

# Show me the data you didn't consider!

## Reducing unknown unknowns in data science

Claus Thorn Ekstrøm  
UCPH Biostatistics

[ekstrom@sund.ku.dk](mailto:ekstrom@sund.ku.dk)

Ann Arbor R Users meeting, March 9th

Slides @ [biostatistics.dk/talks/](https://biostatistics.dk/talks/)







The long read

# How statistics lost their power – and why we should fear what comes next

The ability of statistics to accurately represent the world is declining. In its wake, a new age of big data controlled by private companies is taking over – and putting democracy in peril

p-value





**I CAN'T KEEP  
CALM**

**because of  
SOCIOLOGY**





# WE'RE OUT

» Britain votes to quit the EU » Pound goes into freefall

● **Britain waking to an EU exit**

THE 2016 REFERENDUM JUDGMENT, JULY 1, 2016  
Britain, said, this morning, at 5am, of the day, the 24th of June, a historic referendum day for the British people. The result, 52% to 48%, was a clear and decisive victory for the 'Leave' campaign. The result, which was announced at 10pm, was a surprise to many, but it was a result that the 'Leave' campaign had fought hard to achieve. The result, which was announced at 10pm, was a surprise to many, but it was a result that the 'Leave' campaign had fought hard to achieve. The result, which was announced at 10pm, was a surprise to many, but it was a result that the 'Leave' campaign had fought hard to achieve.

● **Leave daims win in huge poll**



# SEE EU LATER!



**Trump on brink of White House as rocks Hillary the world edge poll**

**Washington Standard**

**BILLIONAIRE PULLS OFF HUGE UPSET TO BECOME PRESIDENT**

**HE PLEDGES TO BIND WOUNDS OF DIVISION AS CLINTON CONCEDES**

**PUTIN AND MAY SEND MESSAGES OF CONGRATULATION**



**Clinton train crash: horror, several dead**

# TRUMP TRIUMPH SHOCKS WORLD

REPORTS AND ANALYSIS

COMMENT



# Quiz



# The life of a data scientist

Data scientists, according to interviews and expert estimates, spend from **50 percent to 80 percent** of their time mired in this more mundane labor of collecting and preparing unruly digital data, before it can be explored for useful nuggets.

-- "For Big-Data Scientists, 'Janitor Work' Is Key Hurdle to Insight" - The New York Times, 2014

An iceberg floating in a blue ocean under a blue sky. The tip of the iceberg is above the water line, while the much larger base is submerged. The text 'In reality' and 'BIG DATA' are overlaid on the submerged part of the iceberg. To the right of the iceberg, there are two lines of text: 'WHAT WE KNOW...' and 'THE REST...', both preceded by a left-pointing arrow.

◀ WHAT WE KNOW...

◀ THE REST...

In reality  
**BIG DATA**



# Beyond the hype

Many of the existing problems with small data are also applicable to big data. ...

The problems do not disappear because the data sizes becomes larger. They become **worse**.

# State of the art?

## Statistical analysis

All of the data were analyzed with data processing software and figures with Microsoft excel 2007.

-- Tayefe *et al*, Advances in Bioresearch, 2014



A close-up photograph of an elderly person's hand, showing wrinkled skin and short, slightly discolored fingernails. The hand is resting on a dark, textured wooden surface. To the right of the hand is a glass filled with a dark amber liquid, possibly whiskey. The background is blurred, showing a white, textured surface. The text "The RESCueH project" is overlaid in white, sans-serif font across the middle of the hand.

