# Reproducible data science to inform outbreak response

## Report from the North-Kivu Ebola outbreak

**Dr Thibaut Jombart** (@TeebzR)
London School of Hygiene and Tropical Medicine
Imperial College London
R Epidemics Consortium (RECON)

UserR!2019
Toulouse
11 July 2019

# Context

Outbreak analytics
RECON
Ebola in North Kivu

# On the emergence of "outbreak analytics"



https://doi.org/10.1098/rstb.2018.0276

- **DoB**: Polonsky et al. (2019) Phil. Trans. R. Soc. B 374

- **Data science** mixing statistics, mathematical modeling, computer simulations, database infrastructure, GIS, genetics, software engineering

- At the crossroad of **public health institutions**, **private sector**, and **academia**

- Aims to **inform response to emergencies in real-time**

- **Lack of available tools**

# RECON: bringing data science into health emergencies



https://www.repidemicsconsortium.org/
https://www.reconlearn.org/

- **Origin**: Hackout 3 (rOpenSci / Imperial College London), Berkeley, 2016

- NGO for open analytics resources for health emergencies and humanitarian crises

- ~35 members, 200-300 subscribers

- Packages: 10 on CRAN, 15-20 in development

- Events: short courses, workshops, hackathons

- **Deployments** to support response to emergencies

# Ebola in North-Kivu & Ituri, DRC







- Largest Ebola epidemic in DRC, 2nd largest in the world

- August 2018 - today:
  - >2400 cases (confirmed / probable)
  - 67% deaths

- Difficulties due to military conflicts
  - Threats to local population
  - Threats to response staff and facilities

- **First deployment of an analytical cell** as part of the Emergency Operations Centre

# Outbreak analytics cell: aims and challenges



Courbe épidemique: cas confirmés et probables par date de notification



Projections de nouveaux cas sur 21 jours - Mangina



- Multiple (messy) data sources, no global database

- Independent updates of different databases

- Needs: data cleaning, visualisation, in-depth analyses, forecasting

- Routine versus *ad-hoc* analyses

- Need for regular results updates and traceability

- Bad internet, different platforms, low R literacy

# Tidier rmarkdown workflows with *reportfactory* : use case



## Original requirements

- Handle multiple `.Rmd` reports
- Handle multiple (dated) versions of the same report
- Separate data, scripts, `.Rmd` sources, outputs
- Generates time-stamped outputs
- Update all reports in one go
- Handle dependencies on packages
- Non-invasive: use of standard `.Rmd`, no config file
- Easy to use: accessible by people new to R
- Offline: does not require internet
- Portable: work on any platform

# What does the *reportfactory* do?



How to handle analysis reports with different versions, each compiled several times?
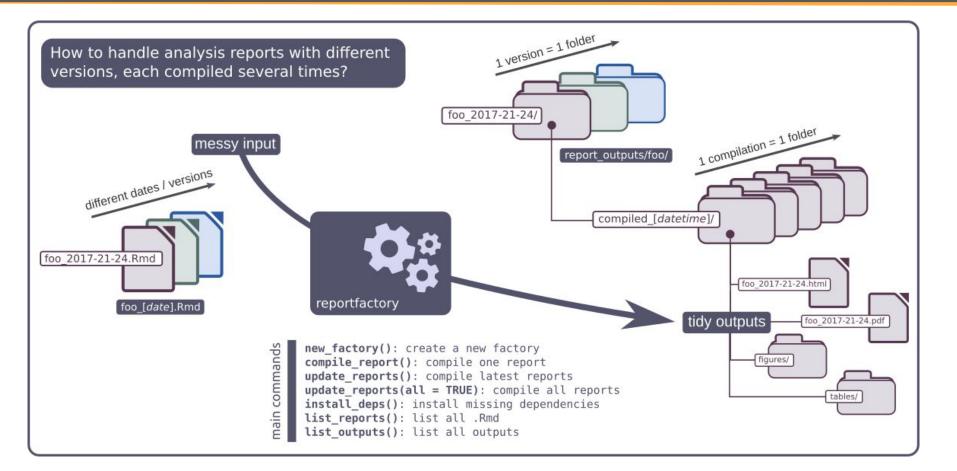
messy input

