



## EXPLORING DATA SCIENCE WITH »R«

Created: 2021-05-19 Mi 19:47

## WHAT WILL YOU LEARN?

---

Two ways of looking at data

---

Data science  $\leftarrow$  modeling  $\cup$  coding

---

Data science examples using R

---

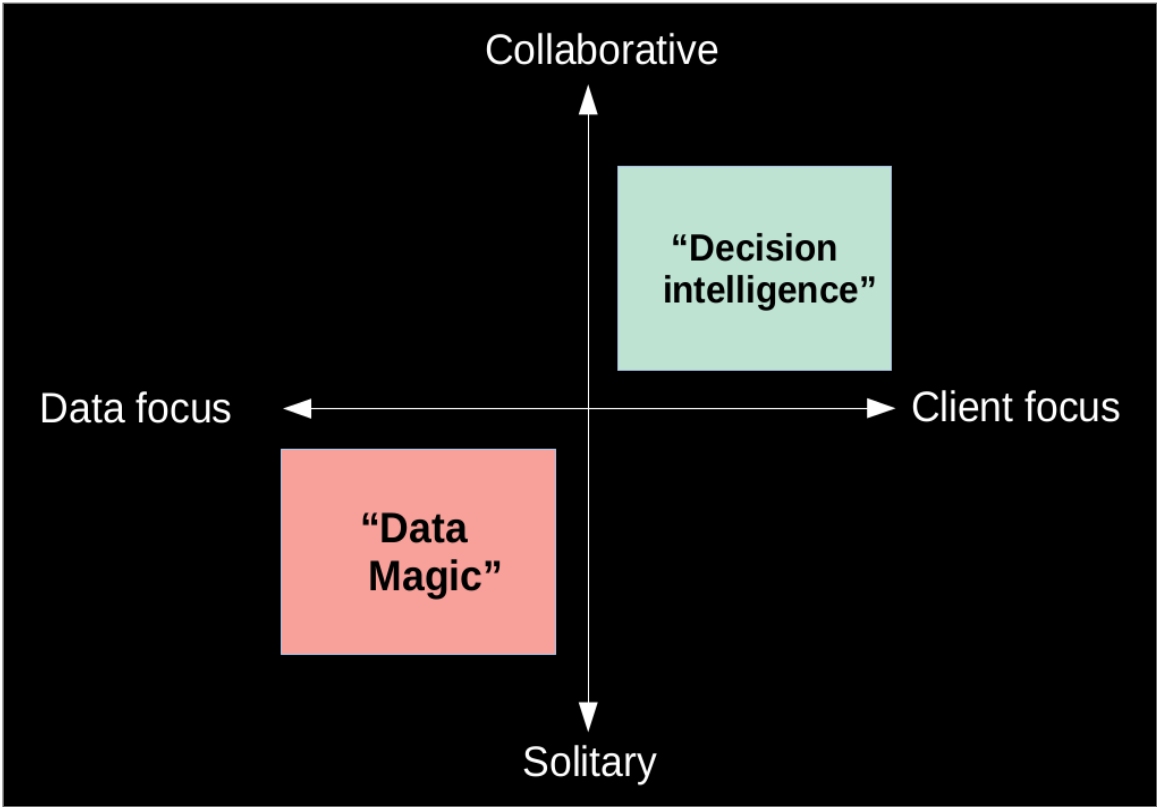
R vs. Python for data science

---

## WHAT CAN YOU DO?

- Answer polls, like "**do you know R?**"
- Leave comments/questions in the chat
- Download the slides (PDF) **here**
- Learn R, e.g. **here**

# TWO WAYS OF LOOKING AT DATA



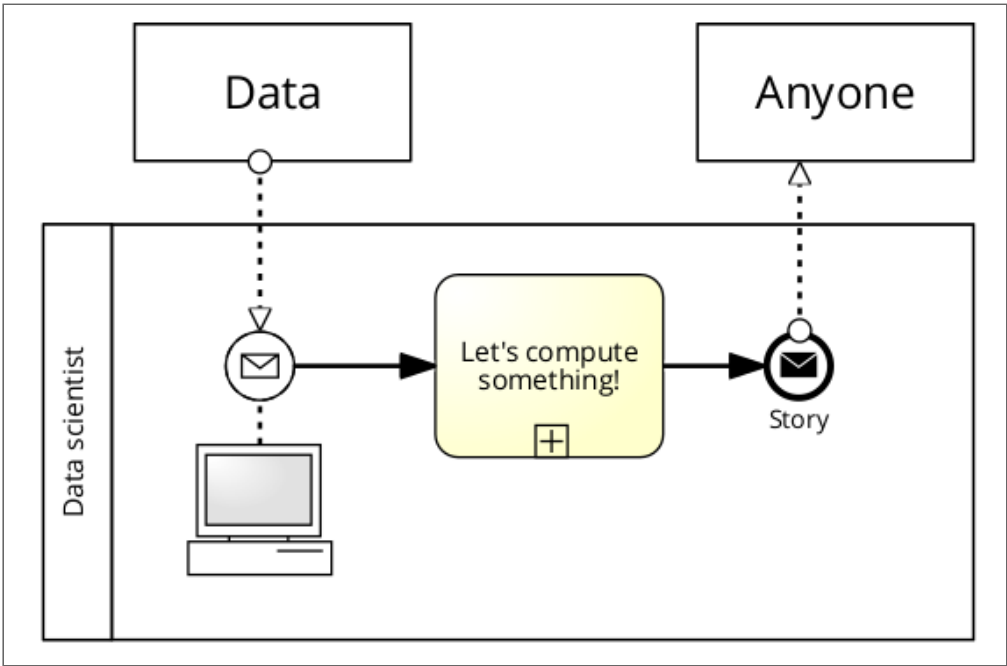
## SCENARIO 1: "DATA MAGIC"

- We've got data!
- We've got computers!
- Let's compute something!