# The RESCueH project

# Timeline follow back

	day1	day2	day3
1	18	NA	NA
2	14	NA	99
3	20	17	40
4	23	14	17

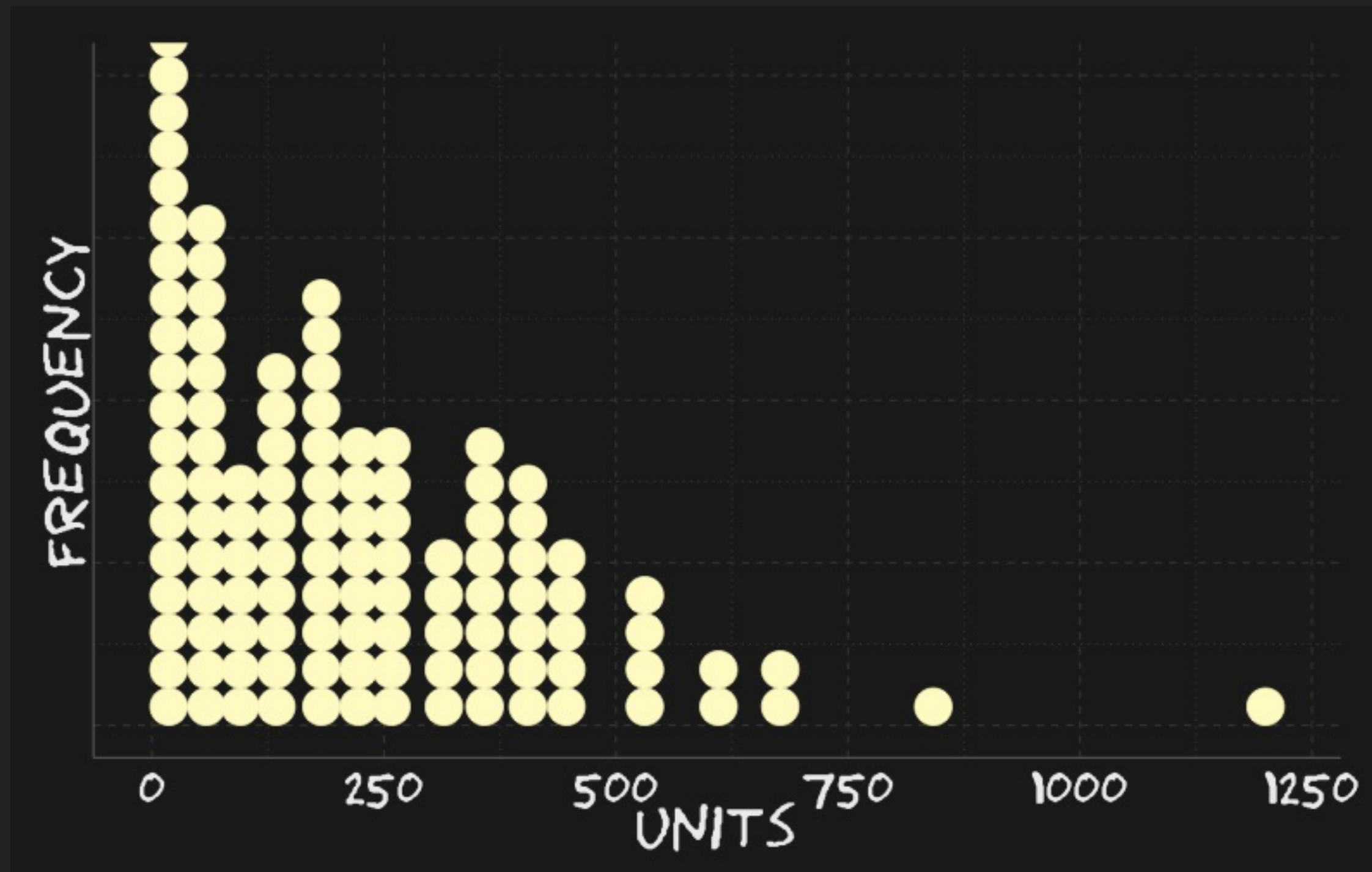


# Timeline follow back

	day1	day2	day3
1	18	NA	NA
2	14	NA	99
3	20	17	40
4	23	14	17

	day1	day2	day3
1	18	NA	NA
2	14	NA	99
3	20	17	40
4	23	14	17
5	10	24	2
6	19	88	8

# Monthly Alcohol units





# Reproducible research

What **didn't** we check?

# Reproducible research

What **didn't** we check?

