# Open Data in Practice

http://training.theodi.org/InPractice

David Tarrant · @davetaz

# Day 3: Session 1
# The 80/20 of open data

# DATABLOG
## Facts are sacred

# Anyone can do it. Data journalism is the new punk

Can anyone be a data journalist? **Simon Rogers** on what we can learn from a 1977 diagram

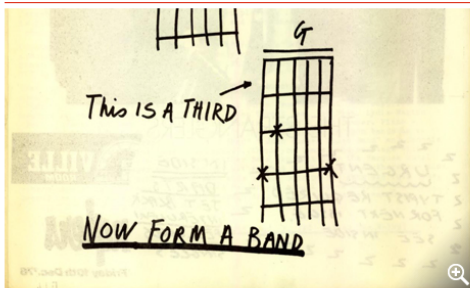• Another view: What data can and cannot do by Jonathan Gray

Page two of Sideburns, January 1977

Posted by
Simon Rogers
Thursday 24 May 2012
13.00 BST
theguardian.com
Jump to comments (8)

Article history

**Media**
Data journalism · Open journalism

**Technology**

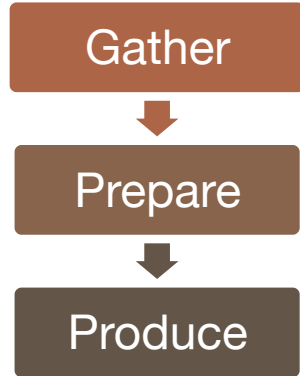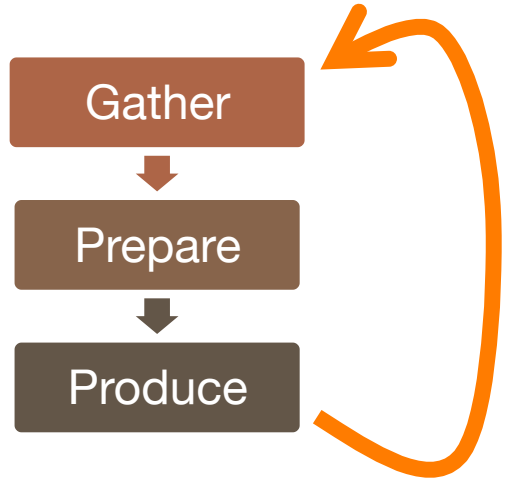❝ This is a chord… this is another… this is a third. NOW FORM A BAND

The data percolator

# Data percolation: A model of data preparation and analysis



```
Gather
  ↓
Prepare
  ↓
Produce
```

# Data percolation: A model of data preparation and analysis

# How should I budget my time?

| | |
|---|---|
| 20% | Gather |
| 60% | Prepare |
| 20% | Produce |

Finding a story is a creative process.

Let it percolate!

SENT DATA

BREAKING NEWS

RECURRING EVENTS

SHARE DATA

THEORIES TO BE EXPLORED

WHAT TO COMPARE OR SHOW CHANGE

WHAT DOES THE DATA MEAN

WHAT OTHER DATA SETS TO USE WITH IT?

SPREAD SHEETS

SPREAD SHEETS

DATA IN WRONG FORMAT

MERGED CELLS

UNNECESSARY COLUMNS OF DATA

DATA MEASURED IN DIFFERENT UNITS

PERFORM CALCULATIONS ON THE DATA

RECALCULATE IF NEEDED

SANITY CHECK THE RESULTS

# PERFORM CALCULATIONS ON THE DATA

RECALCULATE
IF NEEDED

SANITY CHECK
THE RESULTS

# OUTPUT

OUR GRAPHICS
TEAM

STORY

FREE VIZ
TOOLS

JUST PUBLISH

GOOGLE FUSION
TABLE

# How should I budget my time?

**Gather**

↓

**Prepare**

↓

**Produce**

1.1 **FIND** reliable data sources

1.2 Understand your **RIGHTS**

1.3 Visualise and **UNDERSTAND** your data

2.1 **CLEAN** your data

2.2 **TRANSFORM** it where useful

2.3 **COMBINE** it with other data sets

3.1 **REDUCE** and find the story

3.2 Think and understand the **CONTEXT**

3.3 Do your results pass a **SENSE-CHECK**?

# Time planning

Gather

Produce

Prepare

2.1 **CLEAN**

2.2 **TRANSFORM**

2.3 **COMBINE**

2.4 **ENRICH**

2.5 **ANALYSE**

Validating and cleaning data
Inception points: cross filters, outliers and comparisons

# Prepare (Stage 1)

Prepare

2.1 **CLEAN**

2.2 **TRANSFORM**

2.3 **COMBINE**

2.4 **ENRICH**

2.5 **ANALYSE**

# Introducing Open Refine



http://openrefine.org

# Exercise

Validating and cleaning data exercise

Focus on sections:

    2) Multiple representations

    4) Summation records

    5) Mixed use of numerical scales

# Prepare

Prepare

2.1 **CLEAN**

2.2 **TRANSFORM**

2.3 **COMBINE**

2.4 **ENRICH**

2.5 **ANALYSE**

# Enriching data and using pivot tables

# opencorporates

The Open Database Of The Corporate World

# Aggregator/Enabler

## We have information on
## 70,597,888 companies

[                                    ] **SEARCH**

○ search companies   ○ search officers

### Filter by jurisdiction

| | |
|---|---|
| 1,298 | Abu Dhabi (UAE) |
| 144,755 | Alaska (US) |
| 40,157 | Albania |
| 899,455 | Arizona (US) |
| 46,537 | Aruba |
| 165,582 | Bahamas |
| 99,185 | Bahrain |
| 88,563 | Bangladesh |

### Just released:
### OpenCorporates API v0.3

Corporate network data, financial accounts, complex filters, and more. Read more

Get data access to over
60 million companies

### Announcing Open LEIs

Today, OpenCorporates announces a new sister website, Open LEIs, a user-friendly interface on the emerging Global Legal Entity Identifier System. Read more

**OPENLEIs**
A BETA VIEW ON THE LEI SYSTEM

### New! Just added: Open corporate network data

Read more about this important new feature

# Exercise



Enriching a dataset containing company names (e.g. transactions) with company data from OpenCorporates

# Sense-checking

The best way to sense-check is to get a second pair of eyes to help you.

Any stories of common mistakes you'd like to share?

Infographics: The good the bad and the ugly
Charts: What to use and when
Colour, popout and order

# Session 4
# Infographics

# Recap

Gather

Produce

Prepare

3.1 **REDUCE**

3.2 **CONTEXT**

3.3 **SENSE**

Infographics (separate presentation)

Thank-you