.0	1	0	53.1000	S	First	woman	False	C	Southampton	yes	False	Futrelle, Mrs. Jacques Heath (Lily May Peel)
4	0	3	male	35.0	0	0	8.0500	S	Third	man	True	NaN	Southampton	no	True	Allen, Mr. William Henry



# Data Auditing

What are the numerical  
and categorical columns?

Check basic information about the dataframe

- Number of records,
- Whether there are null values,
- Etc

datatype: `df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 892 entries, 0 to 891
Data columns (total 16 columns):
survived      892 non-null int64
pclass        892 non-null int64
sex           892 non-null object
age           715 non-null float64
sibsp         892 non-null int64
parch         892 non-null int64
fare          892 non-null float64
embarked      890 non-null object
class         892 non-null object
who           892 non-null object
adult_male    892 non-null bool
deck          204 non-null object
embark_town   890 non-null object
alive         892 non-null object
alone         892 non-null bool
name          892 non-null object
dtypes: bool(2), float64(2), int64(4), object(8)
memory usage: 99.4+ KB
```

# Data Auditing

Check some basic statistics about columns:

- Numerical columns: `df.describe()`
- Object columns: `df.describe(include=[np.object])`

<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.describe.html>

	survived	pclass	age	sibsp	parch	fare
<b>count</b>	892.000000	892.000000	715.000000	892.000000	892.000000	892.000000
<b>mean</b>	0.384529	2.307175	29.720517	0.522422	0.381166	32.201737
<b>std</b>	0.486757	0.836750	14.490914	1.102264	0.805706	49.665589
<b>min</b>	0.000000	1.000000	0.420000	0.000000	0.000000	0.000000
<b>25%</b>	0.000000	2.000000	NaN	0.000000	0.000000	7.917700
<b>50%</b>	0.000000	3.000000	NaN	0.		
<b>75%</b>	1.000000	3.000000	NaN	1.		
<b>max</b>	1.000000	3.000000	80.000000	8.		

	sex	embarked	class	who	deck	embark_town	alive	name
<b>count</b>	892	890	892	892	204	890	892	892
<b>unique</b>	4	3	3	3	7	7	2	891
<b>top</b>	male	S	Third	man	C	Southampton	no	Behr, Mr. Karl Howell
<b>freq</b>	574	644	491	538	60	643	550	2

# Data Auditing

Check correlation amongst variables: `df.corr()`

	<b>survived</b>	<b>pclass</b>	<b>age</b>	<b>sibsp</b>	<b>parch</b>	<b>fare</b>	<b>adult_male</b>	<b>alone</b>
<b>survived</b>	1.000000	-0.339932	-0.079229	-0.035960	0.080874	0.257012	-0.555222	-0.204758
<b>pclass</b>	-0.339932	1.000000	-0.367110	0.083789	0.019245	-0.548667	0.092448	0.137065
<b>age</b>	-0.079229	-0.367110	1.000000	-0.301632	-0.185857	0.096148	0.278219	0.196753
<b>sibsp</b>	-0.035960	0.083789	-0.301632	1.000000	0.414985	0.159654	-0.253892	-0.583247
<b>parch</b>	0.080874	0.019245	-0.185857	0.414985	1.000000	0.216221	-0.350200	-0.582176
<b>fare</b>	0.257012	-0.548667	0.096148	0.159654	0.216221	1.000000	-0.181997	-0.271540
<b>adult_male</b>	-0.555222	0.092448	0.278219	-0.253892	-0.350200	-0.181997	1.000000	0.403131
<b>alone</b>	-0.204758	0.137065	0.196753	-0.583247	-0.582176	-0.271540	0.403131	1.000000

# Examples





# [1] Misspelling and Inconsistency

.....	<b>Suburbs</b>
.....	burwood.
.....	springvale
....	Burwood
....	Springvae
....	East Melbourne
.....	E. Melbourne
.....	.....

- *What are the problems?*
- *How can you detect these data problems?*
- *How can you resolve them?*

# [1] Misspelling and Inconsistency

- Inconsistency
  - common cases:
    - upper vs. lower case
    - inconsistency in domain value representation, e.g., 0 vs. No, 1 vs. Yes
  - detecting and fixing
    - investigate unique domain values (`unique()`)
    - make the representation consistent, e.g., replace
- Misspelling
  - investigate unique domain values (`unique()`)
  - string matching
    - calculate domain value frequencies (`value_counts()`)
    - for all values, find matches for the infrequent values
    - replace infrequent values with the best match (if it exists) from the more frequent values.

## [2] Irregularities

.....	<b>Entry Time</b>
.....	12/13/2010
.....	1/1/2014
....	45/2/2010
....	20/3/2011
