

Statistical uncertainty

Pärt Prommik, PhD

Ülo Maiväli, PhD

Plan for today

- Why do we need statistical uncertainty?
- Why we need R (and tidyverse)?
- Practice

id y x
a 0 2
b 1 1
c 2 0

Data is the new gold/oil



INDIA NEWS

'Data is the new oil, new gold,' says PM Modi in Houston

Forbes

Nov 15, 2019, 08:15am EST | 58,895 views

Data Is The New Oil -- And That's A Good Thing



PARTNER CONTENT JORIS TOONDERS, YONEGO

DATA IS THE NEW OIL OF THE DIGITAL ECONOMY



Data is the new oil, the New Gold of the Digital Era!



IS BIG DATA THE NEW BLACK GOLD?

The Economist

Regulating the internet giants

The world's most valuable resource is no longer oil, but data

ARAB NEWS



Dr. May Alabdaly

14 March 2021

EUROPEAN
BUSINESS REVIEW

Data is the new gold. This is how it can benefit everyone – while harming no one

C E O
T O D A Y

Is Data The New Gold?

Forbes
AFRICA

Data Is The New Gold

FUTURES COT

PUBLIC TECHNOLOGY AND INNOVATION NEWS



MATT WATTS, NOVEMBER 17, 2021 | 2 MIN READ

Why data is the new oil



...when it is well-analysed

Harvard Business Review

Data Is Useless Without the Skills to Analyze It

by Jeanne Harris



Published in Towards Data Science

Is Data Really the New Oil in the 21st Century?

Exploring the strengths and limitations of this metaphor in the information age.



Data was, Analytics is, the New Oil

Technology innovation and the decreasing cost of data processing are fuelling the demand for actionable intelligence across sectors.

E

Without Good Analysis, Big Data Is Just a Big Trash Dump

Cisco UK & Ireland Blog > Innovation



Innovation

Data is neither the new oil nor the new gold...



Stu Higgins
October 9, 2018

...and is only valuable when gathered correctly and interpreted accurately

Forbes

Data, Like Oil, Is Pretty Useless In Its Raw State

LinkedIn

Big Data is Useless without Visual Analytics

towards
data science

Data is not the new oil

About the reality of working with data

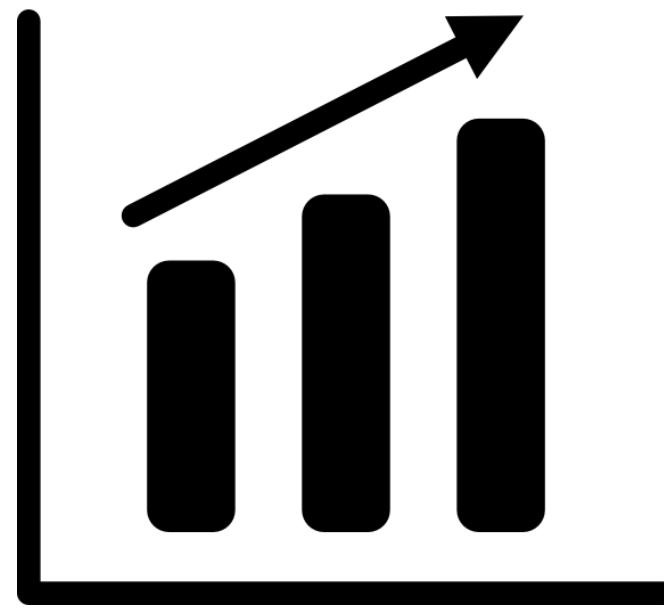
Summary: data is not the new oil

The process of working with data is messy, requires careful planning, engineering, and research, and contains a lot of unknowns and pitfalls.

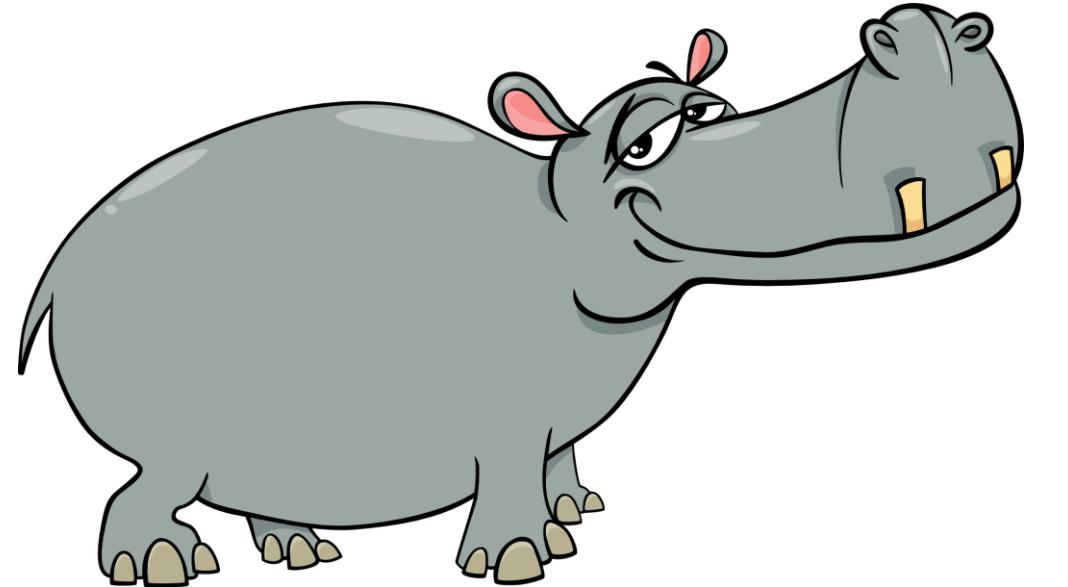
Most importantly, it is not always clear how to leverage the data, since data by itself is too noisy to provide value. Equating data to oil is neglecting this messy and complicated reality.