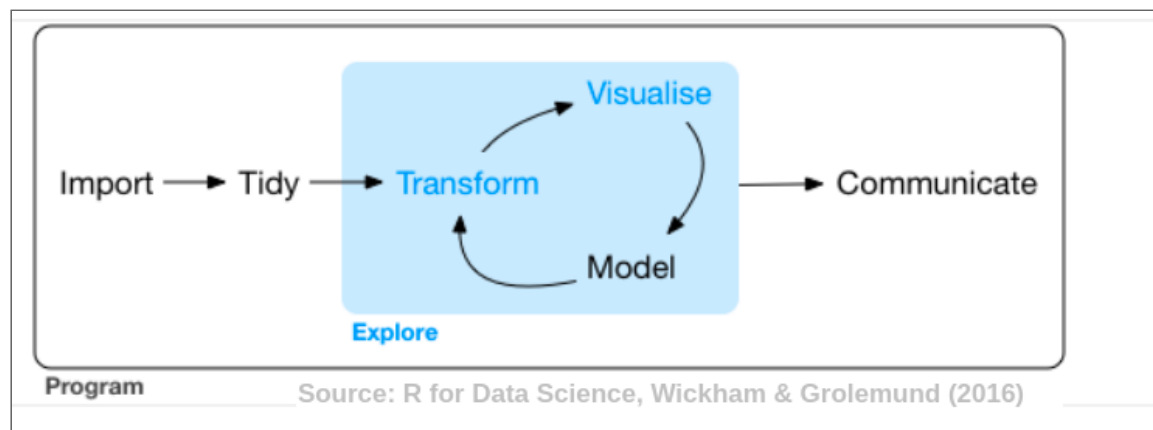
## GOOD QUESTION, BAD SETUP

***What's the story behind the data?***

"DATA MAGIC" WORKFLOW



## EXPLORATORY DATA ANALYSIS





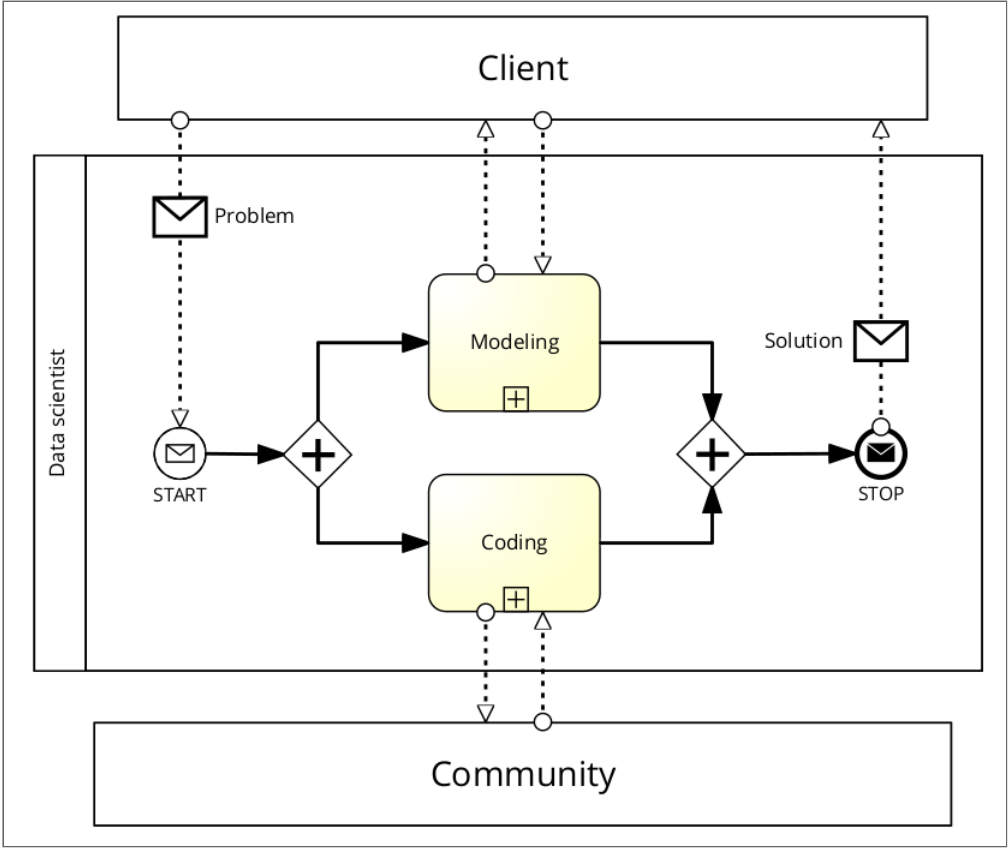
## **SCENARIO 2: "DECISION INTELLIGENCE"**

- A client got a problem!
- A client needs a solution!

