# DMV
# Data Types

Suzanne Little

# Data types

- Recall: Data is collected information (a working definition)
- Structured vs Unstructured
- Quantitative vs Qualitative
- Discrete vs Continuous
- Four levels of data
- Some special data types to watch for

# Data

Example: a *person* (**object** or **entity** or **instance** or **record** or **row**) has **attributes** (or **features** or **descriptors** or **variables** or **columns**)

- ○ Name
- ○ Passport number
- ○ Birth place
- ○ Eye colour
- ○ Shoe size

# Structured vs Unstructured

tables, organised, observations,

Row is instance, Column is attribute

Examples:
    company records
    scientific observation

Easier for Machine Learning to work with (kinda)

| 1 | Total salaried emp | 1995 | 1996 | 1997 |
|---|---|---|---|---|
| 32 | Chile | 69.40000153 | 70.09999847 | 70.40000153 |
| 33 | Colombia | 66.19999695 | 66.5 | 64.90000153 |
| 34 | Costa Rica | 71.40000153 | 71.19999695 | 69.90000153 |
| 35 | Croatia | | 71.40000153 | 74.09999847 |
| 36 | Cuba | 84 | 84.30000305 | 83.59999847 |

No hierarchy or arrangement

Raw signals that need processing

Examples:
    tweets & social media posts
    server logs
    media (images, video, etc)

More challenging to work with. How to turn into "Structured"?

**DCU** ✔
@DublinCityUni

Following

Wishing all of our new and returning students the very best of luck on their first day of lectures! 📊📝

1:28 AM - 24 Sep 2018

13 Retweets 71 Likes

# Special types of data to watch for

- Temporal (or Time Series)
- Geographic (or Spatial)
- Documents, Images, Video, Audio, 3D
- "Raw" data - unstructured and (sometimes) incidental

# Qualitative          vs          Quantitative

| | |
|---|---|
| Quality, Label, Trait | Quantity, Measurement |
| Categorical | Numerical |
| Limited mathematical functions | "All the maths!" (well most) |
| Examples: | Examples: |
|     Country of origin |     Shoe size |
|     Gender |     Temperature |
|     Favourite Colour |     Bank balance |

# Quantitative

## Discrete               vs               Continuous

only certain values are valid

ie: there are gaps

usually from counting

Examples:
    Number of times attended
    Number of crimes reported

theoretically any value is possible

depends on measuring device ability

usually from measurements

Examples:
    Cholesterol level
    Time required to complete task

# Data types

Structured vs Unstructured

Quantitative vs Qualitative

Discrete vs Continuous

Four levels of data measurement

1. Nominal
2. Ordinal
3. Interval
4. Ratio

# NOIR (Stanley Stevens)

<div style="writing-mode: vertical">Categorical</div>

**Nominal** (name, label, category)

Gender, Department, Language

Not described by numbers
No maths except equality & set membership
mode but not mean or median

<div style="writing-mode: vertical">Qualitative</div>

**Ordinal** (labels plus order)

Temperature (very hot, hot, warm, mild)
Medals (Gold, Silver, Bronze), Scale (Likert - 1
to 10), colour
Can be arranged by order but not added or
subtracted, median but not mean

<div style="writing-mode: vertical">Measurement</div>

**Interval** (numbers with proportionate spaces)

We can now talk about "difference" (+/-)

Income, Shoe size,
Temperature ($^o$C, $^o$F)

"defined interval between values but lacks a zero
point"

<div style="writing-mode: vertical">Quantitative</div>

**Ratio** (also numbers but with zero)

Can now multiply & divide

Age, Amount of rainfall, Book sales,
Temperature (in Kelvin), [normally counting]