Data is meaningless without analysis

id	y	x
a	0	2
b	1	1
c	2	0



Who run faster, elephants or hippos?



vs



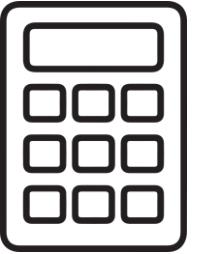
130,000

440,000

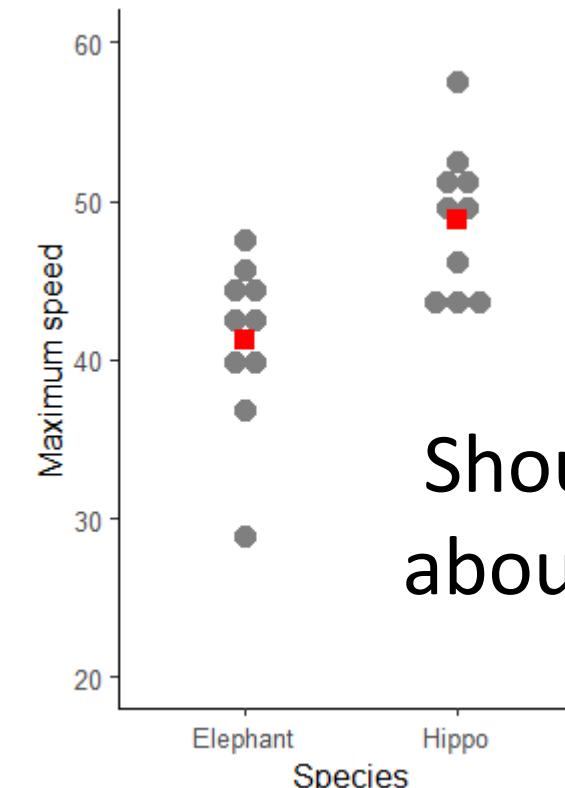
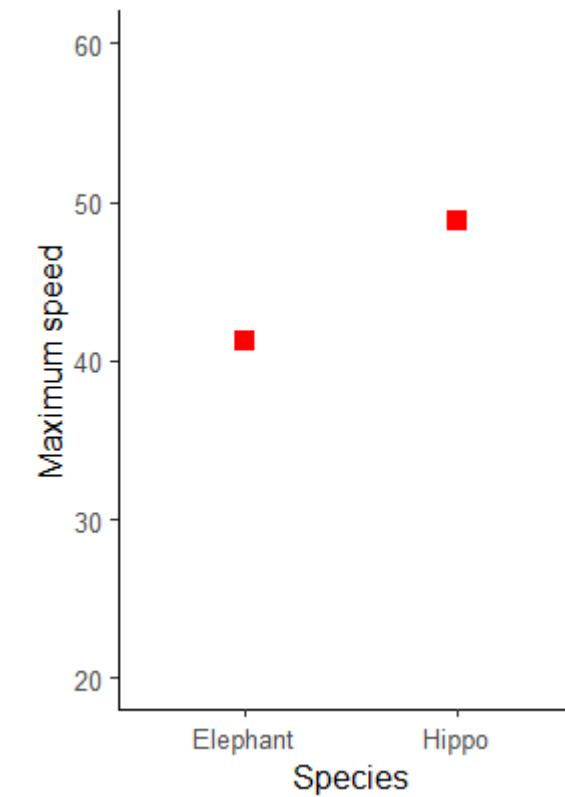
Calculating point estimates

id	y	x
a	0	2
b	1	1
c	2	0

	Species	Max Speed
Hippo	47.8	
Hippo	48.1	
Hippo	33.5	
Hippo	55.0	
Hippo	45.0	
Hippo	43.9	
Hippo	43.3	
Hippo	50.5	
Hippo	53.8	
Hippo	56.2	
Elephant	42.8	
Elephant	42.3	
Elephant	37.3	
Elephant	43.7	
Elephant	40.3	
Elephant	50.8	
Elephant	39.4	
Elephant	35.4	
Elephant	37.5	
Elephant	35.4	