....	2/14/2014
.....	25/12/2014
.....	.....

- *What are the problems?*
- *How can you detect these data problems?*
- *How can you resolve them?*

# [2] Irregularities

- Common cases:
  - invalid dates
  - domain dependent value, value not valid for a specific domain, e.g., negative value for number of passengers
- Detecting
  - investigate unique domain values (`unique()`)
  - investigate value ranges for the column
  - type casting, e.g., parse date string to datetime object, catch exceptions when it is not a valid date format
- Fixing
  - refer to documentation if it exists, to see whether these values have special meaning
  - replace
  - remove



# [3] Integrity Constraint Violation

....	year_built	time_settled
.....	2010	12/13/2010
.....	2010	1/1/2014
....	2010	45/2/2002
....	2010	20/3/2011
....	2021	2/14/2014
.....	.....	.....

- *What are the problems?*
- *How can you detect these data problems?*
- *How can you resolve them?*

# [3] Integrity Constraint Violation

- Common case:
  - highly dependent on context, e.g.,
    - sold date vs. advertised date,
    - one field is the sum of the other two,
    - land size must be larger than building size, etc.
- Detecting
  - highly dependent on the domain and problems
- Fixing
  - swap
  - remove, etc.

# [4] Duplicates

Index	Name	Gender	D.O.B	Mobile	Address	.....
10	John	N/A	7/9/1985	0412685210	N/A	.....
145	John Walter	M	7/9/1985	0412685210	2 Yale st., burwood	.....
200	John Walter	Male	7/9/1985	0412685210	2 Yale street, burwood	....
268	Walter, John	Male	7/9/1985	0412685210	2 Yale street, burwood	....
450	John Walter	-	1985	-	2 Yale street, burwood	.....
.....	.....		.....			.....

- *What are the duplications in the example and how to detect and resolve these data duplications?*

# [4] Duplicates

- Common cases:
  - complete duplication
  - duplicate due to field missing
    - different record have different piece of info
- Detecting
  - identifying keys to check duplicates
    - fixing other quality issues
    - try different keys
- Fixing
  - combine information/merge
  - remove duplicates



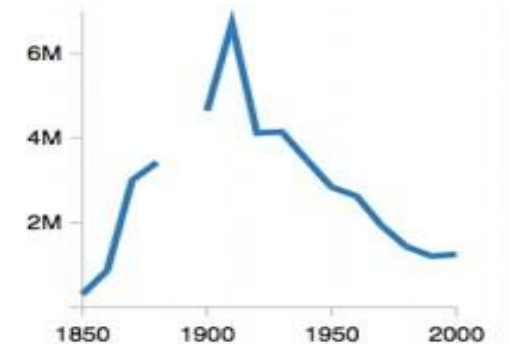
# [5] Missing Values

Adv. Price	Bedrooms	Land (sqm)	Condition	Suburb	Source
800000	2	-	Old	CBD	realestate
80	2	250	*	CBD	domain
1100000	-	-	Fair	Burwood	realestate
*	3	800	*	Dandenong	domain
Contact Agent	-	500	Fair	Burwood	realestate
.....		.....		.....	.....

- *How to detect these missing value records and how to fix them?*

# [5] Missing Values

- Detecting
  - investigate unique domain values (`unique()`)
  - investigate value range, cautious about extremely small and large values
  - domain analysis
- Fixing
  - imputation
    - mean and mode
    - regression (find variables that are closely related (e.g., `df.corr()`)
    - dummy value
  - removal
  - all depends on the situation and needs justification



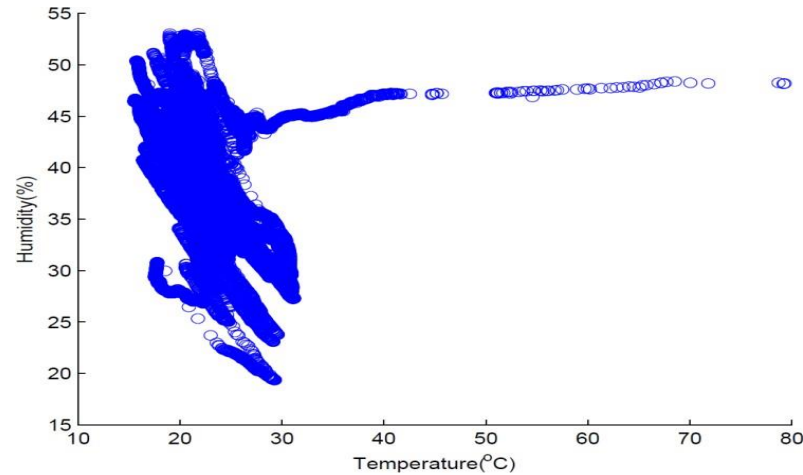
# [6] Outliers

Adv. Price	Bedrooms	Land (sqm)	Condition	Suburb	Source
800000	2	-	Old	CBD	realestate
80	2	250	Fair	CBD	domain
1100000	-	550	Fair	Burwood	realestate
500000	8	800	New	Dandenong	domain
500000	-	500	Fair	Burwood	realestate
.....		.....		.....	.....

- *How to detect these outliers? Can we directly remove outliers once found from tools, e.g., boxplot?*

# [6] Outliers

- Common cases: numerical field
- Challenge:
  - not easy to find
- Detecting
  - range of values `df.describe()`
  - Graphical tools, e.g., boxplot (default using a IQR rule)
  - $3\sigma$  edit rule
  - Good to do some comparison between results found by different identifiers
- Fixing
  - Similar to handling missing values



# Recap: Learning Outcomes

Week 4

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

- Explain open data and linked open data
- Explain how to access to new data sources through APIs
- Identify how different APIs work
- Inspect data quality problems in datasets and recommend solutions to fix them
- *Use data wrangling operations in Python (in Tutorials)*