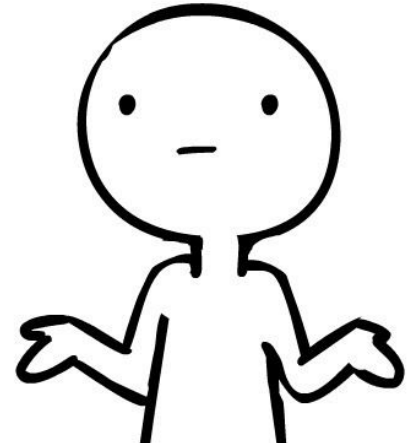Zero has meaning - no negatives
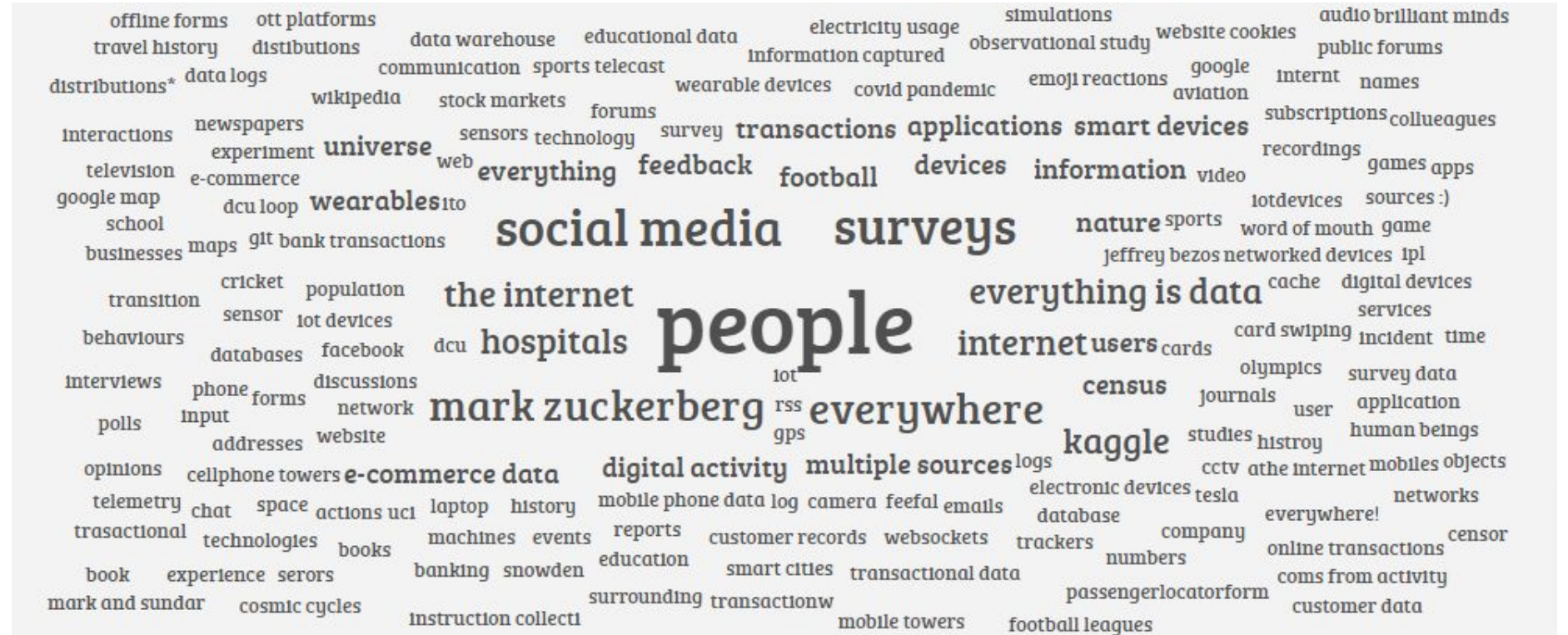
# Why do we care?

Type of data determines:
- What statistics are possible/meaningful
- How data can be processed and/or stored
- Which machine learning model can be used
- Which visualisation method to use

let's try identifying types …
vevox.app

# Data Sources

*Suzanne Little, School of Computing, DCU*

# 2022/2023 Class answers

# 2020/2021 Class answers



Where does data come from?