different dates / versions

foo_2017-21-24.Rmd

foo_[date].Rmd

reportfactory

1 version = 1 folder

foo_2017-21-24/

report_outputs/foo/

compiled_[datetime]/

1 compilation = 1 folder

tidy outputs

foo_2017-21-24.html

foo_2017-21-24.pdf

figures/

tables/

main commands

```
new_factory(): create a new factory
compile_report(): compile one report
update_reports(): compile latest reports
update_reports(all = TRUE): compile all reports
install_deps(): install missing dependencies
list_reports(): list all .Rmd
list_outputs(): list all outputs
```

# *reportfactory* : basic structure

Creating a new factory: `new_factory()`



.`Rmd` files

Outputs:
html files,
figures etc

Data

.`R` scripts

Open the
factory

Anchor
(for file paths)

**Other functionalities**

- List / install dependencies: `list_deps() / install_deps()`

- List reports: `list_reports()`

- Compile all recent reports: `update_reports()`

- Compile specific report: `compile_report()`

- Archive old reports: `archive_reports()`

- … : **contributions welcome!**

  **(join us, we have cookies)**

# Data standardisation using *linelist*

```
x %>% clean_data()
```

Capitalisation
Accents
Separators
Dates

| 'ID | Date of Onset. | GENDER_ | Épi.Case_définition | messy/dates |
|---|---|---|---|---|
| khdntz | 2018-01-09 | male | Confirmed | that's 24/12/1989! |
| hmckhn | 2018-01-09 | male | suspected | // 24//12//1989 |
| ekjmyd | 2018-01-09 | Female | confirmed | that's 24/12/1989! |
| kmoczh | 2018-01-04 | MALE | suspected | female |
| kftifx | 2018-01-02 | FEMALE | suspected | // 24//12//1989 |
| qyipse | 2018-01-09 | Male | PROBABLE | 01-12-2001 |
| zprzec | 2018-01-03 | male | suspected | NA |
| bgsmfn | 2018-01-06 | Female | suspected | that's 24/12/1989! |
| syfnfd | 2018-01-05 | Female | confirmed | 01-12-2001 |
| aekdlv | 2018-01-07 | FEMALE | not a case | female |
| kcejly | 2018-01-05 | Female | Confirmed | that's 24/12/1989! |
| jyxnhl | 2018-01-11 | female | confirmed | // 24//12//1989 |

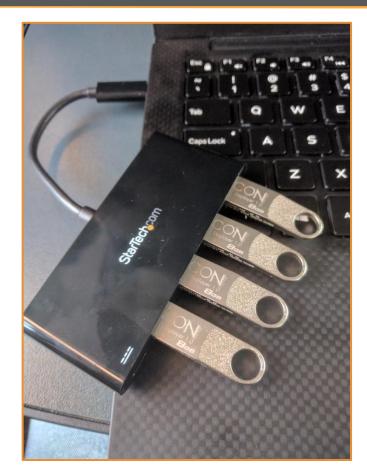| id | date_of_onset | gender | epi_case_definition | messy_dates |
|---|---|---|---|---|
| khdntz | 2018-01-09 | male | confirmed | 1989-12-24 |
| hmckhn | 2018-01-09 | male | suspected | 1989-12-24 |
| ekjmyd | 2018-01-09 | female | confirmed | 1989-12-24 |
| kmoczh | 2018-01-04 | male | suspected | NA |
| kftifx | 2018-01-02 | female | suspected | 1989-12-24 |
| qyipse | 2018-01-09 | male | probable | 2001-12-01 |
| zprzec | 2018-01-03 | male | suspected | NA |
| bgsmfn | 2018-01-06 | female | suspected | 1989-12-24 |
| syfnfd | 2018-01-05 | female | confirmed | 2001-12-01 |
| aekdlv | 2018-01-07 | female | not_a_case | NA |
| kcejly | 2018-01-05 | female | confirmed | 1989-12-24 |
| jyxnhl | 2018-01-11 | female | confirmed | 1989-12-24 |