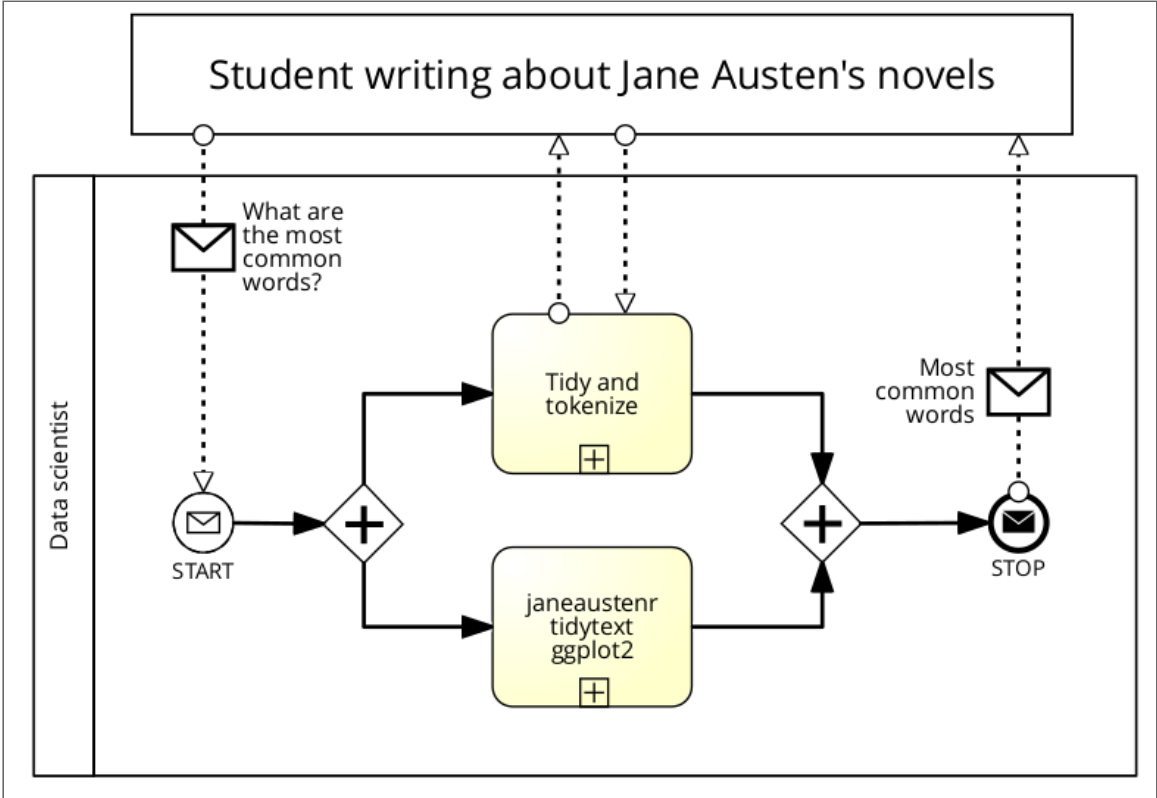
## BETTER QUESTION, GOOD SETUP

***What are your options?***

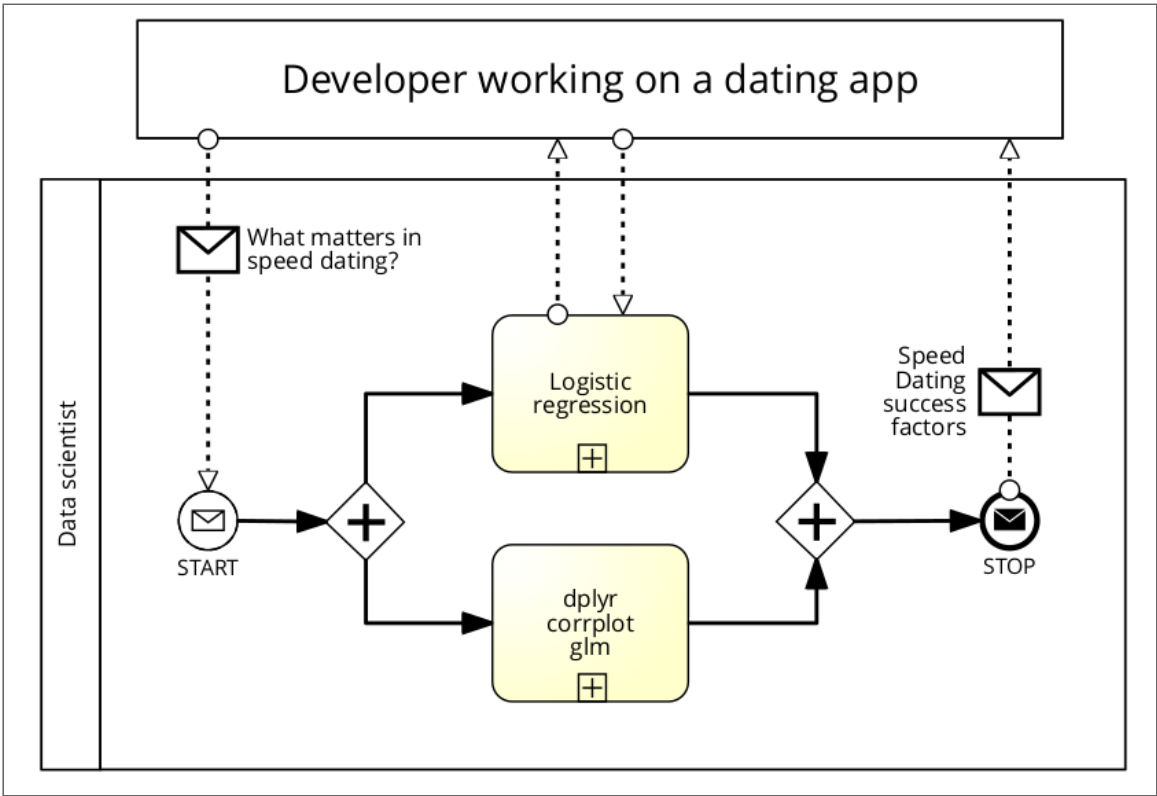
# PROBLEM-BASED COLLABORATIVE EXPLORATORY DATA ANALYSIS



EXAMPLE: TEXT MINING



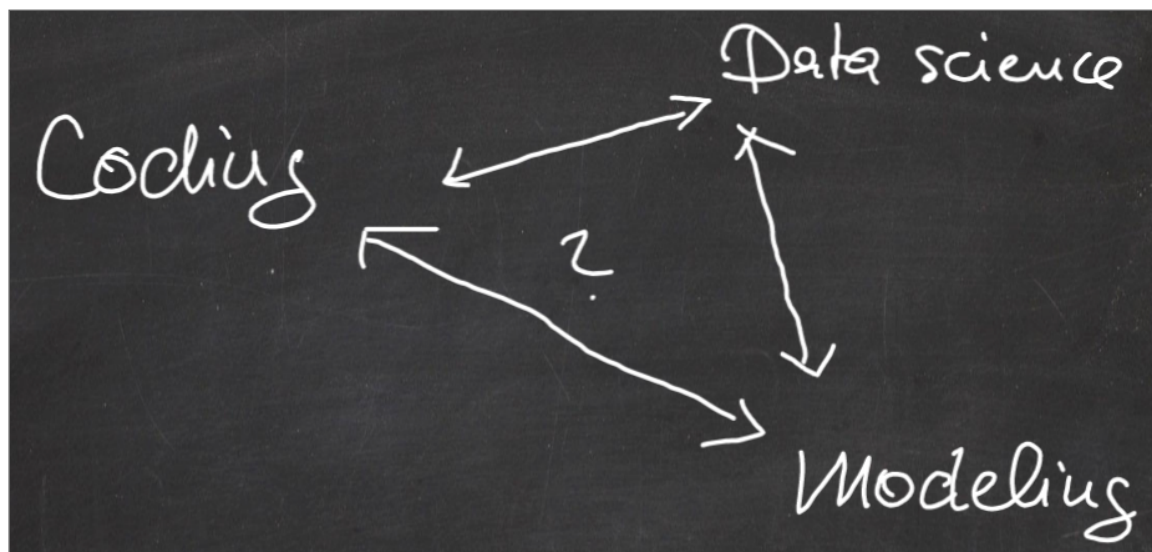
EXAMPLE: SPEED DATING



## SUMMARY

- "Data magic" vs. "decision intelligence"
- Data/solitary vs. Client/collaborative
- Text mining and recommender systems
- Modeling/coding as collaborative activities

**DS**  $\leftarrow$  **MODELING**  $\cup$  **CODING**



## CODING SKILLS

- Algorithmic thinking