Type your answer here...    Submit

20 characters remaining

sales data  linkedin
user browsing habits  pollution data  cars  individuals  employees data  observation , survey  hand written doc  any observations
vevox  cutsomer habits  measurements  stocks  phone  video,audio,database  weather  web browsing data  internet browsing  literacy  annual reports  newspapers
spending habits  scriptures  survey  ecommerce sites  ecommerce  erp systems  comms networks  geographical data
train data  my fridge  videos  census  zoom  songs  stock exchange
books  everything  reports  internet  transactions  social media, iot  malware data
satellites  consumer surveys  logfiles  any machine  diaries
exam results  immages  cookies  web  sensors  social media  google nest  humans  iot devices data
from humans  facebook  patents
people create it!!  latin  heath records  payroll
bus ticket  information  ads  library  everywhere  database  encyclopedia  the internet
mobile phone  feedback survey  publications  cctv  us  polls
audio files  smart phones  research papers  phone records  world information
history  human mind  digital media  surveys  news  call records  everthing
citizen data  phone conversations  people  iot  letters  banks  wearable  surveillance video
smartphones  cctv data  camera  cloud personal diaries
smart devices  databases  bio-data  my oven  fitbit  environment  any recorded activit  server data
insurance data  flights  our own trail  recorded information  tracking
browsing data  web, media, images  2nd year students  product information  applications,  social media,sensors  internet,databases
consumer habits  data is everywhere  hospital monitors  multimedia  files, documents, pa  files
photograph  sensors, social medi  humans and machines  smartphones,e-commer
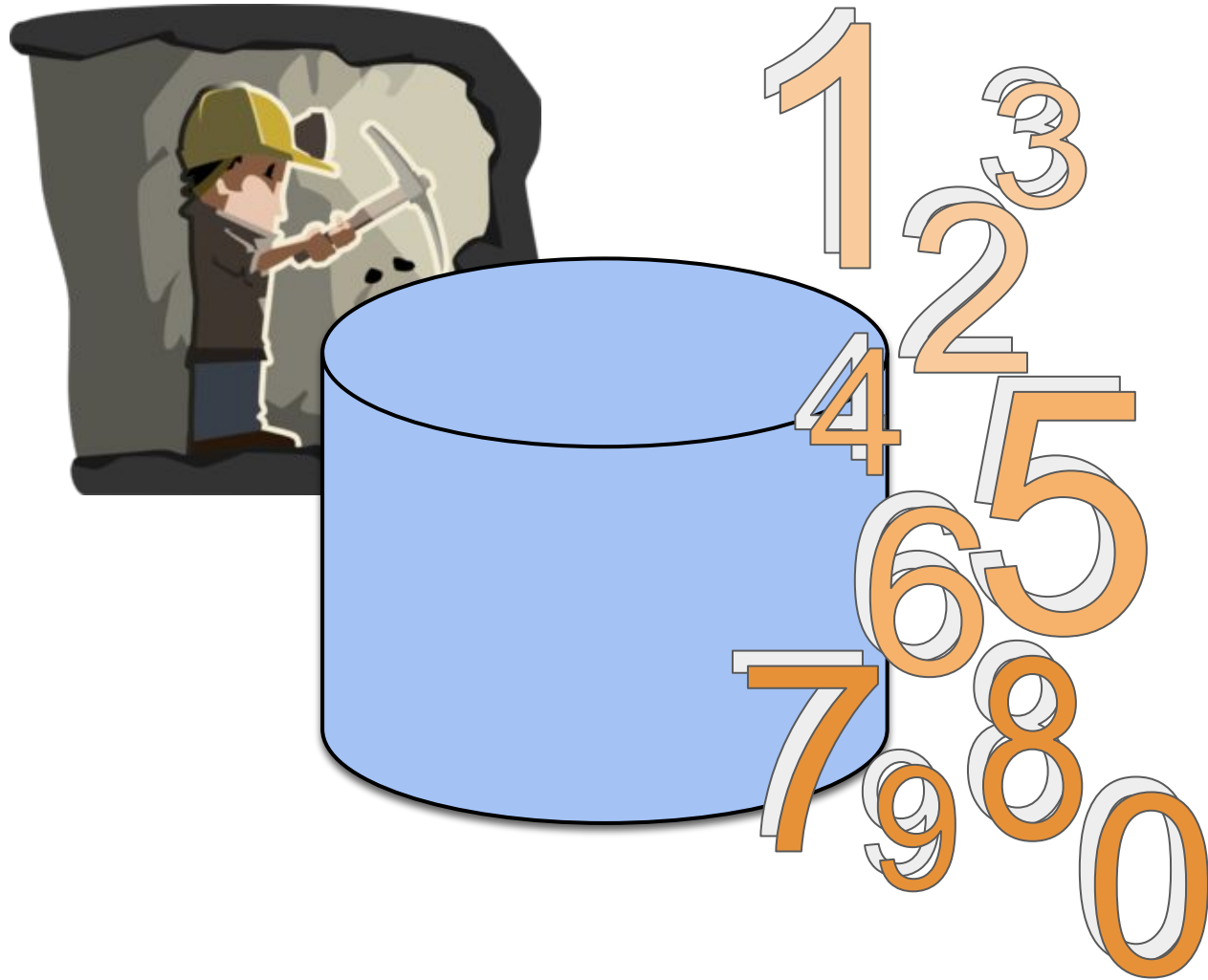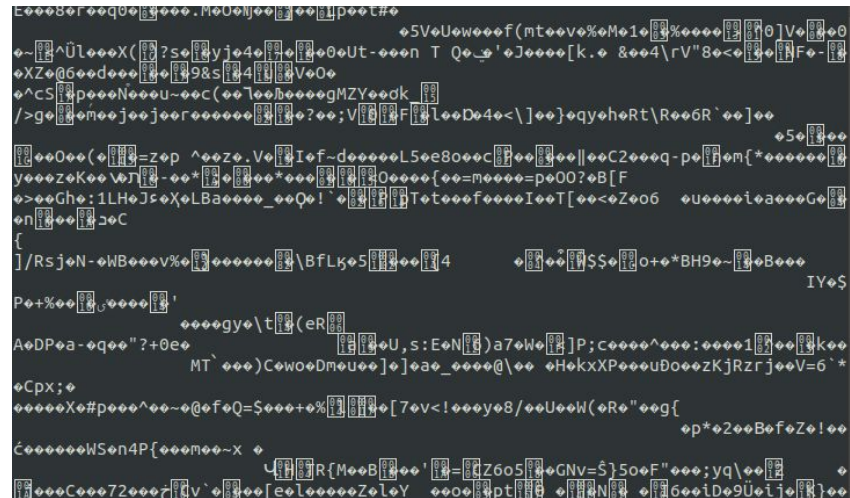
_Suzanne Little, School of Computing, DCU_