# Dictionary-based cleaning using *linelist*

`x %>% clean_data(wordlists = rules)`

Typos
Re-levelling
Variable-specific
rules

| 'ID | Date of Onset. | GENDER_ | Épi.Case_définition |
|-----|------|---------|---------------------|
| hlywxf | 2018-01-10 | m | ConFRImed |
| zgsjfx | 2018-01-05 | man | NA |
| nbmrvn | 2018-01-08 | female | NA |
| fasshf | 2018-01-02 | male | suspected |
| wlfhgk | 2018-01-03 | f | Not.a.Case |
| qdmhyp | 2018-01-08 | NA | Confirmed |
| ywntgm | 2018-01-03 | male | not a case |
| vlpamu | 2018-01-04 | male | PROBABLE |
| fqigws | 2018-01-02 | MALE | Not.a.Case |
| vrzpkj | 2018-01-06 | Female | confirmed |
| gsbjak | 2018-01-06 | f | female |
| zozxjp | 2018-01-11 | f | male |

## rules

| change | to | variable |
|--------|-----|----------|
| m | male | gender |
| f | female | gender |
| man | male | gender |
| .missing | unknown | .global |
| confrimed | confirmed | epi_case_definition |
| female | unknown | epi_case_definition |
| male | unknown | epi_case_definition |

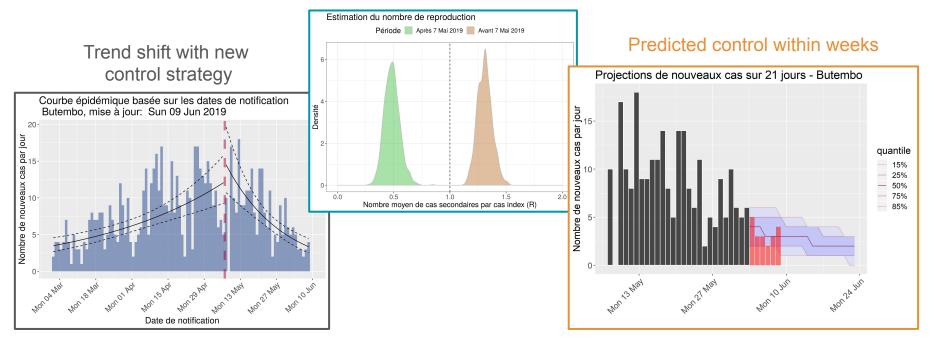| id | date_of_onset | gender | epi_case_definition |
|-----|------|--------|---------------------|
| hlywxf | 2018-01-10 | male | confirmed |
| zgsjfx | 2018-01-05 | male | unknown |
| nbmrvn | 2018-01-08 | female | unknown |
| fasshf | 2018-01-02 | male | suspected |
| wlfhgk | 2018-01-03 | female | not_a_case |
| qdmhyp | 2018-01-08 | unknown | confirmed |
| ywntgm | 2018-01-03 | male | not_a_case |
| vlpamu | 2018-01-04 | male | probable |
| fqigws | 2018-01-02 | male | not_a_case |
| vrzpkj | 2018-01-06 | female | confirmed |
| gsbjak | 2018-01-06 | female | unknown |
| zozxjp | 2018-01-11 | female | unknown |

## The RECON deployer

- USB stick with latest R, Rtools, Rstudio for Windows, MacOSX, Linux

- Local package repository - instance of *nomad*: https://github.com/reconhub/nomad

- ~2000-3000 CRAN packages

- ~10-20 github packages

- Cheatsheets

- Website: https://github.com/reconhub/deployer

# Making a difference

Showing what works
Join the movement

# Join the movement!



**Outbreak analytics**
- Still an emerging field
- Funding and training gaps
- Data scientists needed!

**The good stuff**
- Help respond to health emergencies and humanitarian crises
- Work with visible impact
- Exciting data challenges
- Lots of potential for capacity building: the **next generation of data scientists needs to be in-country**

# Thanks to

Get these slides

LONDON SCHOOL of HYGIENE &TROPICAL MEDICINE

Imperial College London

RECON