- Fixed at run-time
- Procedure oriented
- Depends on data structures
- Object-oriented programming (OOP)
- Functional Programming (FP)



## MODELING SKILLS

- Heuristic thinking
- Adaptive at run-time
- Pattern oriented
- Cognitive bias
- Analogy, induction
- Petri nets, BPMN, UML

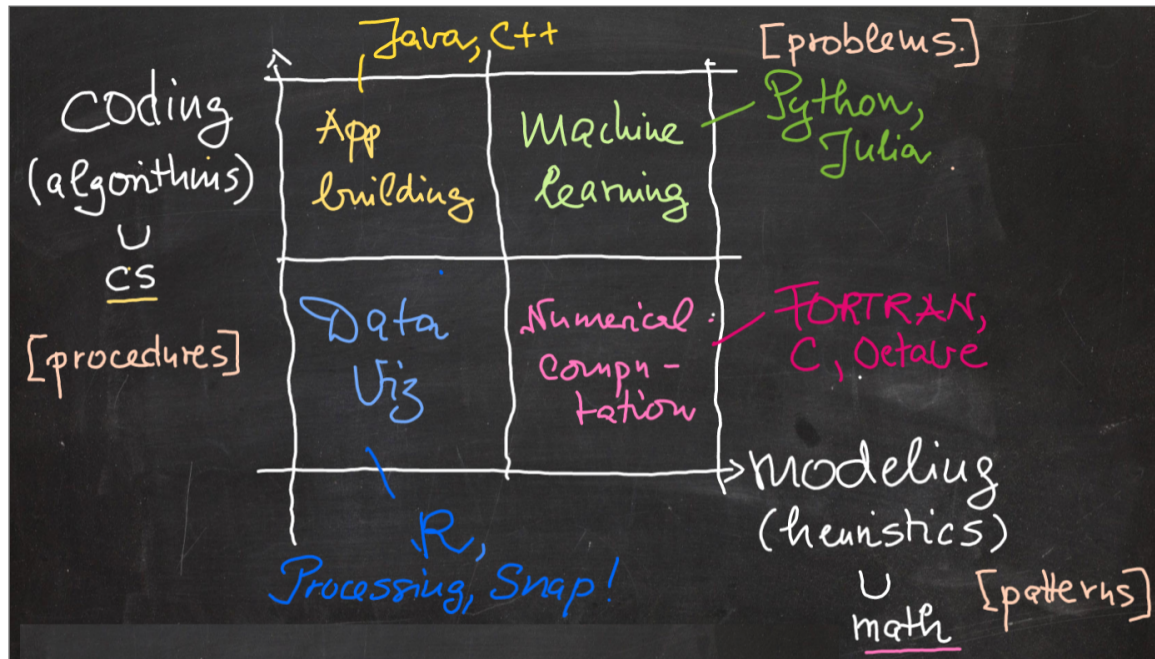
## FOUR SCENARIOS

- Data Visualization
- Application Building
- Numerical Computation
- Machine Learning

## PROGRAMMING LANGUAGES

- Apps: Java, C++
- Numerics: FORTRAN, C, Octave
- DataViz: R / Python, Snap!, Processing
- Machine Learning: Python / R, Julia