# Data sources

- Files
- Databases
- "The Internet"
- Open Data

# Data sources: Files



- Text or Binary
- Open or Proprietary
- Tabulated data - CSV, TSV, DB
- Other text data
  - JSON -- JavaScript Object Notation
  - XHTML -- web pages
  - KML (https://developers.google.com/kml/documentation/kml_tut?csw=1)
  - many other XML-based! (YAML …)
  - Specialist data formats (GDP, ASX etc.) -- may be proprietary
- List of file formats … (https://en.wikipedia.org/wiki/List_of_file_formats)

# Data sources: Databases

- Traditional relational db: Oracle, MySQL, Postgres, etc.
  - Tables ("relations") of rows and columns
  - Unique key per row
  - Links between rows ("foreign key")
  - Optimise structures (the database schema)
  - Stored procedures (queries) to speed up responses
  - Most commonly use SQL - Structured Query Language
    - `SELECT CustomerName,City FROM Customers;`
    - `SELECT CustomerName,Age FROM Customers WHERE City='Dublin';`
- In memory databases: SAP Hana (http://hana.sap.com/abouthana.html)
- NoSQL, document, column, graph, etc.

# Data sources: the Internet

- Crawlers or spiders
    - Scraping data from semi-structured sources
        - Parse HTML
        - Match Patterns to extract data
        - Identify links (repeat)
- URL
    - Files and databases on the web
    - Many libraries and apps will accept either a local path or url
- How many file formats?

# Open data

Public data, shared and freely available

Why?

Why not?

# Examples of open data

Some examples of projects that use open data are:

- [Plantwise - Lose less, feed more](https://www.plantwise.org/) (https://www.plantwise.org/)
- [Humanitarian OpenStreetMap](https://www.hotosm.org/) (https://www.hotosm.org/)
- [OpenGLAM](https://openglam.org/) (https://openglam.org/)

Or on a less elevated topic … the [Great British Public Toilet Map: open geospatial data](#)!

# Where to find open data?

- [https://data.gov.ie/](https://data.gov.ie/)
- [http://www.dublindashboard.ie/](http://www.dublindashboard.ie/)
- [https://www.google.com/publicdata/directory](https://www.google.com/publicdata/directory)
- [https://www.freecodecamp.org/news/https-medium-freecodecamp-org-best-free-open-data-sources-anyone-can-use-a65b514b0f2d/](https://www.freecodecamp.org/news/https-medium-freecodecamp-org-best-free-open-data-sources-anyone-can-use-a65b514b0f2d/)

Also lots of datasets for learning data science:
- [https://www.kaggle.com/datasets](https://www.kaggle.com/datasets)
- [https://github.com/datasets](https://github.com/datasets)

# Exercise

Data.gov is the portal for the US Government's Open Data. Browse the portal and find a dataset to answer the following questions.

1.  What format is the dataset available in?
2.  How many features (attributes or columns) does the data have?
3.  Are the features mostly categorical (qualitative) or numerical (quantitative)?
4.  What is a question that this dataset could help you answer? (you don't need to provide the answer!)

# What's on Loop or will be on github.com?

Material from last week plus Formal Data Management Lifecycles document

Slides & Notes on Data Types

Notes on Files

→ Exercise: Data Formats - Files

Notes on Open Data

→ Exercise: using open data

Linked Data (RDF & SPARQL) → Includes Exercise: using SPARQL on DBPedia