



# Open data driven policy making

## AN2: Data and analysis

Please fill in the register  
[bit.ly/csl-register](http://bit.ly/csl-register)

# Welcome

Commissioned by



Civil Service  
Learning

Delivered by



Dr David Tarrant

The Open Data Institute

Lucy Knight

Local Government



# Aim

To equip policy makers with the knowledge and skills they need to effectively use open data in policy making



# Introductions

Who are you?

What is your role in relation to data driven policy making?

What would you like to get out of this training?



# Today

Introduction to open data driven policy making

Benefits, caveats and risks of using data in policy making

Planning a data driven policy process

Data driven policy making in practice



Civil Service  
Learning

# Data driven policy making



# Introduction to data driven policy making

## Outcomes

1. Define data, big data and open data.
2. Describe how data is used to inform policy making in different fields.
3. Identify the benefits of using data in policy making.
4. Assess the risks, caveats and limitations of using data in policy making.



# Exercise

## What is data?

In as few words as possible, define ‘data’.

You can use an example as an answer if you wish.

One twist: you cannot use the word **information** in your answer.

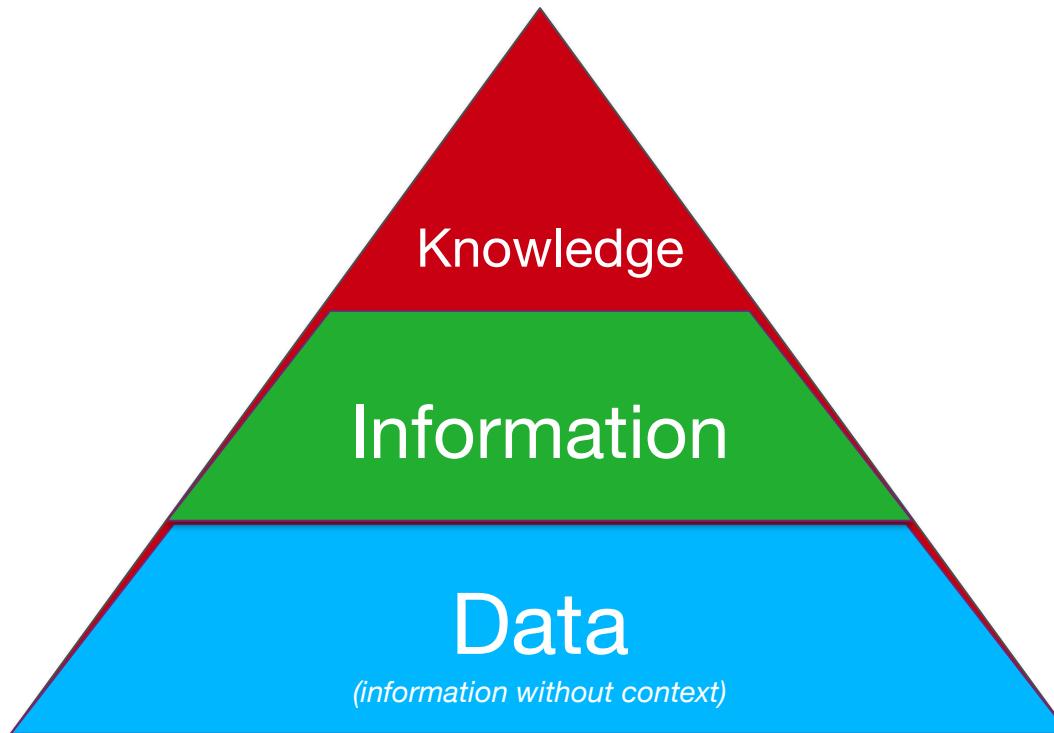


## What is data?

- A collection of facts, information and statistics that can be analysed to develop new knowledge.
- A collection of numbers assigned as values to quantitative variables and / or characters assigned as values to qualitative variables.
- The lowest level of abstraction from which information and then knowledge are derived.



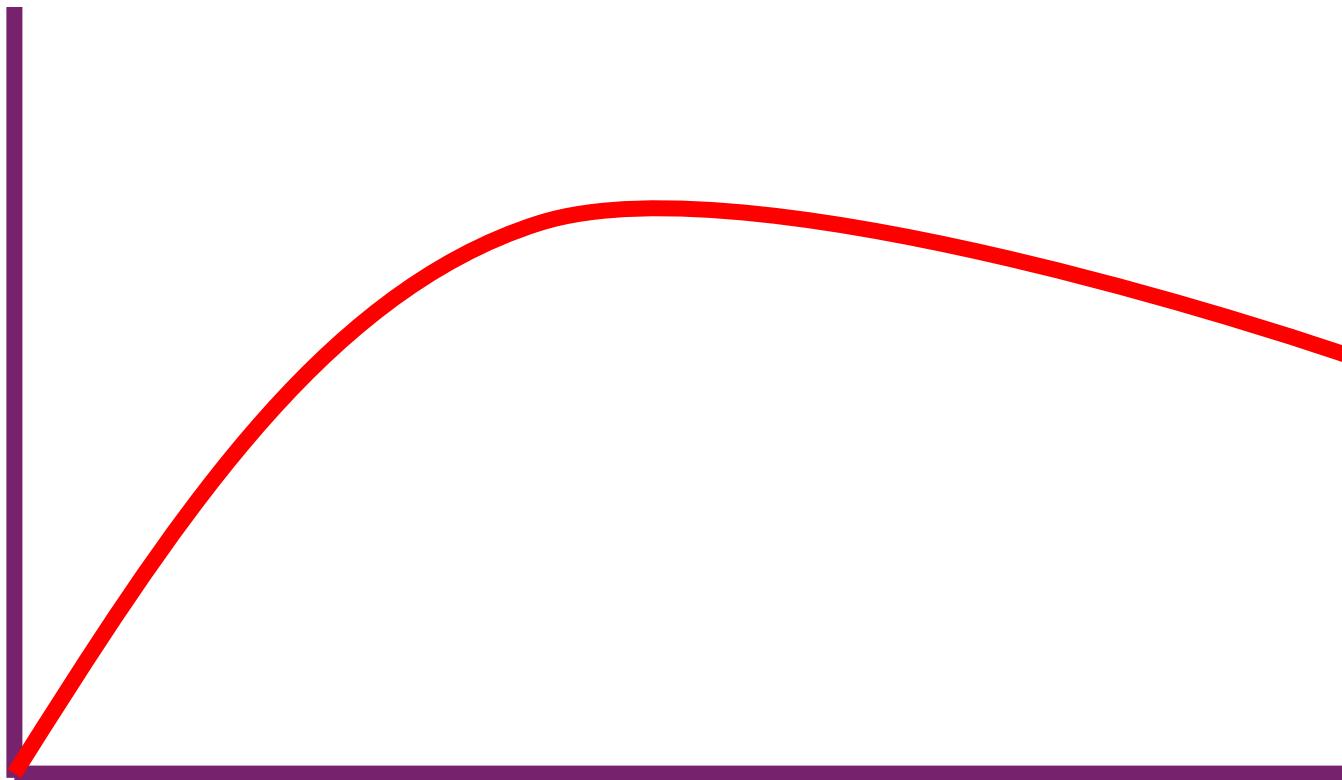
# How is data different from information?







Civil Service  
Learning





## What did you see?

Data: What did you see?

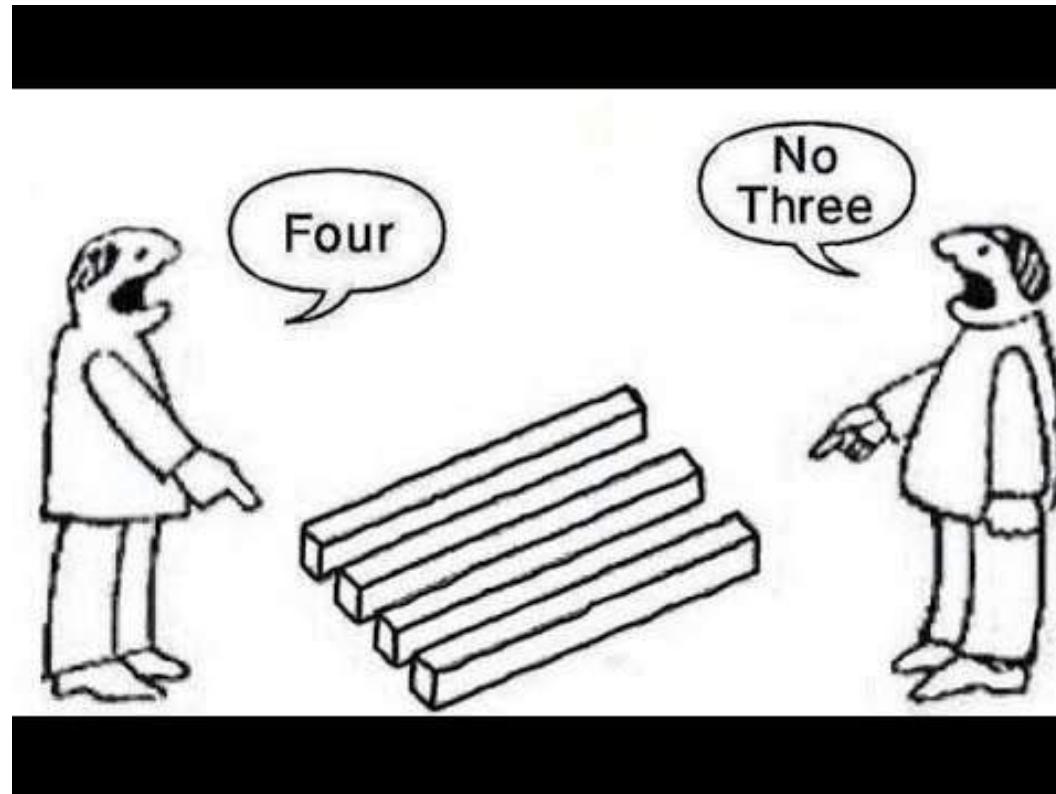
Information: How did you interpret it?

Knowledge: Did it mean (or relate to) anything for you?



## Seeing what isn't there

Even counting is difficult (more on that later)





## Exercise

### What is open data?

In as few words as possible, what is ‘open data’?



A piece of data or content is open if anyone is free to use, reuse, and redistribute it - subject only, at most, to the requirement to attribute and/or share-alike.