	Species	Mean Max Speed
Hippo	48.8	
Elephant	37.7	



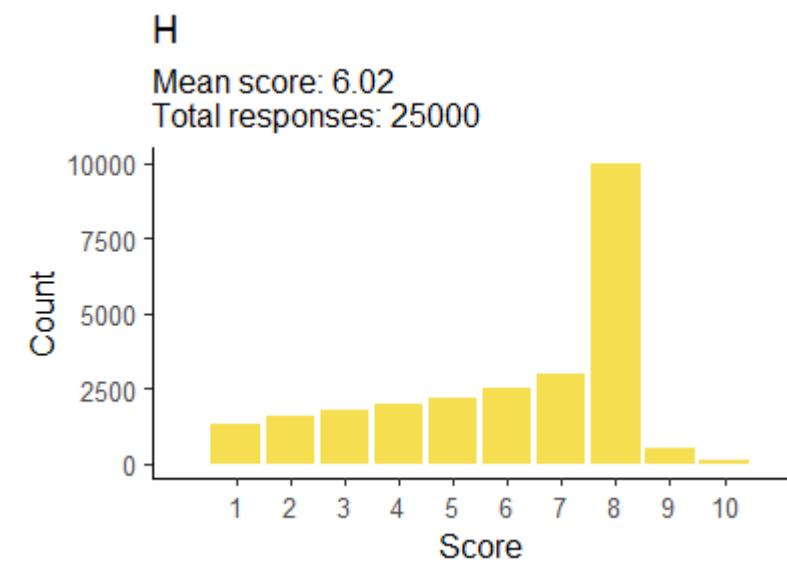
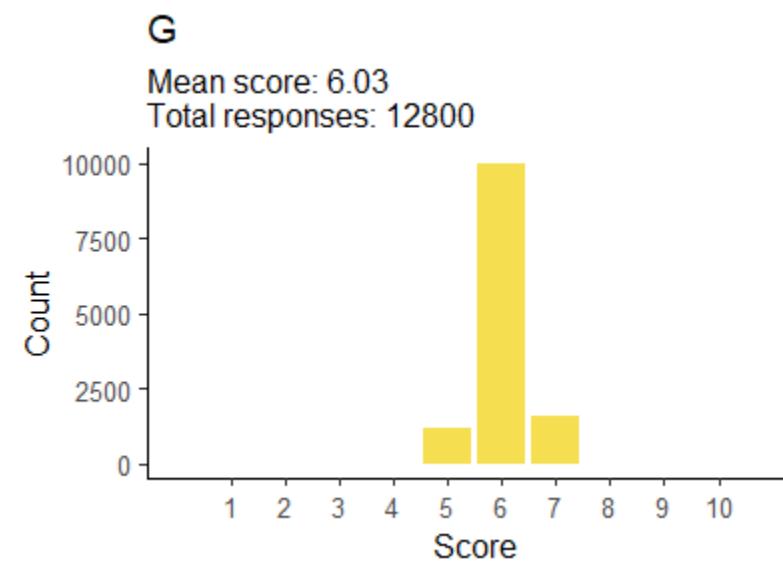
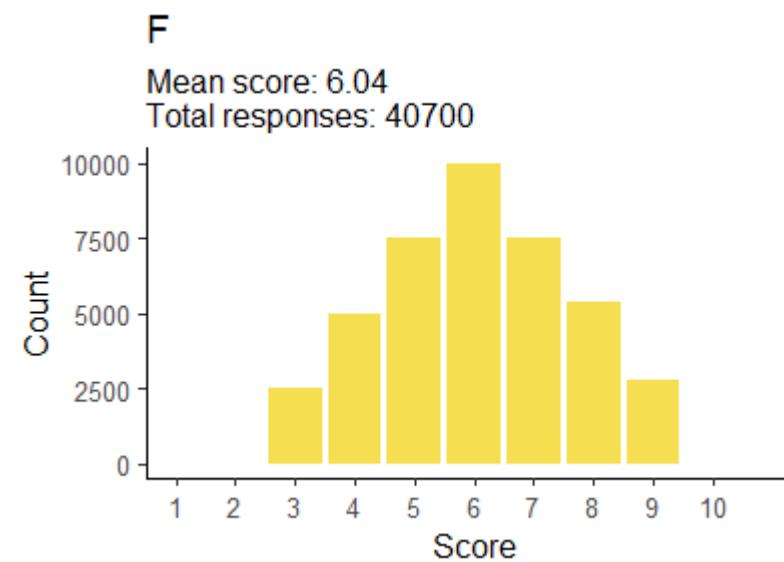
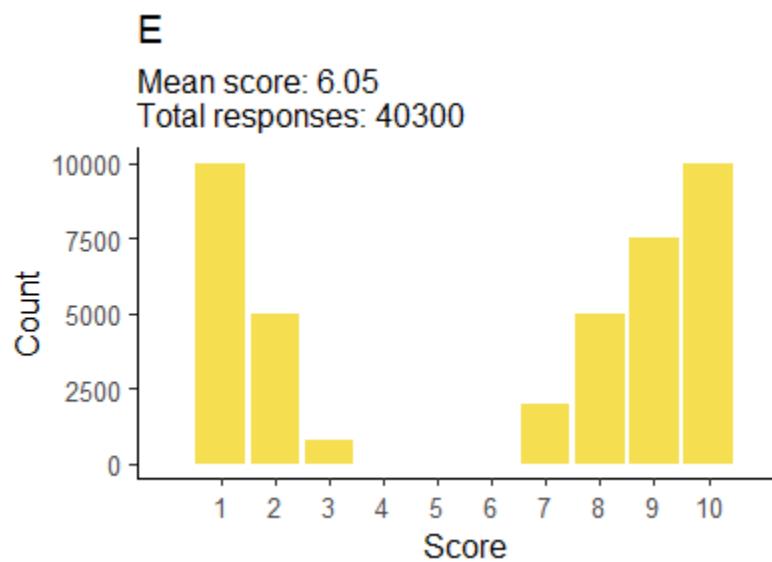