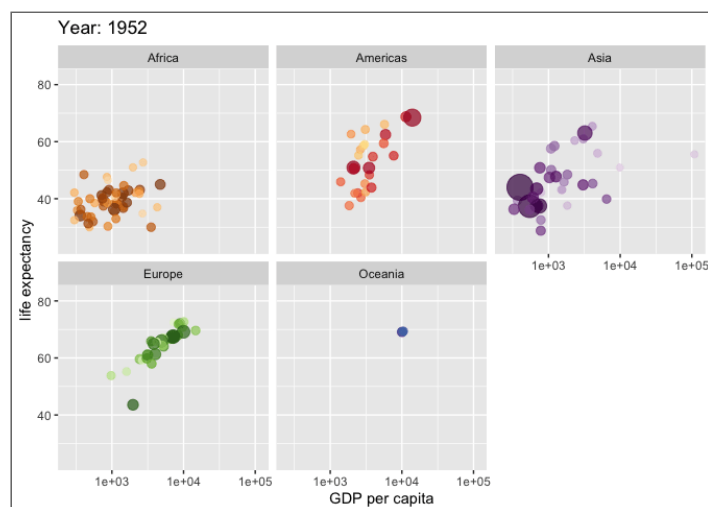
## DATA SCIENCE FRAMEWORK



## SUMMARY

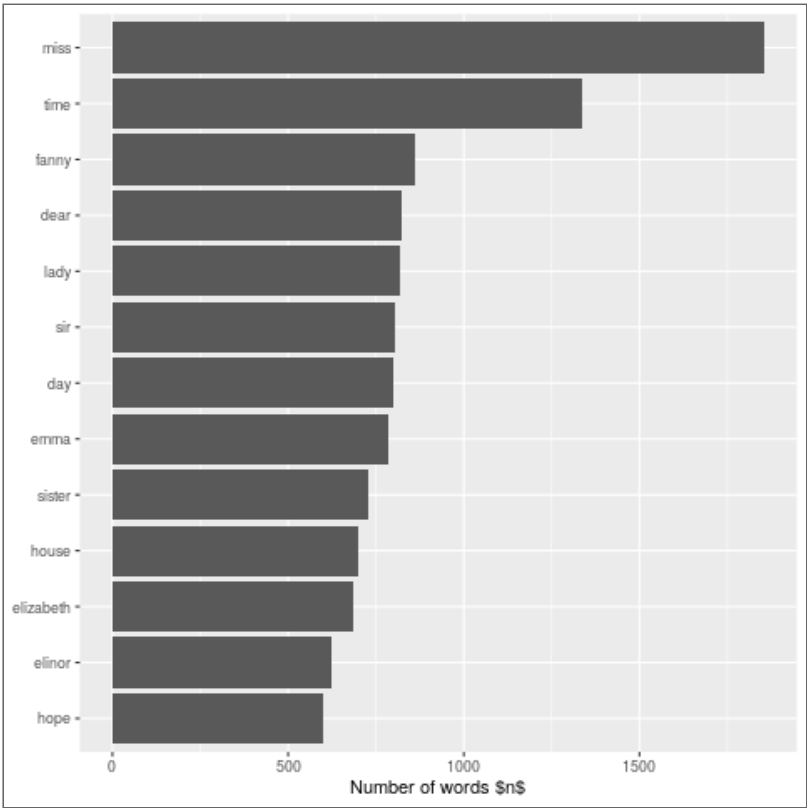
- Coding is algorithmic, modeling is heuristic
- Paradigms: OOP/FP, BPMN/UML
- Scenarios: DataViz, Apps, Numerics, ML
- Languages: R, Python, and many more
- Two pathways: CS/coding vs. DS/modeling

## DATA SCIENCE EXAMPLES USING R



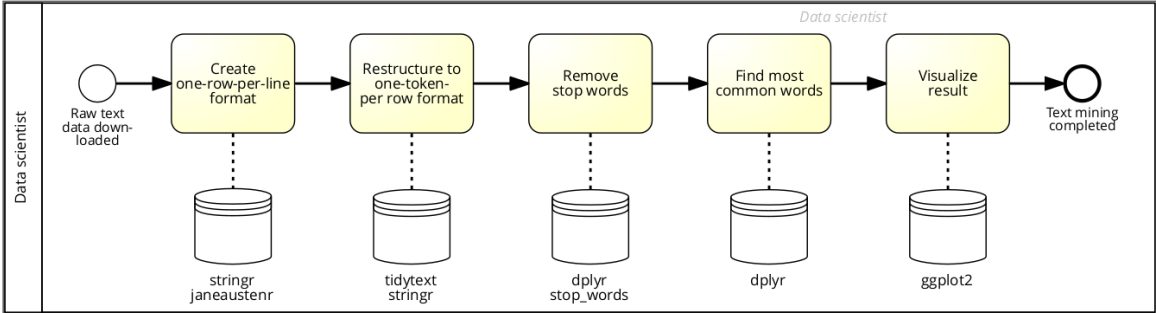
- Text mining Jane Austen
- What matters in speed dating
- Current student projects

# TEXT MINING JANE AUSTEN



The most common words in Jane Austen's novels

# TEXT MINING JANE AUSTEN: STEPS



BPMN diagram showing required R libraries



## TEXT MINING JANE AUSTEN: CODE

```
library(janeaustenr)
library(dplyr)
library(stringr)
library(tidytext)
library(ggplot2)
data(stop_words)

original_books <- austen_books() %>%
  group_by(book) %>%
  mutate(linenumber = row_number(),
         chapter = cumsum(str_detect(text,
                                     regex("^chapter [\\divxlc]",
                                           ignore_case = TRUE)))) %>%
  ungroup()

tidy_books <- original_books %>%
  unnest_tokens(word, text)

tidy_books <- tidy_books %>%
  anti_join(stop_words)

tidy_books %>%
  count(word, sort = TRUE)

tidy_books %>%
  count(word, sort = TRUE) %>%
  filter(n > 600) %>%
  mutate(word = reorder(word, n)) %>%
  ggplot(aes(n, word)) +
  geom_col() +
  labs(x="Most common words in Jane Austen's novels", y = NULL)
```

Source: Silge & Robinson (2017)

# REAL WORLD APPLICATION

**Arxiv Sanity Preserver**  
 Built in spare time by @karpathy to accelerate research.  
 Serving last 153012 papers from cs.[CV|CL|LG|AI|NE]/stat.ML

User:  Pass:  [Login or Create](#)


[Fork me on GitHub](#)

[most recent](#)
[top recent](#)
[top hype](#)
[friends](#)
[discussions](#)
[recommended](#)
[library](#)

**How we do things with words: Analyzing text as social and cultural data**

Dong Nguyen, Maria Liakata, Simon DeDeo, Jacob Eisenstein, David Mimno, Rebekah Tromble, Jane Winters  
 7/2/2019 [cs.CL](#)

1907.01468v1 [pdf](#)  
[show similar](#) | [discuss](#)




In this article we describe our experiences with computational text analysis. We hope to achieve three primary goals. First, we aim to shed light on thorny issues not always at the forefront of discussions about computational text analysis methods. Second, we hope to provide a set of best practices for working with thick social and cultural concepts. Our guidance is based on our own experiences and is therefore inherently imperfect. Still, given our diversity of disciplinary backgrounds and research practices, we hope to capture a range of ideas and identify commonalities that will resonate for many. And this leads to our final goal: to help promote interdisciplinary collaborations. Interdisciplinary insights and partnerships are essential for realizing the full potential of any computational text analysis that involves social and cultural concepts, and the more we are able to bridge these divides, the more fruitful we believe our work will be.