Civil Service  
Learning

# Data that anyone can access, use and share.

The Open Data Institute



Open data is data that is published in an open format, is machine readable and is published under a license that allows for free reuse.



# Machine readable data?

A	B	C	D	E	F	G	H
1							
2		Traffic counts on different roads			120000		
3							
4	Road	Cars					
5	M1	100000					
6	M2	50555					
7	M3	25772					
8	Total	176327					
9							
10							
11							
12							
13							

,,,  
,Traffic counts on different roads,,,  
,,,  
,Road,Cars,Lorries,Busses  
,M1,100000,10000,2000  
,M2,50555,20000,1000  
,M3,25772,15478,1500  
,Total,176327,45478,4500  
,,,  
,,,  
,,,  
,,,



## Machine readable data?

	A	B	C	D
1	Road	Cars	Lorries	Busses
2	M1	100000	10000	2000
3	M2	50555	20000	1000
4	M3	25772	15478	1500
5		Road,Cars,Lorries,Busses		
6		M1,100000,10000,2000		
7		M2,50555,20000,1000		
8		M3,25772,15478,1500		



# Open data

The important points are:

- licensed openly (e.g. Open Government Licence)
- free to use, not always free to access (currently government open data must be available at no cost)
- users must be able to modify and redistribute the data

Data.gov.uk includes:

- published in open format
- machine readable



# Case study: LIDAR

 **DATA.GOV.UK**<sup>Beta</sup>  
Opening up Government

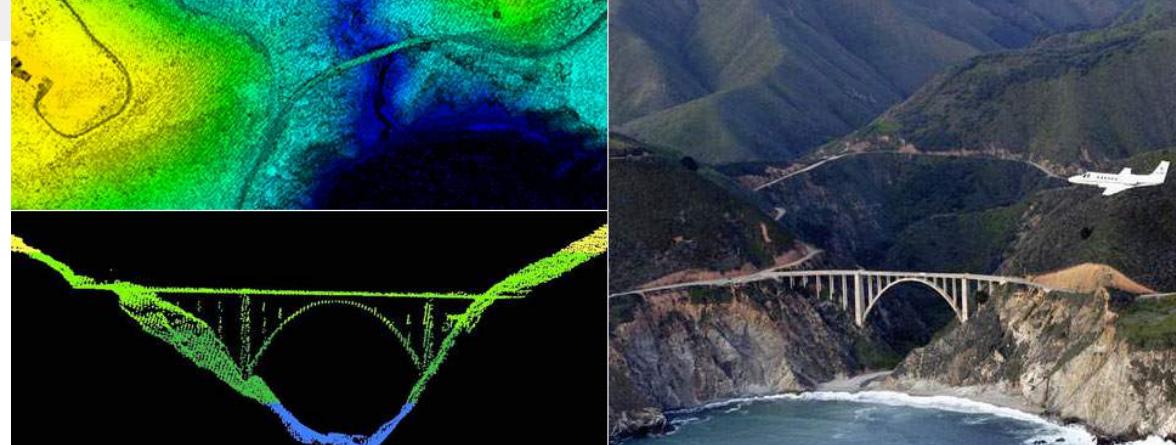
Home

Datasets Map Search Data Requests Publishers Data API Organograms Site Analytics

Home / Datasets / LIDAR Composite DTM - 50cm

## LIDAR Composite DTM - 50cm

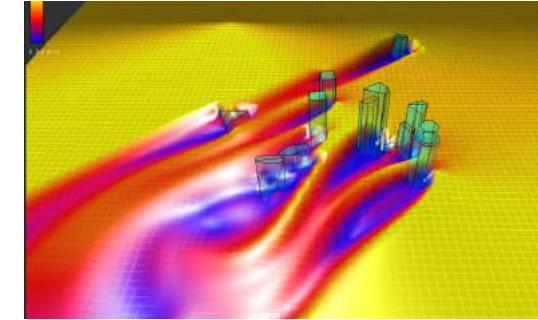
Published by Environment Agency. Licensed under  Open Government Licence.  
Openness rating: ★★★☆☆





## Impact of LIDAR data

Wind modelling



Archology



Educational games (Minecraft)





# Video

**Q: What types and uses of data do Emily and David focus on?**



Civil Service  
Learning



# Review

What types of data do David and Emily focus on?



EMILY focusses on open data and data that is already showing benefit regardless of how big or small it is. Such as LIDAR, Transport, Environmental, Energy and other non-personal data.

DAVID focusses on Big Data and Exhaust Data that come from social interactions and from people. This brings into question lots of ethical and data protection issues and the Big Data Hubris.





# Big data is changing the world

But what is big data?

Collect some ideas together on some post-it notes.  
Do you have any categories emerging?



## Big data is changing the world

Whenever you work with big data you must eliminate the effect of any of these aspects on your result.





# Big data is changing the world

nature

Vol 457 | 19 February 2009 | doi:10.1038/nature07634

## LETTERS

### Detecting influenza epidemics using search engine query data

Jeremy Ginsberg<sup>1</sup>, Matthew H. Mohebbi<sup>1</sup>, Rajan S. Patel<sup>1</sup>, Lynnette Brammer<sup>2</sup>, Mark S. Smolinski<sup>1</sup> & Larry Brilliant<sup>1</sup>

Seasonal influenza epidemics are a major public health concern, causing tens of millions of respiratory illnesses and 250,000 to 500,000 deaths worldwide each year<sup>1</sup>. In addition to seasonal influenza, a new strain of influenza virus against which no previous immunity exists and that demonstrates human-to-human trans-

By aggregating historical logs of online web search queries submitted between 2003 and 2008, we computed a time series of weekly counts for 50 million of the most common search queries in the United States. Separate aggregate weekly counts were kept for every query in each state. No information about the identity of any user was retained. Each



# 97%

Google's claimed accuracy when compared to Centre for Disease Control data.

## google.org Flu Trends

[Google.org home](#)

[Dengue Trends](#)

[Flu Trends](#)

[Home](#)

[United States](#)

[Washington](#)

[Download data](#)

[How does this work?](#)

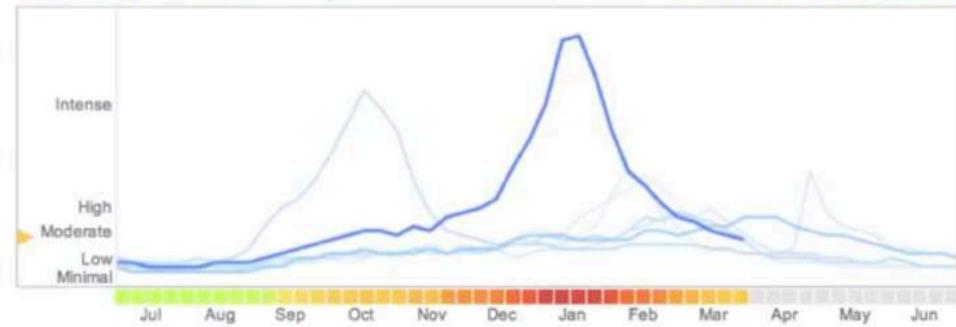
[FAQ](#)

### Explore flu trends - United States

We've found that certain search terms are good indicators of flu activity. Google Flu Trends uses aggregated Google search data to estimate flu activity. [Learn more »](#)

[United States > Washington](#)

● 2012-2013 ● Past years ▾



[States](#) | [Cities](#) (Experimental)





Google

fever aching joints headache

All Shopping Images News Videos More Settings Tools

About 369,000 results (0.52 seconds)

WebMD Symptom Checker helps you find the most common medical conditions indicated by the symptoms **chills, fever, headache** and **joint aches** including Lyme disease, Acute sinusitis, and Aseptic meningitis (adult). ... Osteoarthritis happens when the cartilage in your **joints** breaks down causing pain, stiffness, and swelling.

**Chills, Fever, Headache and Joint aches: Common Related Medical ...**  
[symptomchecker.webmd.com/multiple-symptoms?...chills%7Cfever%7Cheadache%7Cjoints%7C...](http://symptomchecker.webmd.com/multiple-symptoms?...chills%7Cfever%7Cheadache%7Cjoints%7C)

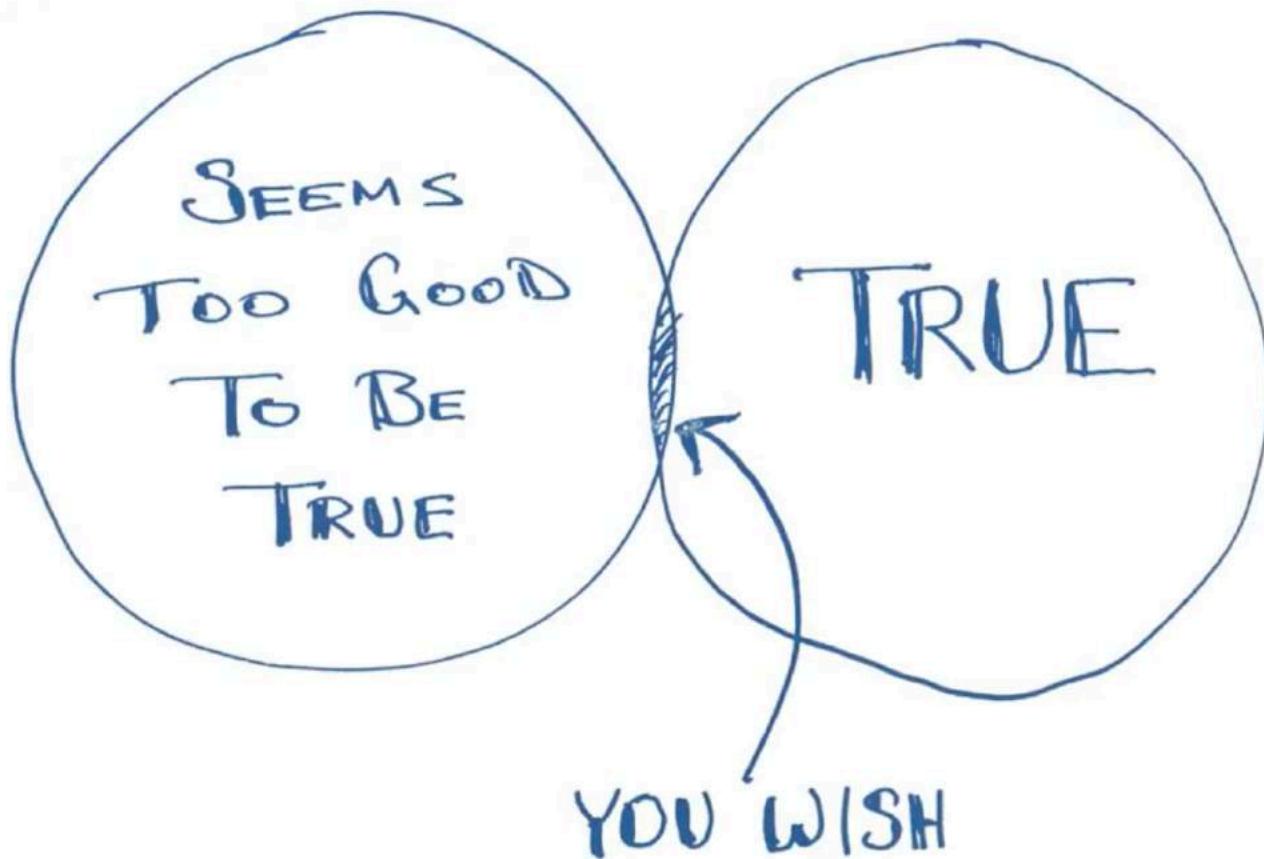
About this result • Feedback



Uncertainty of data

People making flu-related Google searches may know very little about how to diagnose flu?