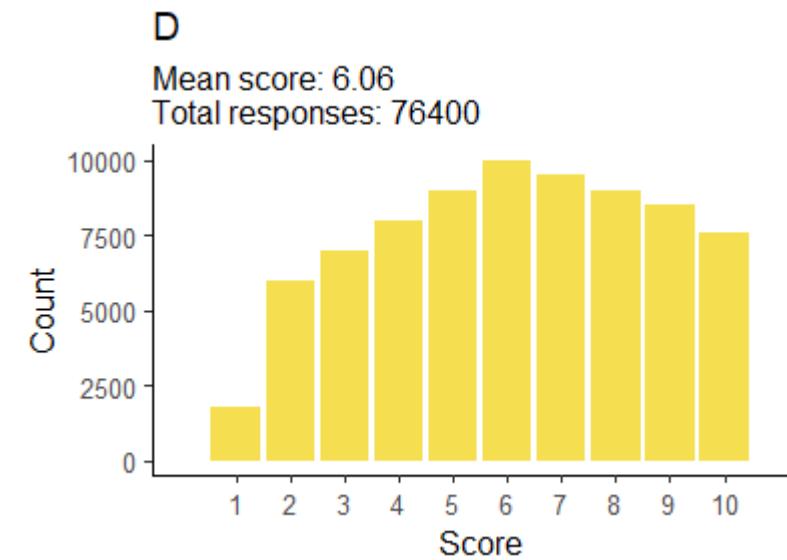
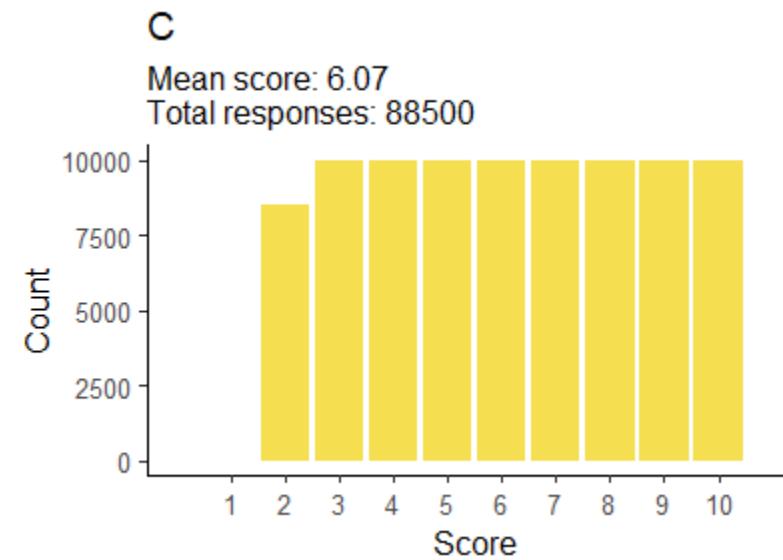
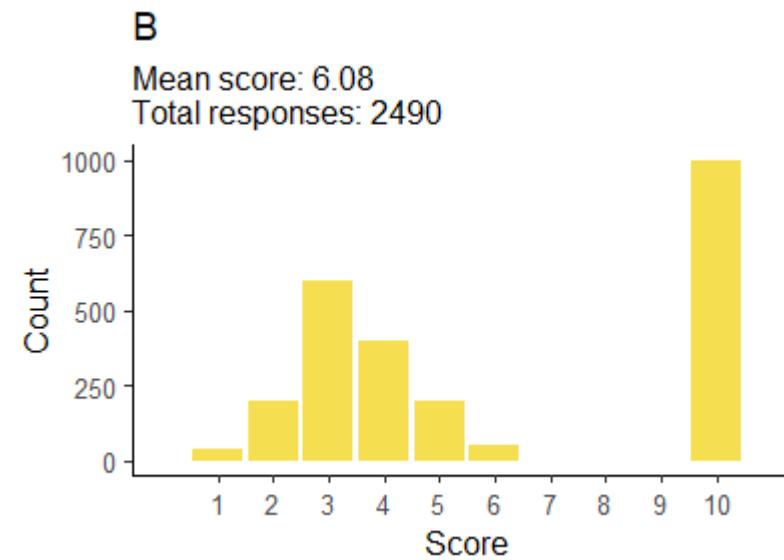
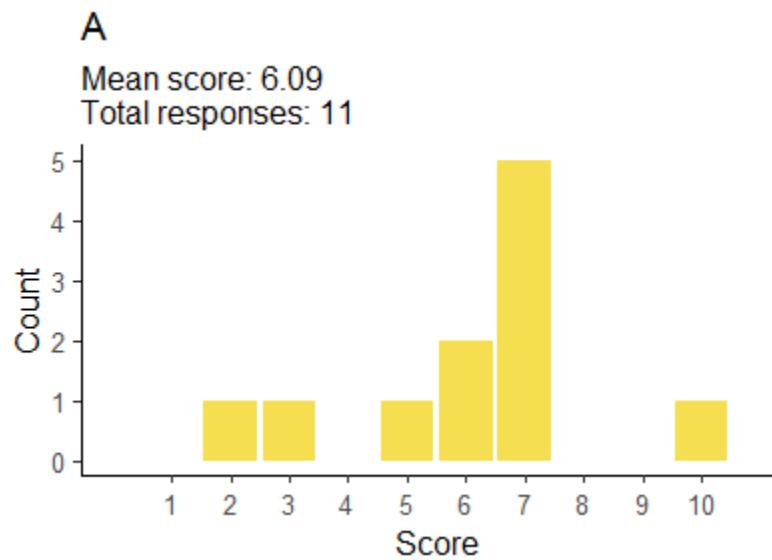
Should we care
about variation?