- Many studies cannot be replicated: time, money, unique
- New technologies increase data sizes
- Merge existing databases into megadatabases
- Genetic data for future analyses

# Reproducible research

What **didn't** we check?

- Many studies cannot be replicated: time, money, unique
- New technologies increase data sizes
- Merge existing databases into megadatabases
- Genetic data for future analyses

You are your worst collaborator.

# dataMaid

```
devtools::install_github("ekstroem/dataMaid")  
library(dataMaid)  
data(toyData)  
clean(toyData)
```

Documentation to be **read** and **evaluated** by a human.

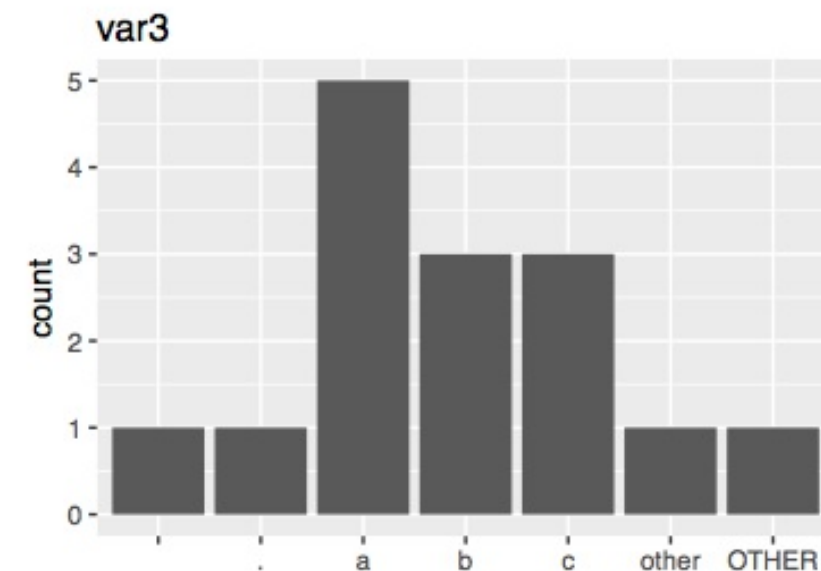
See [github.com/ekstroem/dataMaid](https://github.com/ekstroem/dataMaid) for more info.