Does a search for flu mean they have flu or just an interest in it?





I USED TO THINK  
CORRELATION IMPLIED  
CAUSATION.



THEN I TOOK A  
STATISTICS CLASS.  
NOW I DON'T.



SOUNDS LIKE THE  
CLASS HELPED.  
WELL, MAYBE.



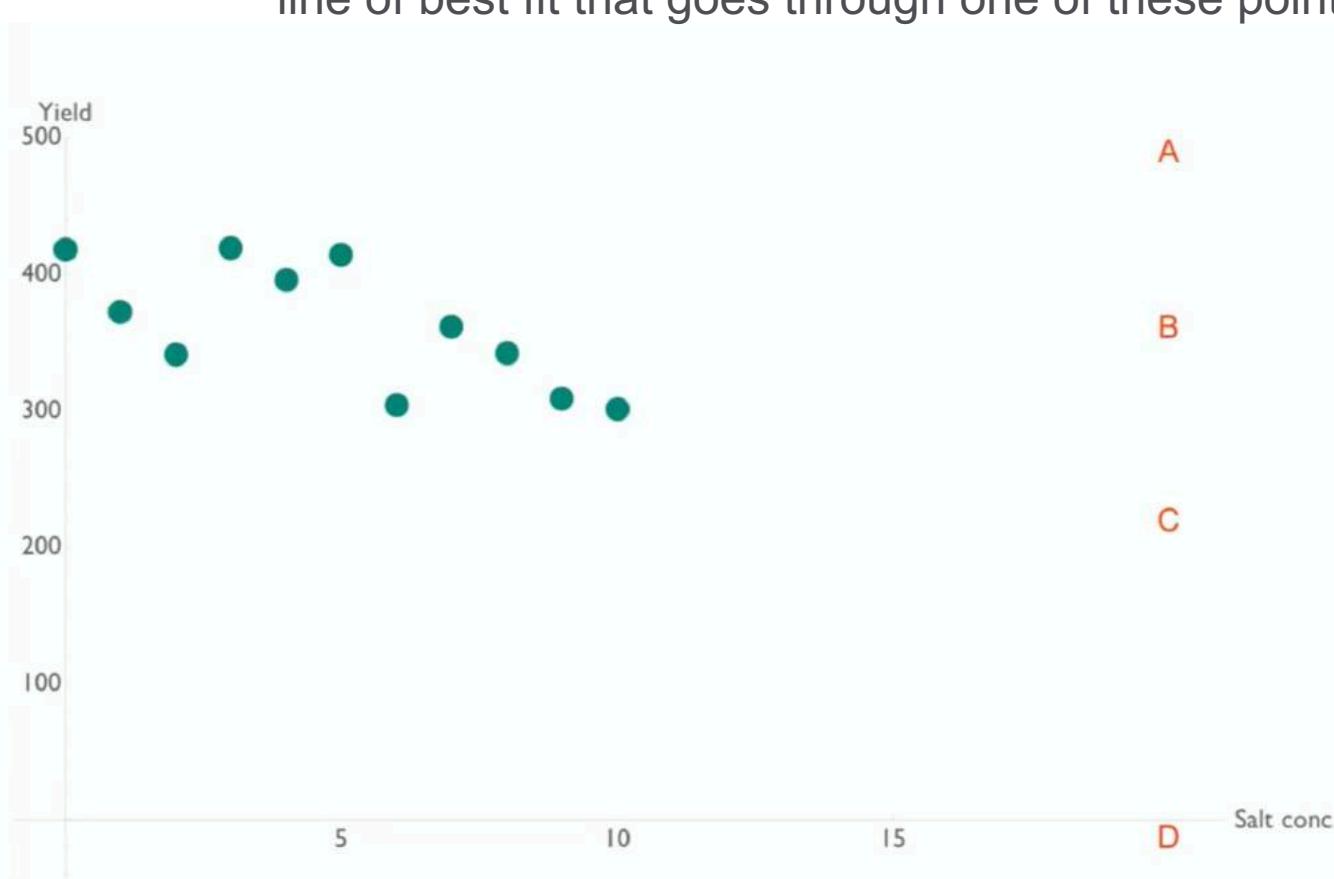
<https://xkcd.com/552/>

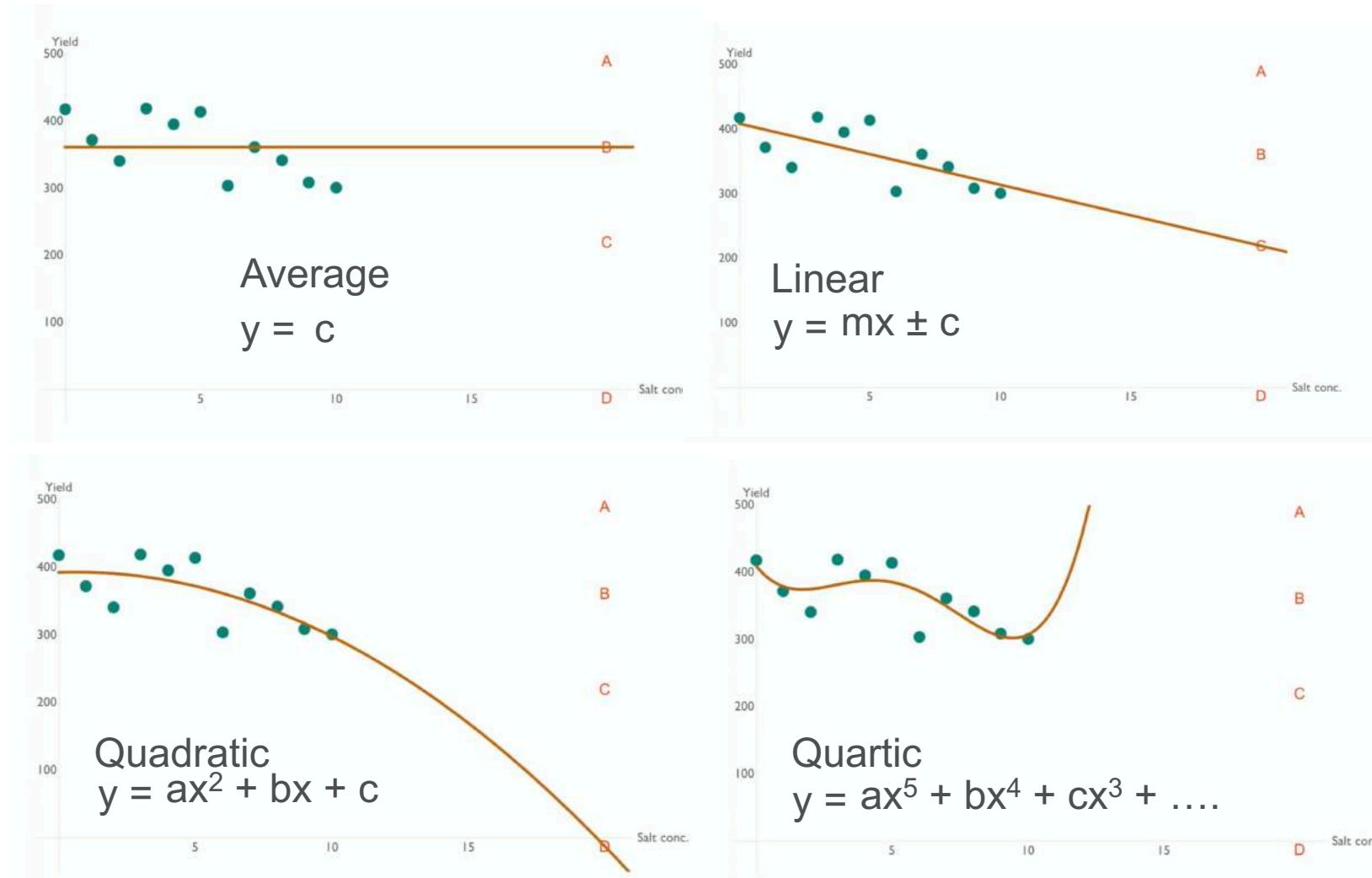


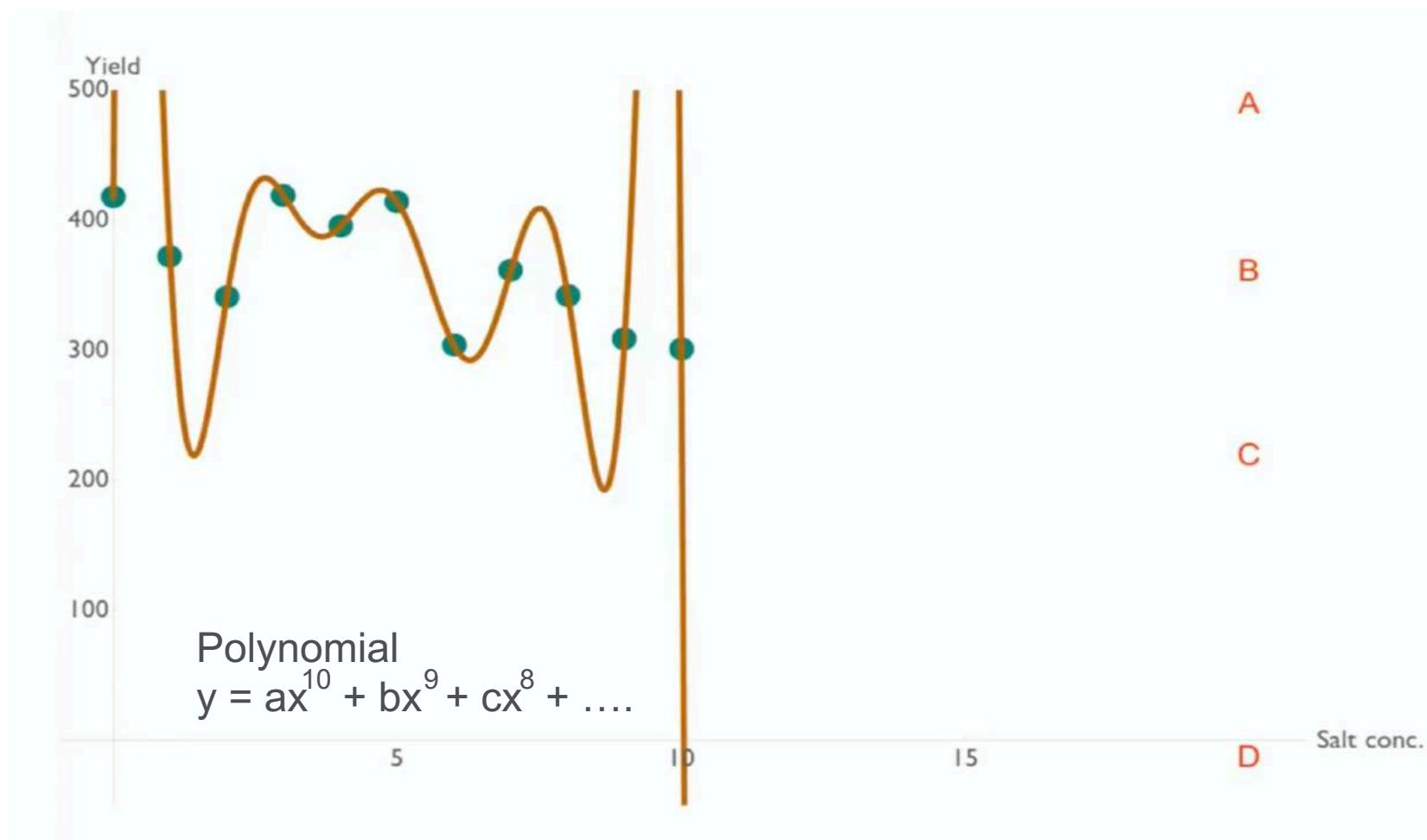
# Overfitting

Plotted here are 10 points in a series that have been generated using a function.

If another 10 points were plotted of the same function which point (A,B,C or D) would they tend towards. Add a line of best fit that goes through one of these points.









## Why did the flu trends example fail?

- Big Data Hubris (complement not supplement)
- Huge veracity problem!
- Used nth degree polynomial (overfitting)
- Solely relied on this data over scientific method

*For more on big data search YouTube for “calling bullshit on big data”*

*Excellent lecture series from the University of Washington*





# Data analysis 101

You each have a set of cards containing numerical values representing something. Can you work out what the dataset represents?

## Data analysis stages:

Data: What do you see?

Information: What can you calculate?

Knowledge: What does it mean to you?



# Averages

- What problems does this dataset have when working out the average?
- How might we solve these problems?
- Is the middle value (with data points ordered) better than the mean?

# Outliers

33,750.00	33,750.00
44,000.00	33,750.00
138,188.00	33,750.00
45,566.67	33,750.00