Choose the best movie for your movienight



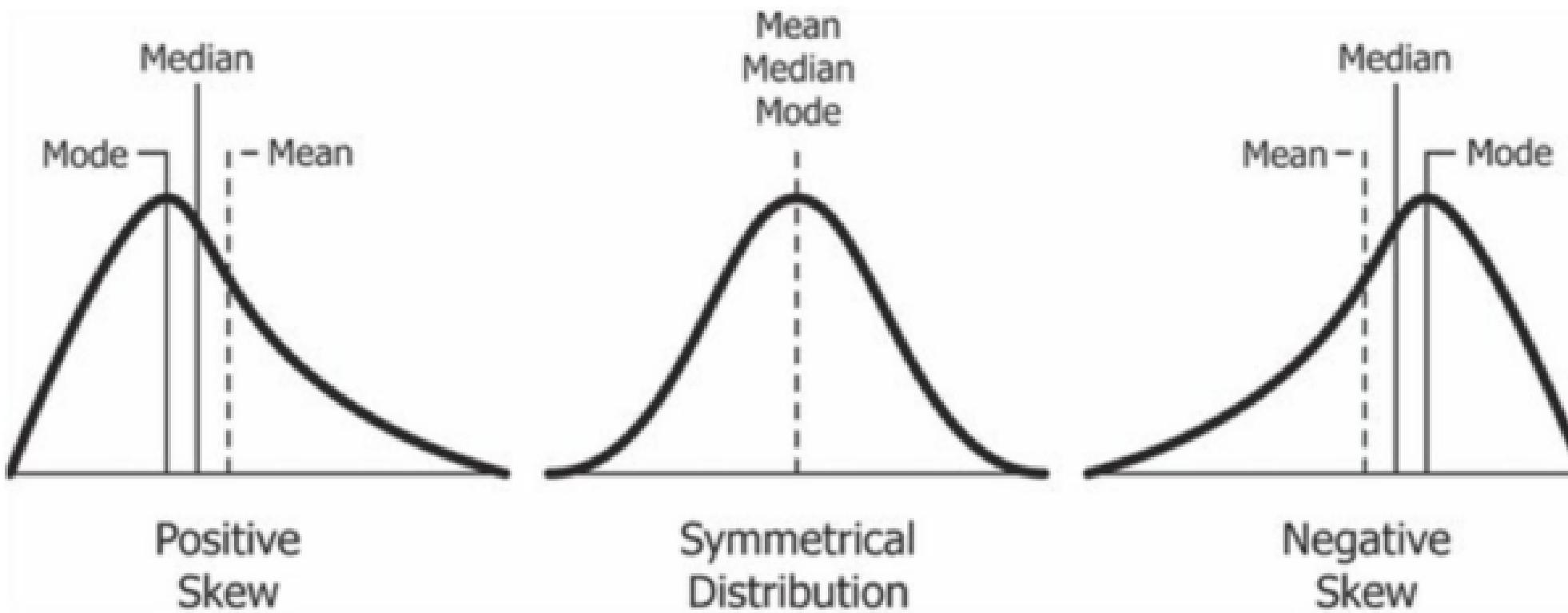
Movie title	Mean IMDB score
A	6.09
B	6.08
C	6.07
D	6.06
E	6.05
F	6.04
G	6.03
H	6.02