**End-to-End Optical Character Recognition for Bengali Handwritten Words**

Farisa Benta Safir, Abu Quwsar Ohi, M. F. Mridha, Muhammad Mostafa Monowar, Md. Abdul Hamid  
 5/9/2021 [cs.CV](#) | [cs.IR](#)

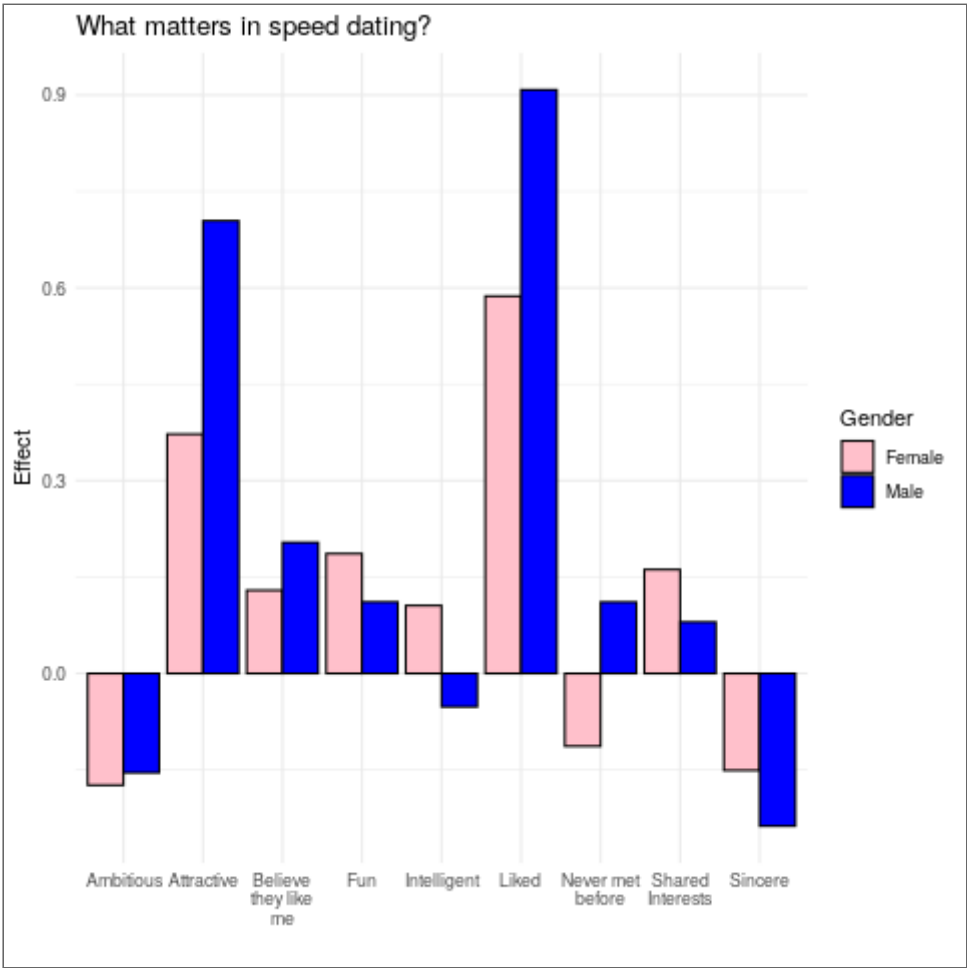
2105.04020v1 [pdf](#)  
[show similar](#) | [discuss](#)



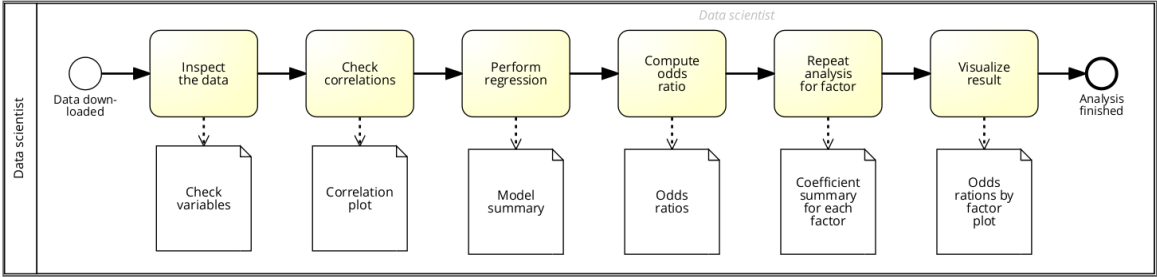
Accepted in "The 4th National Computing Colleges Conference"