44,000.00	44,000.00
141,666.67	44,000.00
292,500.00	44,000.00
5,600,000.00	44,000.04
103,500.00	45,566.67
190,000.00	65,000.00
65,000.00	95,000.00
33,750.00	103,500.00
195,000.00	112,495.50
44,000.04	138,188.00
4,600,000.00	141,666.67
194,375.00	181,500.00
33,750.00	185,000.00
112,495.50	190,000.00
95,000.00	194,375.00
301,999.00	195,000.00
181,500.00	205,000.00
33,750.00	292,500.00
185,000.00	301,999.00
205,000.00	4,600,000.00
44,000.00	5,600,000.00

Average (mean): 518,311.64

Median: ~125k

16% of 25 = 4

16% trimmed mean: 128,109.09  
=trimmean(RANGE,0.16)



# The distribution of data

Divide the data into 6 evenly distributed buckets of values and stack the cards as per the diagram here.

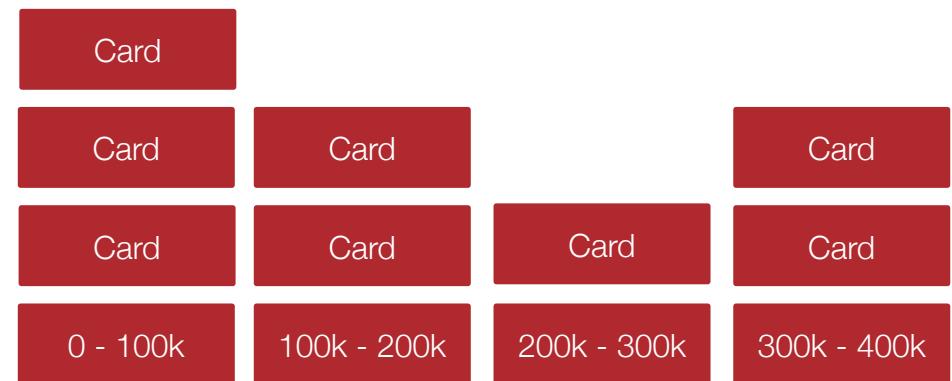
Suggested buckets:

0 – 100,000

100,000 – 200,000

200,000 - 300,000

...

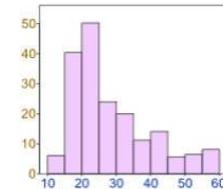




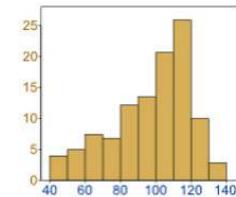
# The distribution of data

Another important aspect to consider is the distribution of data. It is always a good practice to know the distribution of your data before analysing it further. Certain analyses require certain distributions.

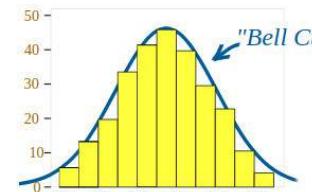
The examples here show different types of distributions of data.



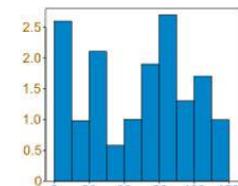
Positively skewed distribution



Negatively skewed distribution



Normal distribution



Non-normal distribution



## 5 number summary (of the trimmed dataset)

Plot 3 more lines

- 1) Median
- 2) Value  $\frac{1}{4}$  through dataset
- 3) Value  $\frac{3}{4}$  through dataset