IMDB score distributions



Choosing correct measures for describing a distribution

Central tendency



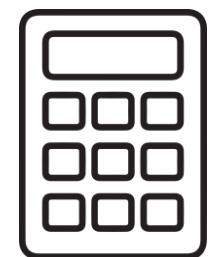
Variation

- Mean
 - Standard deviation
- Median/Mode/Geometric mean
 - Interquartile range
 - Range of percentile
 - Median absolute deviation
 - Geometric standard deviation
 - Minimum and maximum

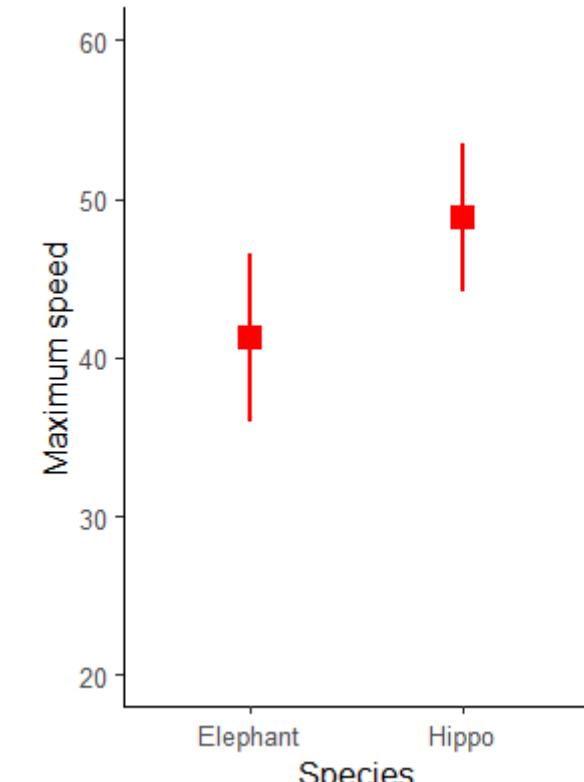
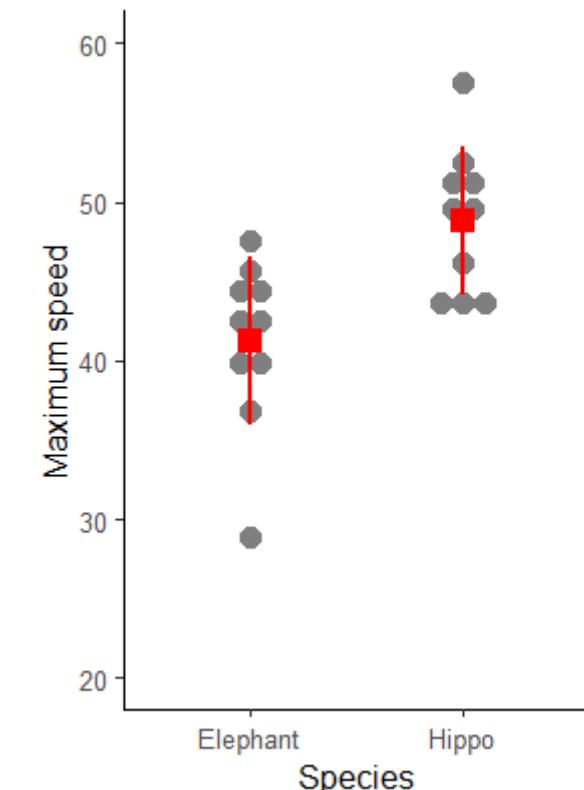
Adding variation is informative

id	y	x
a	0	2
b	1	1
c	2	0

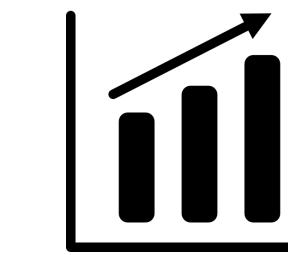
	Species	Max Speed
Hippo	47.8	
Hippo	48.1	
Hippo	33.5	
Hippo	55.0	
Hippo	45.0	
Hippo	43.9	
Hippo	43.3	
Hippo	50.5	
Hippo	53.8	
Hippo	56.2	
Elephant	42.8	
Elephant	42.3	
Elephant	37.3	
Elephant	43.7	
Elephant	40.3	
Elephant	50.8	
Elephant	39.4	
Elephant	35.4	
Elephant	37.5	
Elephant	35.4	



	Species	Mean Max Speed (SD)
Hippo	48.8 (5.3)	
Elephant	37.7 (4.7)	



But what about the difference?