## Arxiv Sanity Preserver (recommender)

# WHAT MATTERS IN SPEED DATING



WHAT MATTERS IN SPEED DATING: STEPS



BPMN diagram showing process & output

# WHAT MATTERS IN SPEED DATING: CODE

```
library(tidyverse)
library(corrplot)

download.file(
  "https://raw.githubusercontent.com/keithmcnulty/speed_dating/master/speed_data_data.RDS",
  "speed_dating_data.RDS")
data <- readRDS("speed_dating_data.RDS")

head(as.data.frame(data,3))

corr_matrix <- data %>%
  dplyr::select(amb, sinc, intel, fun, amb, shar, like, prob, met) %>%
  as.matrix()

M <- cor(corr_matrix, use="complete.obs")
corrplot::corrplot(M)

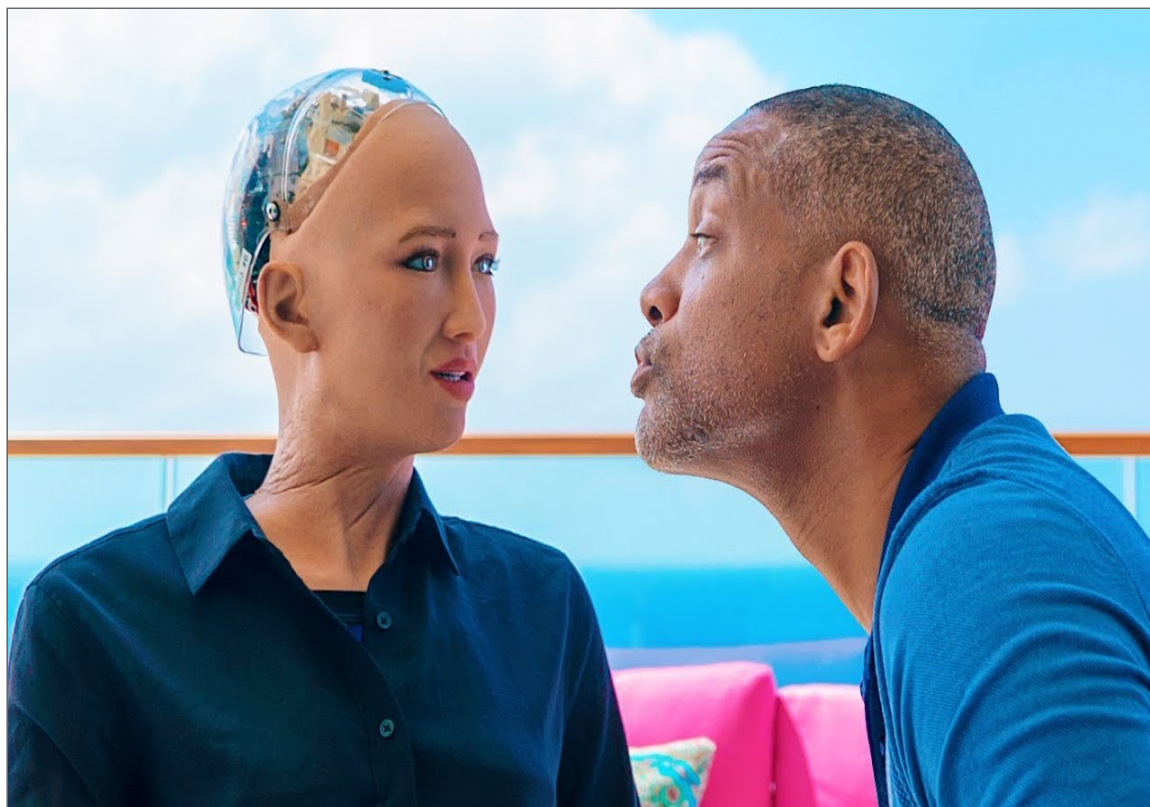
model_m <- glm(dec ~ attr+sinc+intel+fun+amb+shar+like+prob+met,
  data = data %>% dplyr::filter(gender == 1),
  family = 'binomial')
ctable_m <- coef(summary(model_m))
odds_ratio_m <- exp(coef(summary(model_m))[, c("Estimate")])
coef_summary_m <- cbind(ctable_m,
  as.data.frame(odds_ratio_m,
    nrow = nrow(ctable_m),
    ncol = 1))

chart_data <- coef_summary_m %>%
  tibble::rownames_to_column() %>%
  dplyr::left_join(coef_summary_m %>%
    dplyr::add_rownames(), by = "rowname") %>%
  dplyr::select(rowname, odds_ratio_f, odds_ratio_m) %>%
  tidyr::pivot_longer(cols = c("odds_ratio_f", "odds_ratio_m"),
    names_to = "odds_ratio") %>%
  dplyr::mutate(Effect = value - 1,
    Gender = ifelse(odds_ratio == "odds_ratio_f", "Female", "Male"),
    Factor = dplyr::recode(rowname,
      amb = "Ambitious", attr = "Attractive",
      fun = "Fun", intel = "Intelligent",
      like = "Liked", met = "Never met\nbefore",
      prob = "Believe\nthey like\nme",
      shar = "Shared\ninterests",
      sinc = "Sincere"))

ggplot(data = chart_data %>% dplyr::filter(rowname != "(Intercept)"),
  aes(x=Factor, y=Effect, fill=Gender)) +
  geom_bar(stat="identity", color='black', position=position_dodge()) +
  theme_minimal() +
  labs(x="", title = "What matters in speed dating?") +
  scale_fill_manual(values=c("#FFC0CB", "#0000FF"))
```

Source: McNulty (2020)

## REAL WORLD APPLICATION



### "Will Smith Tries Online Dating"

## CURRENT STUDENT PROJECTS

- Avocado sales in different US cities
- Rental prices in German cities
- What is my Netflix consumption like?
- Popularity of different US bills
- Lifestyle habits and weight issues
- Musical consumption during a pandemic
- Influence of Elon Musk tweets on bitcoin
- Influence of Queens Gambit on Chess.com

# WHICH LANGUAGE?

Data visualization	Numerical computation	Application building	Machine learning
R	FORTRAN	Java	R
Snap!	GNU Octave	C++	Python
Processing	C	Kotlin	Julia