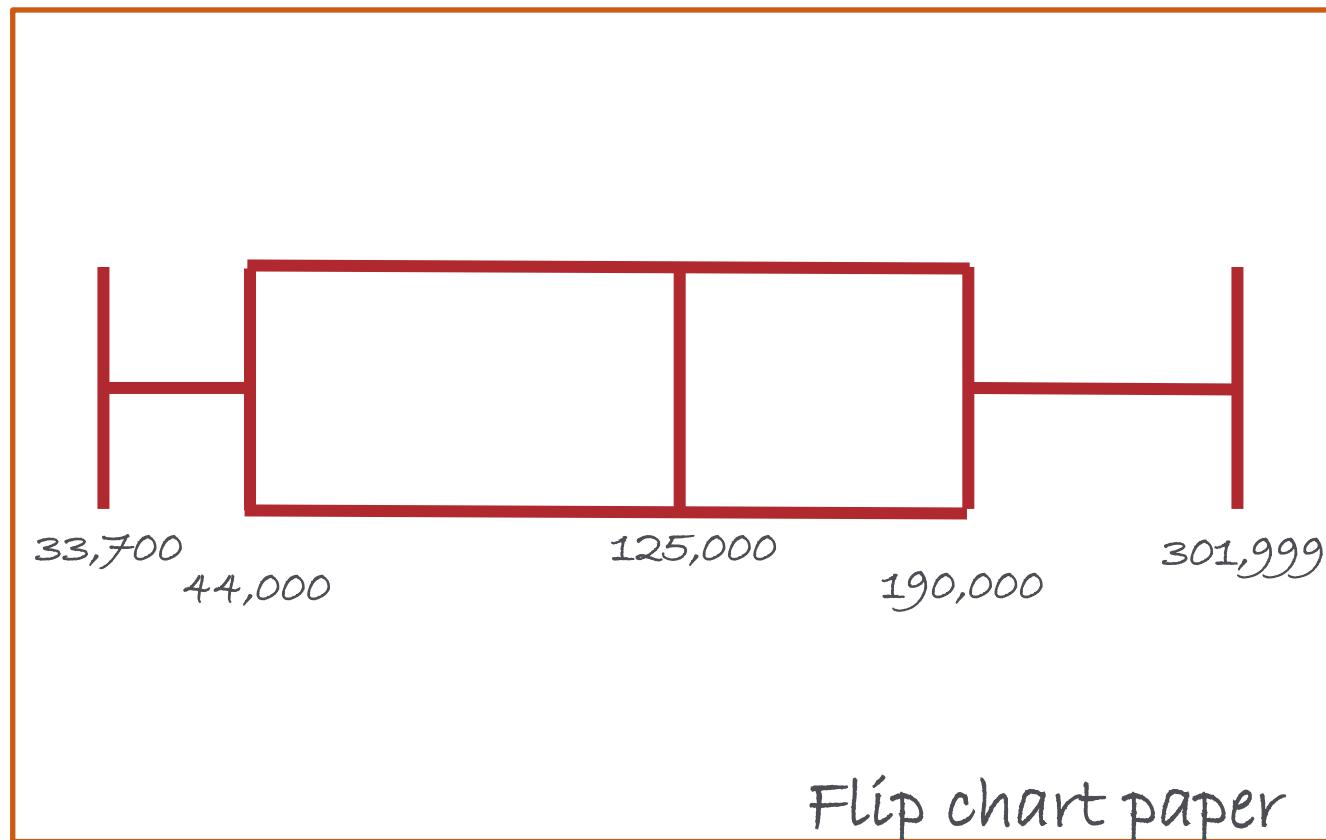
33,700

301,999

Flip chart paper

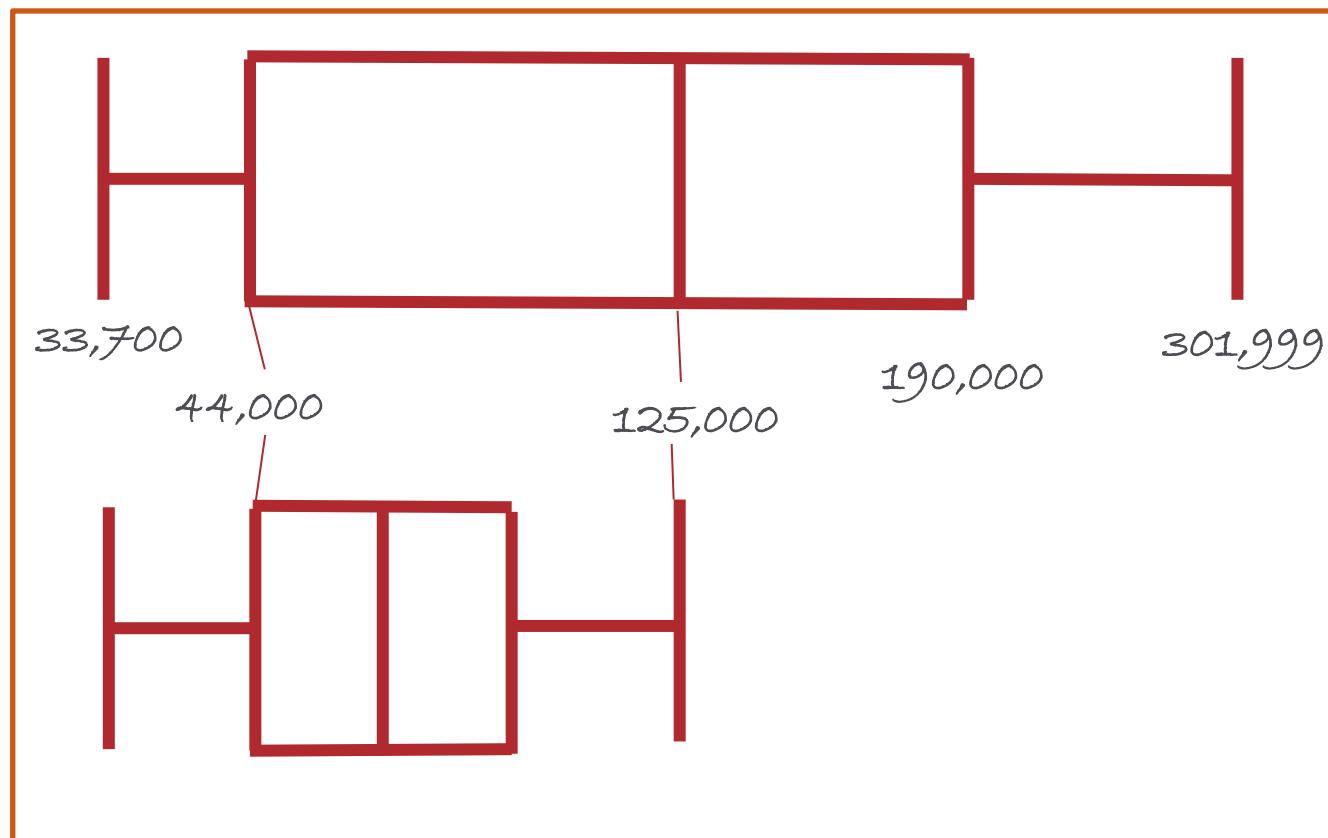


## 5 number summary (of the trimmed dataset)





## 5 number summary (of the trimmed dataset)

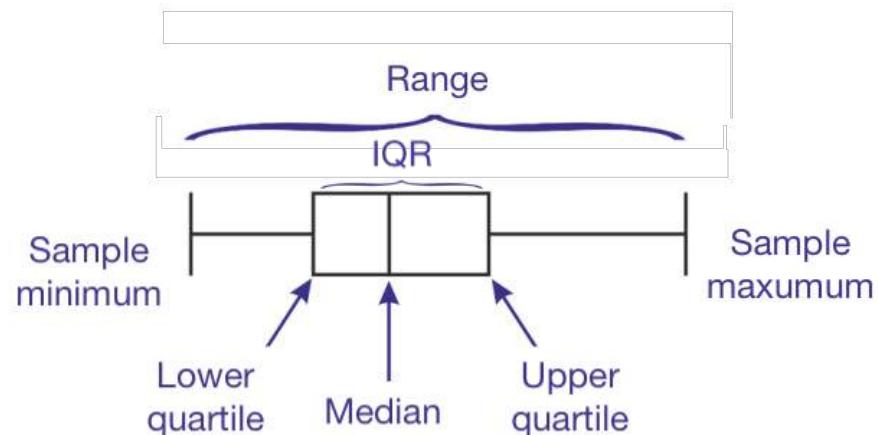




# Five-number summary

The five-number summary

1. the sample minimum (smallest value)
2. the lower quartile (value  $\frac{1}{4}$  through the list)
3. the median (middle value in the list)
4. the upper quartile (value  $\frac{3}{4}$  through the list)
5. the sample maximum (largest value)





# Five-number summary in Tableau

The screenshot shows the Tableau software interface with three main panels:

- Connect:** On the left, a dark sidebar with options to "Connect To a File" (Microsoft Excel, Text file, JSON file, PDF file, Spatial file, Statistical file) and "To a Server" (OData, More...). It also includes a note about working with big data and a "Upgrade Now" button.
- Open:** The central panel, titled "Open", displays four items:
  - A blank white square labeled "LFB\_All\_Data".
  - A blank white square labeled "LFB\_Trial1".
  - A map of London boroughs labeled "London\_boroughs".
  - A scatter plot labeled "Greater Manche...".
- Discover:** On the right, a sidebar with links to "How-to Videos", "Overview", "Intro to the Interface", "Chart Types", and "More how-to videos...". It also features a "Viz of the Day" section titled "THE POTTERVERSE - FAMILY TREE" with links to "Blog - Step and jump into Tableau Public 2018.1", "Sample Data Sets", "Live Training", and "Current Status".



# Summary

Data, it turns out, has shape.  
That shape has meaning.

The shape of data tells you everything you need to know about your data from its obvious features to its deepest secrets.





we can look at what words written down by a social worker or a police officer,



## Machine learning and prediction

Each table has a set of “Top Trump” cards relating to properties in two cities.

Build a decision tree to sort them into “New York” and “San Francisco”.

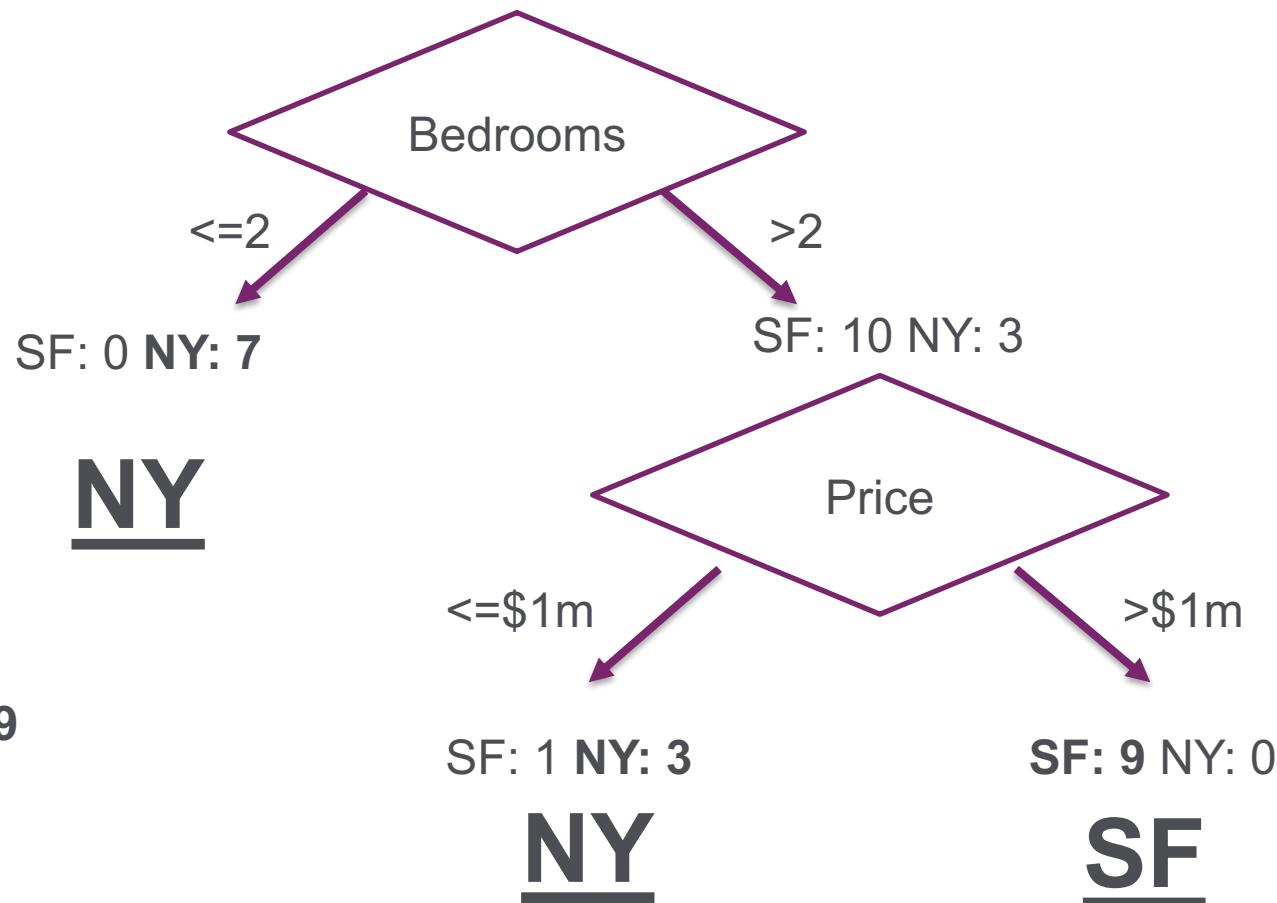
You cannot use the name of the city to sort them.