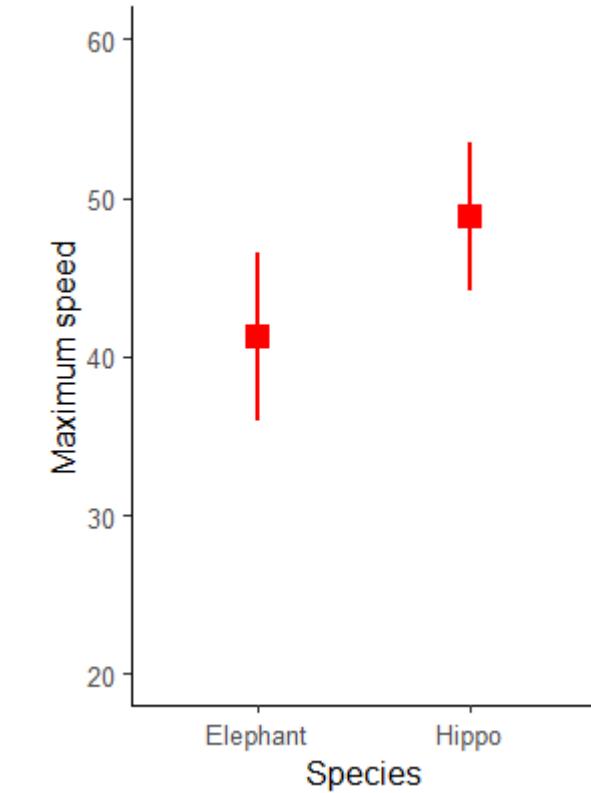
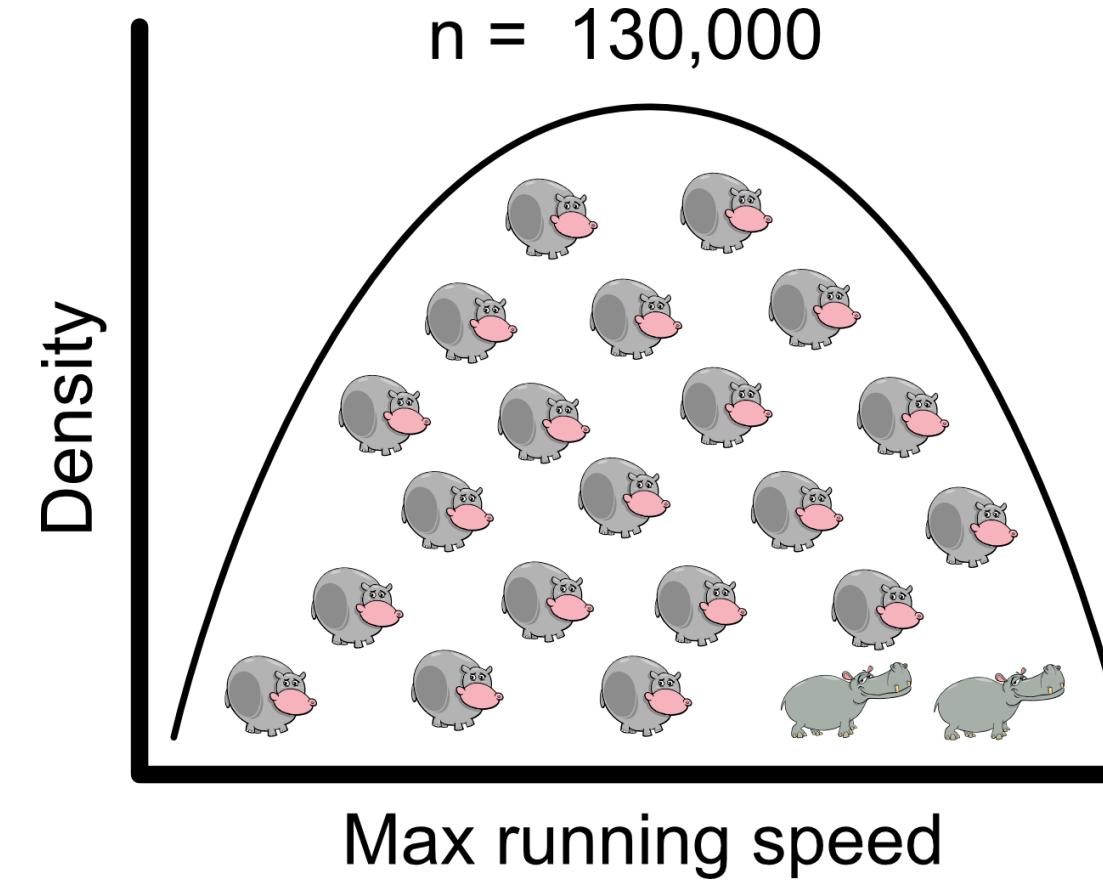


id	y	x	Species	Max Speed
a	0	2	Hippo	47.8
b	1	1	Hippo	48.1
c	2	0	Hippo	33.5
			Hippo	55.0
			Hippo	45.0
			Hippo	43.9
			Hippo	43.3
			Hippo	50.5
			Hippo	53.8
			Hippo	56.2
			Elephant	42.8
			Elephant	42.3
			Elephant	37.3
			Elephant	43.7
			Elephant	40.3
			Elephant	50.8
			Elephant	39.4
			Elephant	35.4
			Elephant	37.5
			Elephant	35.4

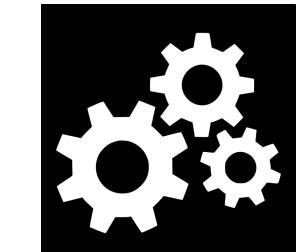
Species	Mean Max Speed (SD)
Hippo	48.8 (5.3)
Elephant	37.7 (4.7)

Hippos run 10.7 k/h faster
(in this sample)

But how certain
we are that
hippos run
faster in these
populations?