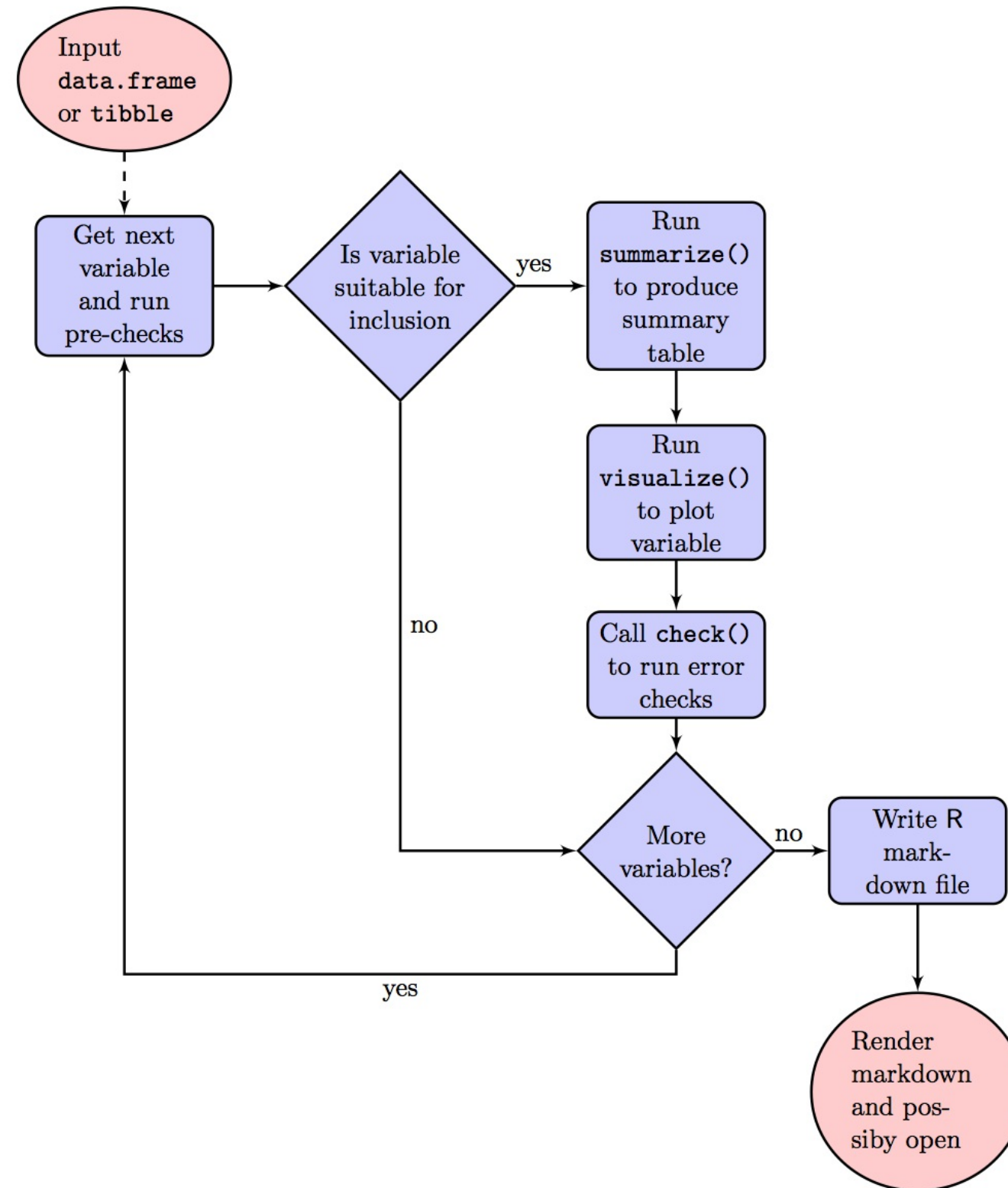
	character	factor	labelled	numeric	integer	logical	Date
Identify miscoded missing values	×	×	×	×	×		
Identify prefixed and suffixed whitespace	×	×	×				
Identify levels with < 6 obs.	×	×	×				
Identify case issues	×	×	×				
Identify misclassified numeric or integer variables	×	×	×				
Identify outliers				×	×		×

## var3

Feature	Result
Variable type	factor
Number of missing obs.	0 (0 %)
Number of unique values	7
Mode	"a"



- The following suspected missing value codes enter as regular values: " ", ".".
- The following values appear with prefixed or suffixed white space: " ".
- Note that the following levels have at most five observations: " ", ".", "a", "b", "c", "other", "OTHER".
- Note that there might be case problems with the following levels: "other", "OTHER".



# Using dataMaid interactively

```
check(toyData$var2)    # Individual check
check(toyData$var2, numericChecks = "identifyMissing")
visualize(toyData$var2)
summarize(toyData$var2)
summarize(toyData$var2,
          numericSummaries = c("centralValue", "minMax"))
```



# Extending dataMaid

```
isID <- function(v, nMax = NULL, ...) {  
  out <- list(problem = FALSE, message = "")  
  if (class(v) %in% setdiff(allClasses(),  
                             c("logical", "Date"))) {  
    v <- as.character(v)  
    lengths <- c(nchar(v))  
    if (all(lengths > 10) &  
        length(unique(lengths)) == 1) {  
      out$problem <- TRUE  
      out$message <- "Warning: Seems to contain IDs."  
    }  
  }  
  out }
```

# Adding the function

```
data("exampleData")
exampleData$names <- sapply(1:300,
  function(i) { paste0(sample(LETTERS, size=10),
                        collapse="") })
clean(exampleData,
  preChecks = c("isID"))
```

Communicate!