95%

Confidence

## Example decision tree





# Approaches



Data first



Knowledge first



# Assumption

What share of income tax paid in the UK is paid by the top 1% of earners?

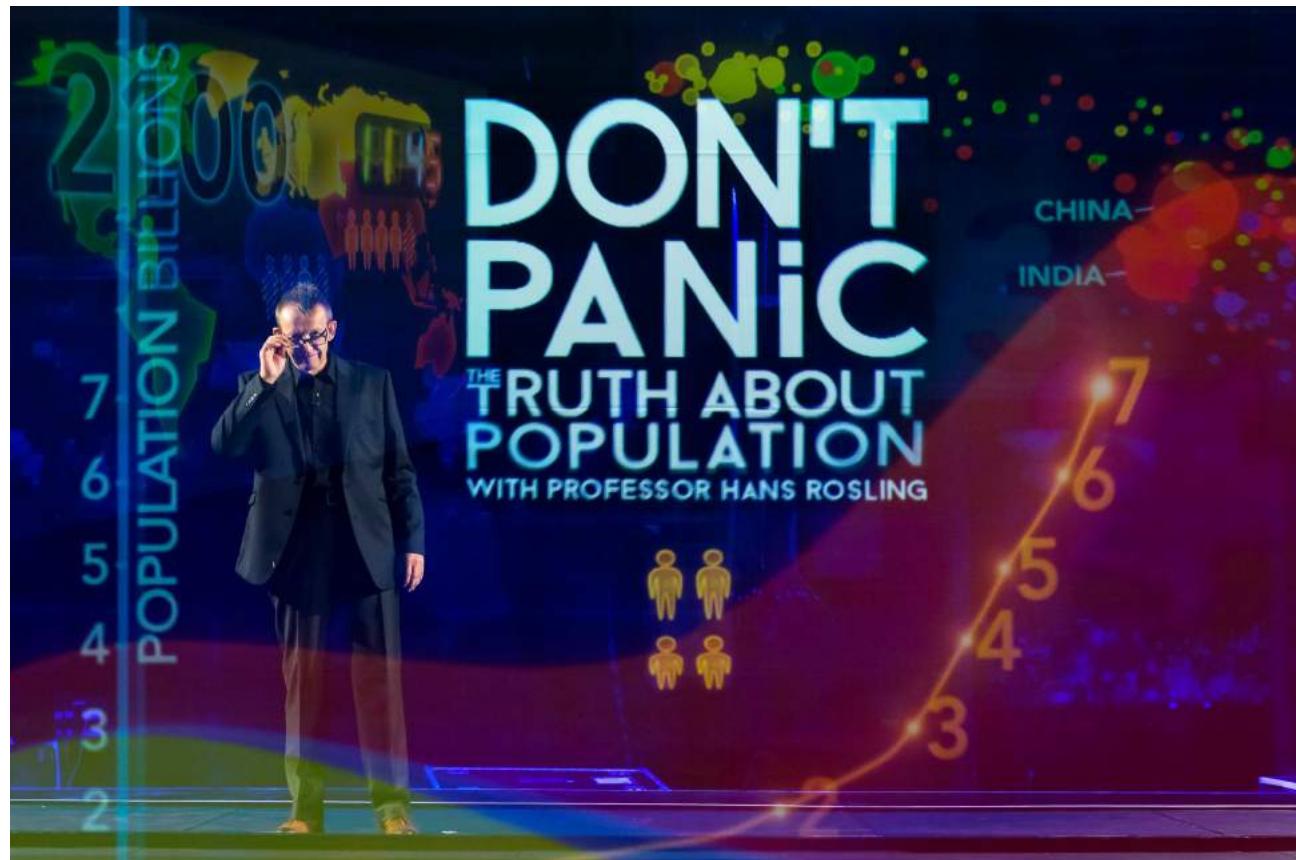
- ◆ A: 5%
- ◆ C: 14%
- ◆ B: 9%
- ◆ D: 17%



# Assumption

What is the average number of children per family in Bangladesh?

- ◆ A: 2
- ◆ B: 3
- ◆ C: 4
- ◆ D: 5



[gapminder.org](http://gapminder.org)



# Data analysis of the houses dataset in Tableau

The screenshot shows the Tableau Public interface with the following details:

- Top Bar:** Standard mode, Show Me button.
- Left Panel (Data Source):** Shows the "houses" dataset with the following fields:
  - Dimensions:** Target, Year Built (selected), Zip.
  - Measures:** Bath, Beds, Elevation, Index, Price, Price Per Sqft, Sqft, Latitude (generated), Longitude (generated), Number of Records, Measure Values.
- Middle Panel:** Columns and Rows sections under Pages, and a Filters section.
- Right Panel (Sheet 1):** A blank canvas with three "Drop field here" placeholder areas. The Marks section is set to Automatic, with options for Color, Size, Text, Detail, and Tooltip.
- Bottom Bar:** Data Source tab (selected), Sheet 1 tab, and other navigation icons.



# Options?

Grow the decision tree until...

1. every classification is perfect?
2. confidence level is above 80%?
3. it is too big to process the data in reasonable time?
4. Until the evaluation and training set have the same confidence?
5. You have used all principal components?



Civil Service  
Learning

# A visual introduction to machine learning

[r2d3.us](http://r2d3.us)

---





# Considerations

Combine the data with existing knowledge but make sure the existing knowledge is correct!

Grow the tree around the principal components (using PCA). Stop when only minor improvements are achieved on the training set. Don't overfit!



Civil Service  
Learning

GoCompare

ALL Details EXACTLY the same

Hi John!

Adjust cover

49 car quotes found, 4 telematics quotes included.

Sort: Annually

**£538.26**

Total Excess: £150

[View details >](#)



Legal  
Assistance  
+£30.99



Breakdown  
Cover  
+£43.99



Personal  
Accident  
✓



Windscreen  
✓



Courtesy Car  
✓

M&S BANK



Legal  
Assistance  
+£26.29



Breakdown  
Cover  
+£31.54



Personal  
Accident  
✓



Windscreen  
✓



Courtesy Car  
✓

**£571.66**

Total Excess: £150

[View details >](#)



Legal  
Assistance  
+£30.99



Breakdown  
Cover  
+£43.99



Personal  
Accident  
✓



Windscreen  
✓



Courtesy Car  
✓

**£573.39**

Total Excess: £150

[View details >](#)



Legal  
Assistance  
+£30.99



Breakdown  
Cover  
+£43.99



Personal  
Accident  
✓



Windscreen  
✓



Courtesy Car  
✓

**£577.87**

Total Excess: £150

[View details >](#)



Civil Service  
Learning

GoCompare

Car: Ford Fiesta Ghia 2002-2008 1.6 Petrol  
Profession: Insurance Director  
Address: Milford Haven (PPI Company :P)

Adjust cover

37 car quotes found, 5 telematics quotes included.

Always resided in UK  
Date of birth: 01/01/1980  
No claims and license: 16 years  
Car kept on drive

Hi Mohammed!

Sort: Annually

 <b>LLOYDS BANK</b>	 Legal Assistance +£30.99	 Breakdown Cover +£43.99	 Personal Accident ✓	 Windscreen ✓	 Courtesy Car ✓	<b>£1,446.32</b> Total Excess: £150	<a href="#">View details &gt;</a>
 <b>HALIFAX</b>	 Legal Assistance +£30.99	 Breakdown Cover +£43.99	 Personal Accident ✓	 Windscreen ✓	 Courtesy Car ✓	<b>£1,458.04</b> Total Excess: £150	<a href="#">View details &gt;</a>
 <b>BANK OF SCOTLAND</b> Decisions well made	 Legal Assistance +£30.99	 Breakdown Cover +£43.99	 Personal Accident ✓	 Windscreen ✓	 Courtesy Car ✓	<b>£1,459.64</b> Total Excess: £150	<a href="#">View details &gt;</a>
<b>M&amp;S BANK</b>	 Legal Assistance ✓	 Breakdown Cover ✓	 Personal Accident ✓	 Windscreen ✓	 Courtesy Car ✓	<b>£1,476.70</b> Total Excess: £150	<a href="#">View details &gt;</a>



## Good or bad ideas?

1. a tool that analyses the sentiment of a user's tweets, assesses whether they are suicidal and alerts friends	4. a risk-assessment tool that uses AI to advise on prison sentences based upon criminal profile analysis
2. automatic pricing algorithm for taxi firm which responds to surges in demand	5. using energy efficiency data and winter fuel allowance data to target efficiency advice
3. using performance data to advise on how to save money in the emergency services	6. publishing genomes of 100,000 individuals for use in public health



# Automated Inference on Criminality using Face Images

Xiaolin Wu

McMaster University  
Shanghai Jiao Tong University

xwu510@gmail.com

Xi Zhang

Shanghai Jiao Tong University  
zhangxi\_19930818@sjtu.edu.cn



“Unlike a human examiner/judge, a computer vision algorithm or classifier has absolutely no subjective baggages [sic], having no emotions, no biases whatsoever due to past experience, race, religion, political doctrine, gender, age, etc., no mental fatigue, no preconditioning of a bad sleep or meal. The automated inference on criminality eliminates the variable of meta-accuracy (the competence of the human judge/examiner) all together.”

—Wu & Zhang(2017)



“Unlike a human examiner/judge, computer vision algorithm or classifier has absolutely no subjective baggages [sic]—no emotions, no biases whatsoever due to experience, race, religion, political affiliation, gender, age, etc., no mental fatigue, no conditioning of a bad sleep or meal. This is the variable of meta-accuracy (the combination of the human judge/examiner) all together.”