## R VS. PYTHON



- Product
- Popularity
- Prediction
- Proposal

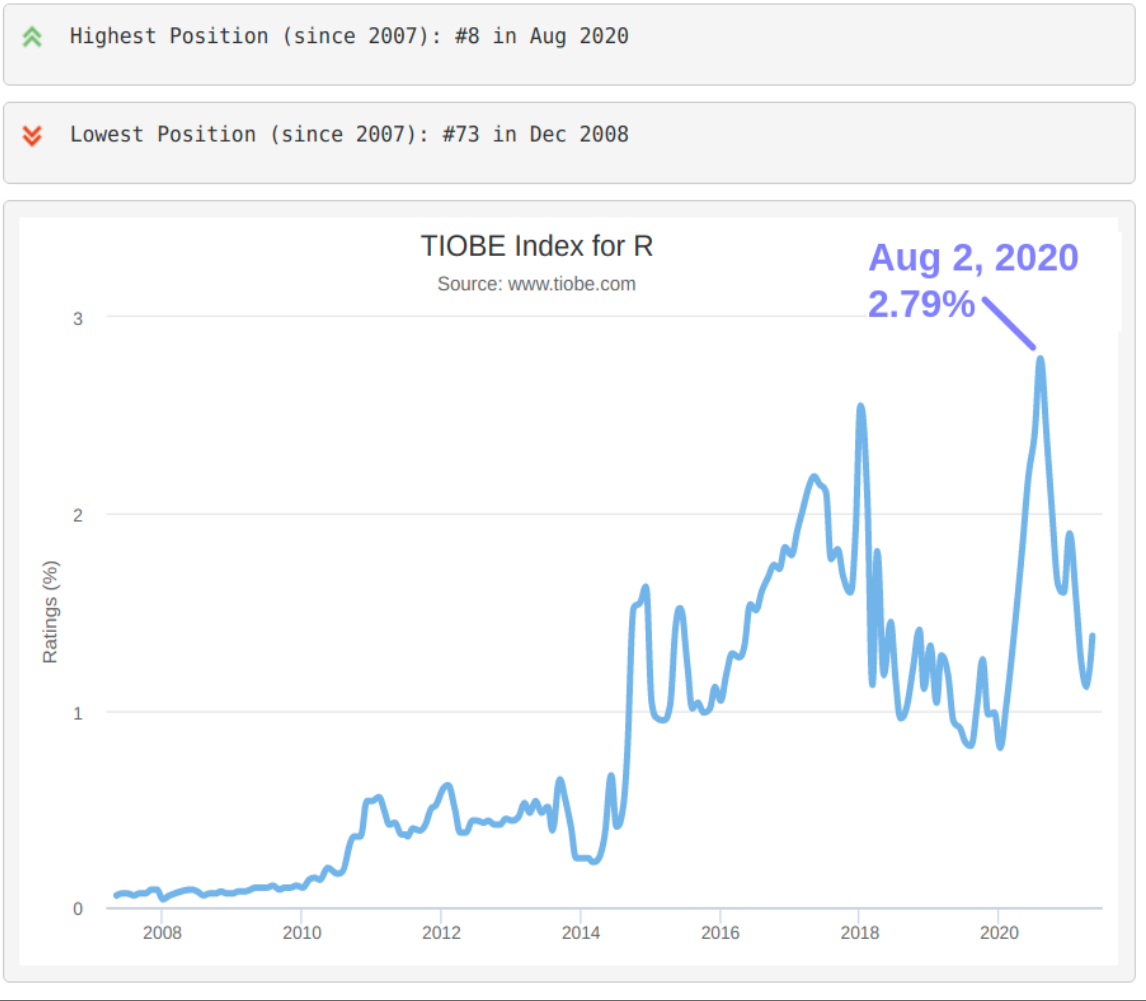
## PRODUCT - DIFFERENCES

	R [4.1.0]	Python [3.9.5]
Release:	1991 (1976)	2008 (1989)
Written in:	FORTRAN, C, R	C
Purpose:	Math & Statistics	Productivity
License:	GPL	<b><u>PSF</u></b>
Libraries:	17,634 ( <b><u>CRAN</u></b> )	137,000 ( <b><u>PyPI</u></b> )
Vectors:	Start at 1	Start at 0
Typedness:	Weak	Strong

## PRODUCT - SIMILARITIES

<b>Access:</b>	Free	Easy to learn
<b>Paradigms:</b>	OOP	Functional
<b>Memory:</b>	Dynamical <sup><u>14</u></sup>	Garbage <sup><u>15</u></sup>
<b>Translation:</b>	Interpreter	REPL/shell
<b>Analytics:</b>	Visualization	Machine Learning
<b>Creativity:</b>	Packages	Portability
<b>Interfaces:</b>	Shiny/Django	SQL/SQLite

POPULARITY - SEARCH



## POPULARITY - NO CONTEST

	R	Python
<b><u>TIOBE Index</u></b>	1.38%	11.74%
<b><u>Loved (Dev)</u></b>	44.5%	66.7%
Dreaded (Dev)	55.5%	33.3%
Wanted (Dev)	5.1%	30.0%
Salary (US)	\$109k	\$120k
Salary (World)	\$59k	\$57k

## PREDICTION - THE WAR IS OVER!



## PROPOSAL - WHAT YOU SHOULD DO

- R
- R + Python
- R + Python + Java
- R + Python + Java + C++
- R + Python + Java + C++ +  $^[\wedge]+\$$

## SUMMARY

---

Data science:	Decision intelligence
Core skills:	Modeling and coding
Applications:	ML / EDA / NLP
R vs. Python:	R + Python

---



## TOOLS



- BPMN: **Signavio** Process Manager
- **R**: Emacs + **ESS** + Org-Mode (+ **Python**)
- Slides: GNU **Emacs** + **Org-Mode** + **reveal.js**
- Computer: Dell 7300 i7 1.9GHz (2019)
- OS: Ubuntu 18.04 LTS Linux (2018)
- Wacom Intuos Art Tablet

**THANK YOU! QUESTIONS? COMMENTS?**



Contact: [birkenkrahe@outlook.com](mailto:birkenkrahe@outlook.com)

## REFERENCES

1. Cox. Translating Statistics to Make Decisions. Apress 2017.
2. Huang, Evans & Chattopadhyay. Deep Learning Without Neural Networks: Fractal-nets for Rare Event Modeling. Preprint. **10.21203/rs.3.rs-86045/v1**
3. Landin. The next 700 programming languages. In: Communications of the ACM 9(3) (1966). **10.1145/365230.365257**
4. McNulty. What Matters in Speed Dating? Online: **towardsdatascience.com** (02/14/2020)
5. Pearl & Mackenzie. The Book of Why: The New Science of Cause and Effect. **New York: Basic Books (2018).**
6. Page. The Model Thinker. Basic Books (2018).
7. Polya. How to solve it. Doubleday (1957).
8. Programming vs Coding - A Short Comparison Between Both. **Online: GeeksforGeeks** (09/08/2020)
9. R: A language and environment fo statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL **https://www.r-project.org/**.
10. Silge & Robinson. Text Mining with R. O'Reilly (2017). **Online: tidytextmining.com**
11. Wickham & Grolemund. R for Data Science: Visualize, Model, Transform, Tidy, And Import

Data. O'Reilly (2016). Online:[r4ds.had.co.nz](https://r4ds.had.co.nz)  
(2016)]]

12. Wing. The Data Life Cycle. In: Harvard Data Science Review 1(1) (2019).  
[doi:10.1162/99608f92.e26845b4](https://doi.org/10.1162/99608f92.e26845b4)