**Utter Bull...  
—Wu & Zhang(2017)**

*Example from “calling bullshit” lecture series from the University of Washington*



## Good or bad ideas?

1. a tool that analyses the sentiment of a user's tweets, assesses whether they are suicidal and alerts friends	4. a risk-assessment tool that uses AI to advise on prison sentences based upon criminal profile analysis
2. automatic pricing algorithm for taxi firm which responds to surges in demand	5. using energy efficiency data and winter fuel allowance data to target efficiency advice
3. using performance data to advise on how to save money in the emergency services	6. publishing genomes of 100,000 individuals for use in public health



# Data ethics

More on the examples used in this section (as well as more examples) -

Guidance

## Data Science Ethical Framework

This framework is intended to give civil servants guidance on conducting data science projects, and the confidence to innovate with data.

---

Published 19 May 2016

Last updated 13 June 2018 — [see all updates](#)

From: [Cabinet Office](#), [Government Digital Service](#), and [The Rt Hon Matt Hancock MP](#)

**This publication was withdrawn on 13 June 2018**

This guidance has been replaced by the [Data Ethics Framework](#).

<https://www.gov.uk/government/publications/data-science-ethical-framework>



# Data ethics canvas



<https://theodi.org/article/data-ethics-canvas/>



# Course materials

<http://training.theodi.org/csl-an2/>

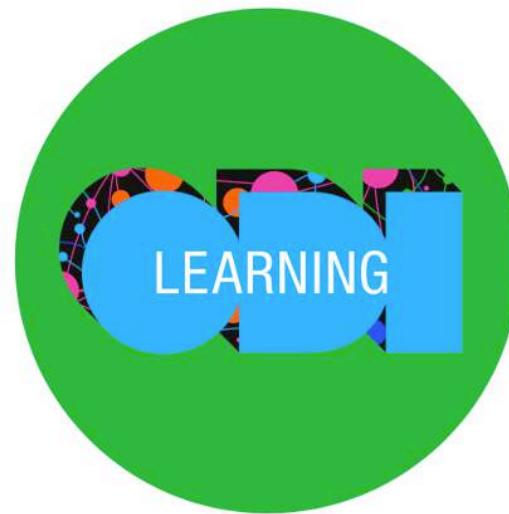
## AN2 - Data and analysis

Welcome. This page accompanies the **AN2 - Data and analysis** course and provides access to slides, exercises and related case studies.

The main aim of the course is to *equip policy makers with the knowledge and skills they need to effectively use open data in policy making.*

Key outcomes of the course include:

- Define data, open data and big data
- Describe how data is used to inform policy making
- Identify the benefits of open data to policy making
- Perform some basic statistical analysis on data to identify the shape and trends in data.
- Assess the ethical risks of using data in policy making
- Perform a practical piece of policy making with real data





# Data analysis in practice

## Outcomes

6. List the stages in carrying out data analysis for policy making.
7. Create a plan for carrying out data analysis for policy making.
8. Review the role of open data in policy making

## Homework

9. Carry out a simple data analysis using a number of tools.
10. Create an interactive data visualisation.
11. Communicate the results of a data analysis to decision makers.



# What has been the impact of closing 10 fire stations in London?



## Data analysis practical

You have been tasked with saving money in the local fire service.

You must propose a solution that saves money while minimising the impact on service delivery.

Which factor do you think is the most significant in analysing performance to identify savings (top of the decision tree)

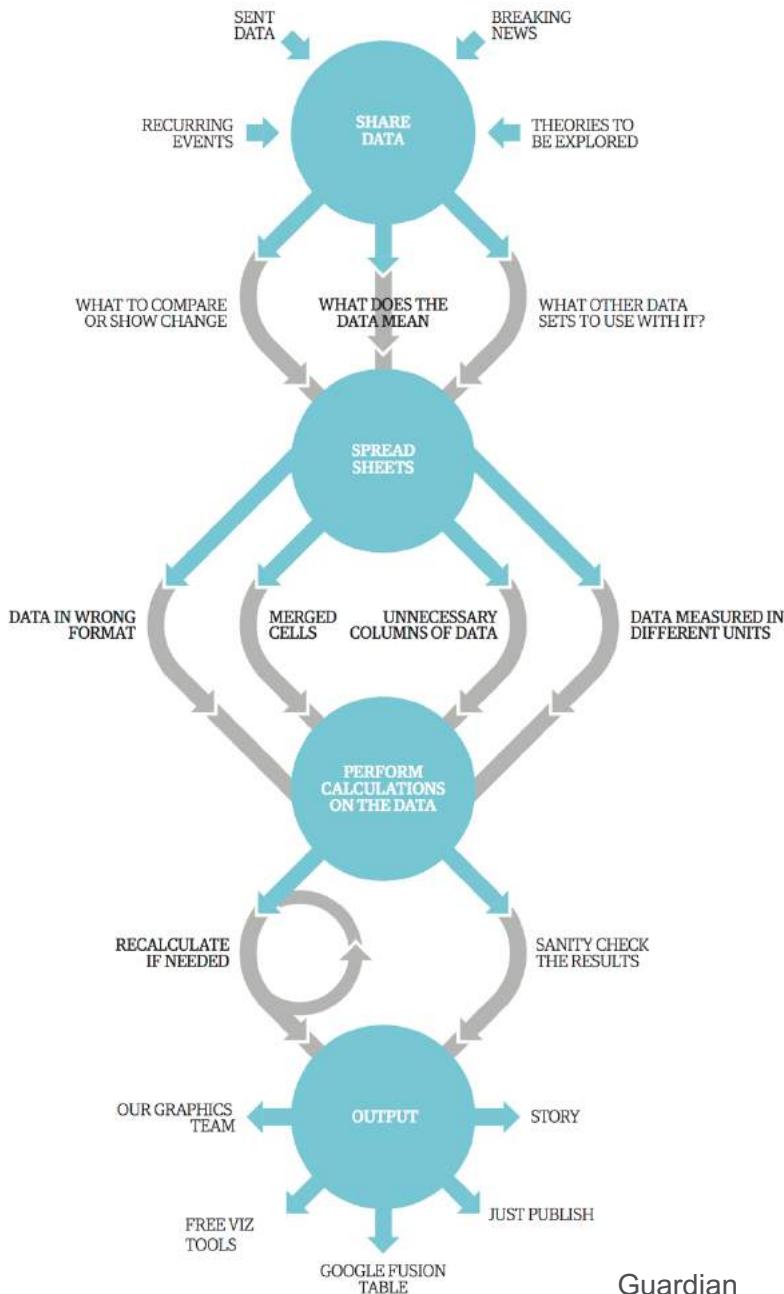


Paul Hudson



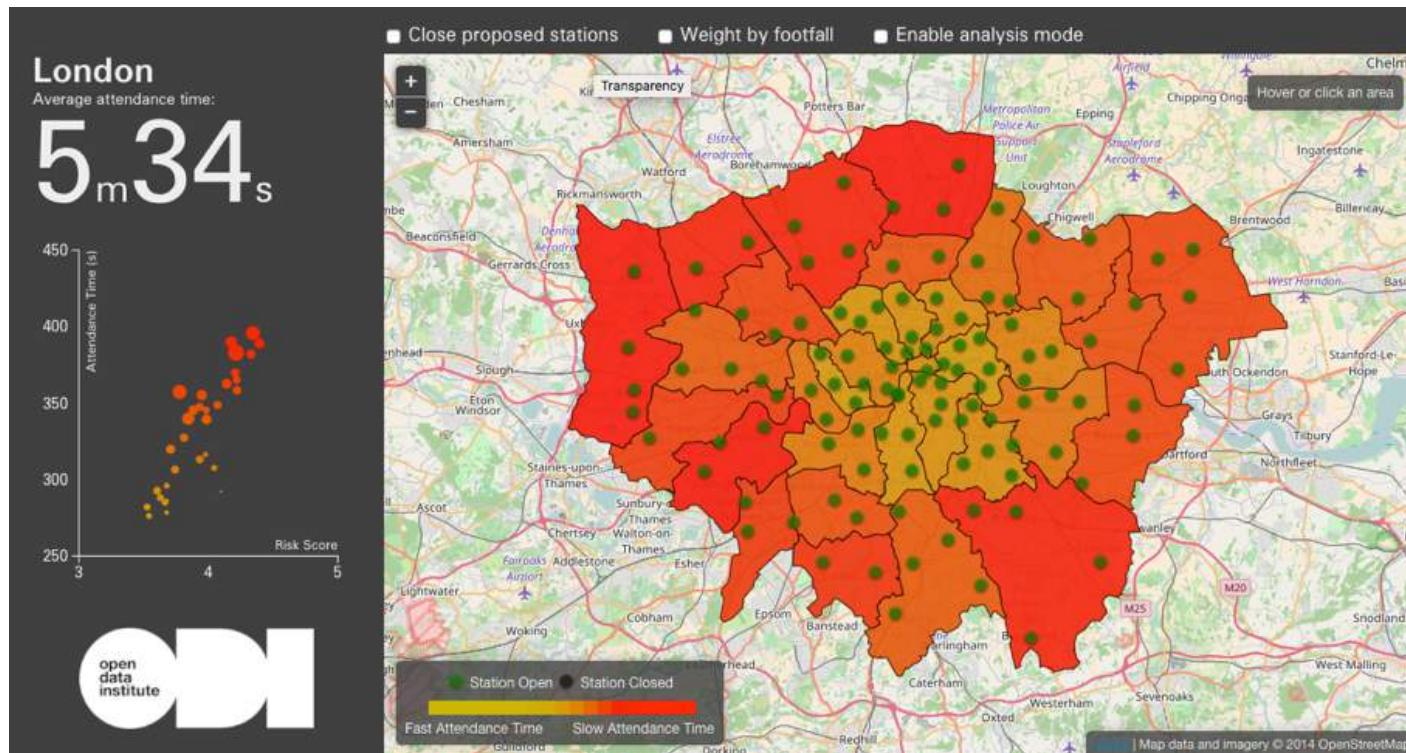
# Using data in analysis

1. Collect data.
2. Establish context of data.
3. Clean the data.
4. Translate, filter and merge data.
5. Initial evaluation.
6. Go back to start?
7. Secondary evaluation and analysis.
8. Sanity check.
9. Create output.





<http://london-fire.labs.theodi.org/explore/>





## Stage 1: collect data

Make a space on your group desk

Put a yellow post-it note in the middle and draw a fire on it. Write a time on it between 4 minutes and 7 minutes.

Scatter a series of all different colour (including yellow) post-it around this one to keep it in the middle. Write a time between 4 and 7 minutes on each one.

## Stage 2: establish context

Each post-it note represents an incident that has been responded to by the London fire brigade.

Each colour represents a different fire stations from which the appliance (fire engine) was sent



## Stage 3: clean the data

Given the context does your data make sense?

## Stage 4: Translate, filter and merge data

Nothing to do...



## Stage 5: initial evaluation

How do we calculate the impact of closing the YELLOW fire station (and specifically the middle incident)?

What is your design?



## Stage 8: Sanity check

Does your design overfit (or use nth degree polynomial?)

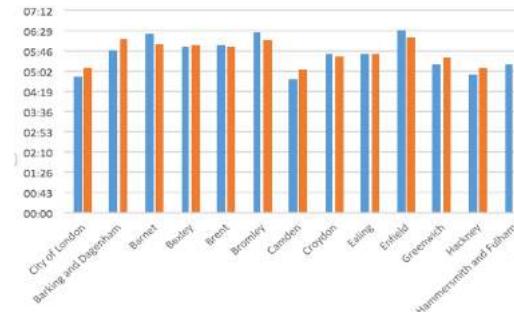
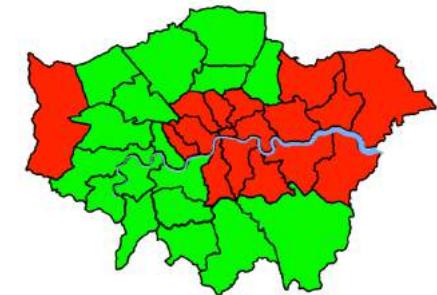
Can your design really model the future or just show how the past was?



## Stage 9: Creating output

There are many options for output. Consider the best one for your purpose.

- A map can be ideal for such geographic datasets (search for London Data Store borough Excel KML).
- A chart showing a comparison of before an after statistics can be easily generated in excel.
- An highly interactive dashboard can be created with [dataseedapp.com](http://dataseedapp.com).





# Review

London Fire Stations analysis was a story started by a participant in an ODI Training course like this one. What is your story? Who and what would you need to help you?

If I'm a policy maker, what do I now need to go and do?

What do I need to go and do next?

What is your main takeaway from the course overall?

What are you going to apply back in the workplace?

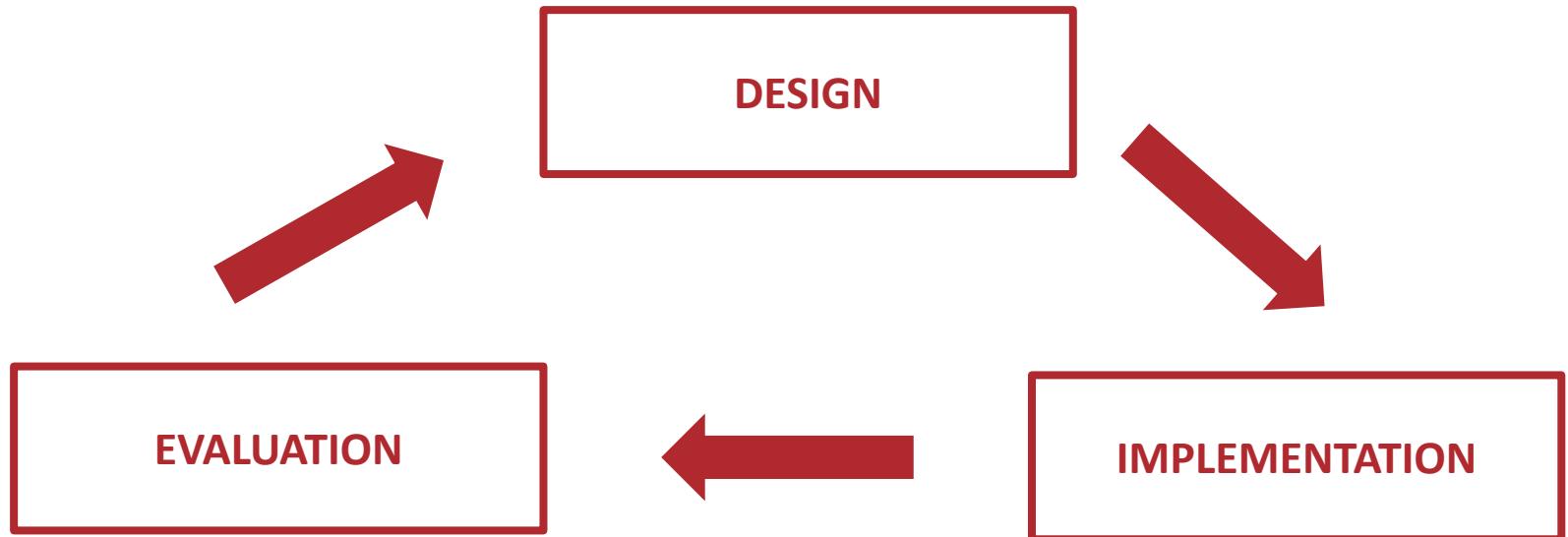


Civil Service  
Learning

# Open data in policy cycles



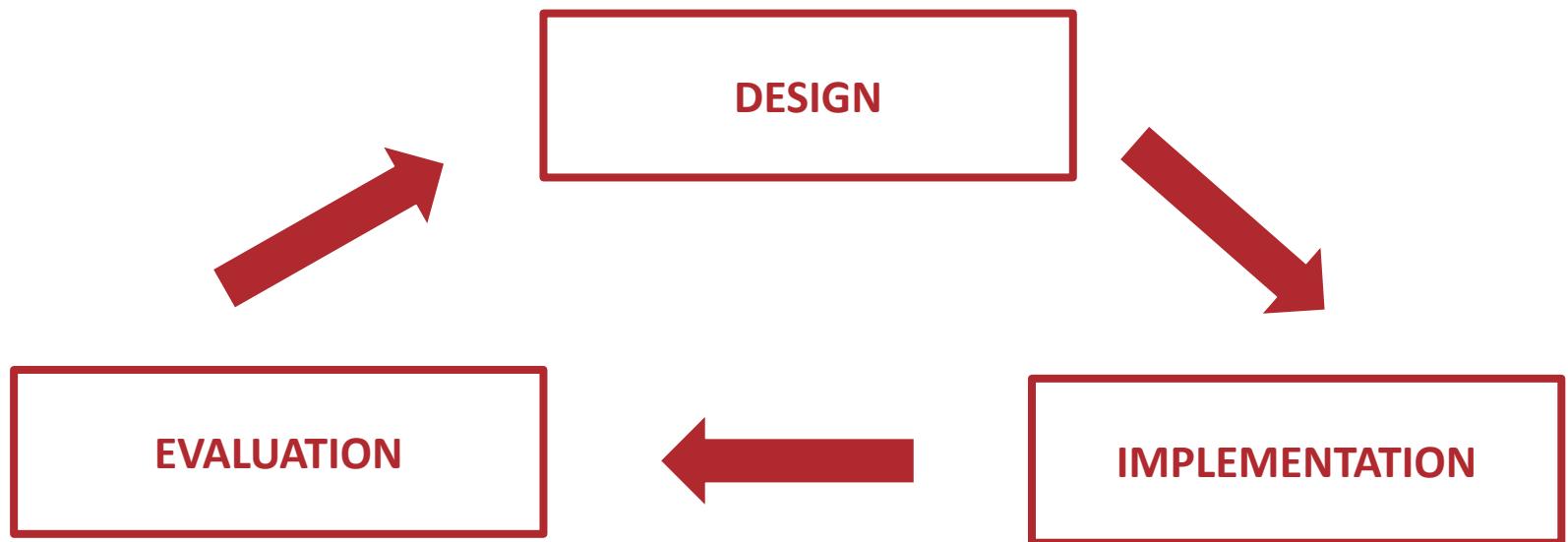
# Public policy cycle





# Question

What role can open data play in each stage of the policy cycle?





# The city of Xalapa's waste management problem

Rubbish collection was frustrating the public as there was no clarity how the service operated.

The main problems:

- Where is the rubbish collected from?
- When is the rubbish collection?
- What day is the rubbish collected?
- What time is the rubbish collected?

Create a policy cycle/s

What role can open data play?





1. Put a GPS on every garbage truck (**design**)
2. Collected the data (**implementation**)
3. Opened the data (**implementation**)
4. Engaged community in design (**design**)
5. Analyzed the maps for new routes (**evaluation & design**)
6. Implemented new routes (**implementation**)
7. Evaluated the impact, repeating 1-3 (**evaluation**)
8. Civil society organizations organized a Hackathon (**design**)
9. Developed an app (exact location, timetables, citizens reports) (**implementation**)
10. Evaluated the solution (**evaluation**)





1. Put a GPS on every garbage truck (**design**)
2. Collected the data (**implementation**)
3. Opened the data (**implementation**)
4. Engaged community in design (**design**)
5. Analyzed the maps for new routes (**evaluation & design**)
6. Implemented new routes (**implementation**)
7. Evaluated the impact, repeating 1-3 (**evaluation**)
8. Civil society organizations organized a Hackathon (**design**)
9. Developed an app (exact location, timetables, citizens reports) (**implementation**)
10. Evaluated the solution (**evaluation**)

Is open data  
just an input  
in this cycle?

If not what  
does this  
change?



# 1. Input data for evidence based policy making

Data is drawn from a number of sources and is analysed to inform policy making.

Examples include:

- environmental impact of third runway at Heathrow
- impact of London fire station closures
- how to regulate peer to peer lending



flickr: lucianf



## 2. Output data for transparency and to encourage citizen interaction

Data is published as part of a transparency or other open government agenda. There is no immediate desire for the data to have any other impact.

Examples include:

- government spending data
- planning application data
- LIDAR data





## 3. Tool to change behaviour

Where the data is the catalyst for change required by the policy.

Examples include:

- plastic bag usage data
- waste and emissions data
- pay gap data
- mobile coverage data
- broadband speed data



US Fish and Wildlife Service



Civil Service  
Learning

## Thank you

We hope you enjoyed this  
experience brought to you by

Delivered by



Civil Service  
Learning

Now over to you! What are you going to do differently?