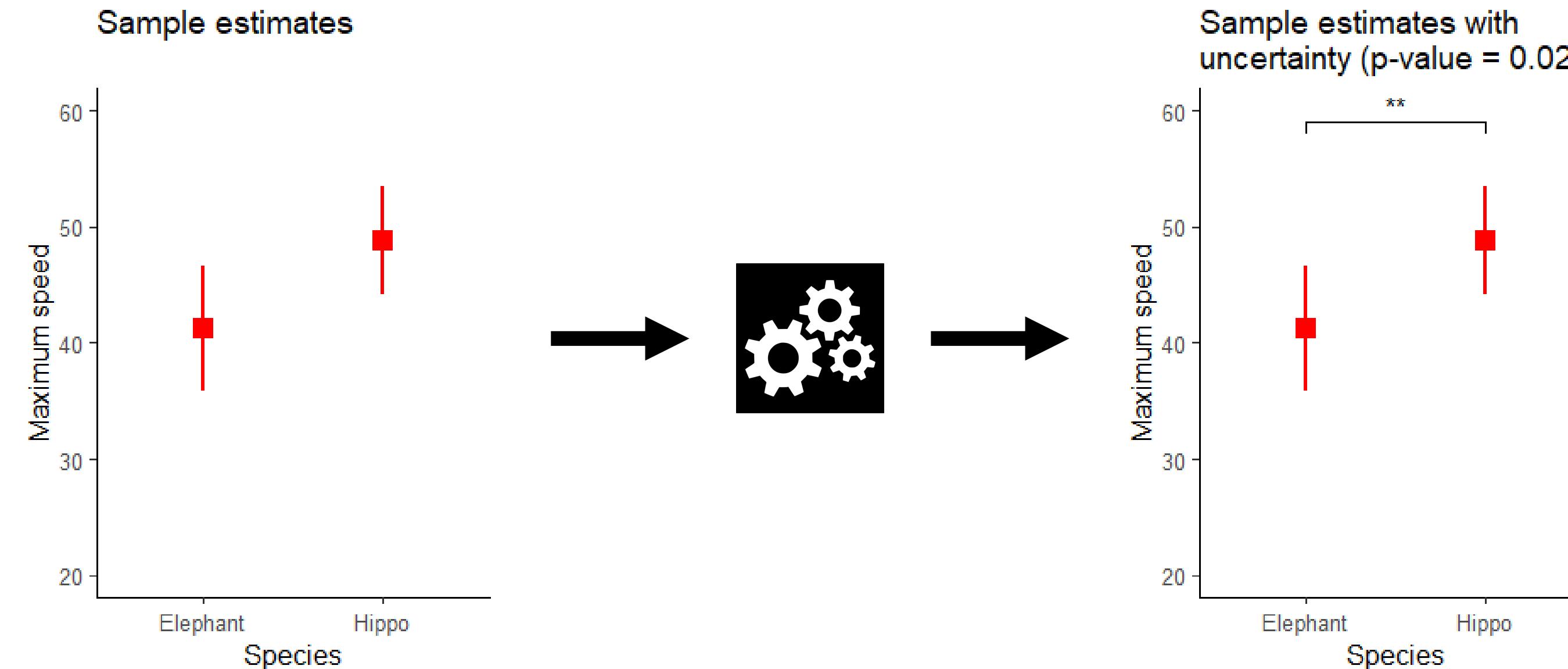


This is why we
need statistics



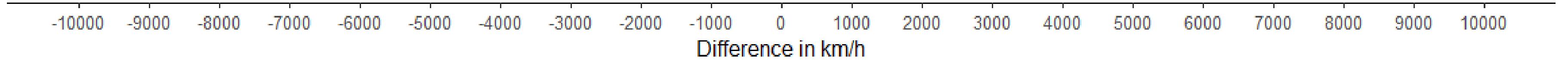
Inferential statistics is about using **sample data**
to make generalizations about a **population**.

Let's add uncertainty (use statistics)



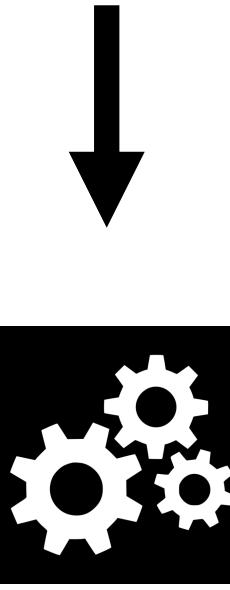
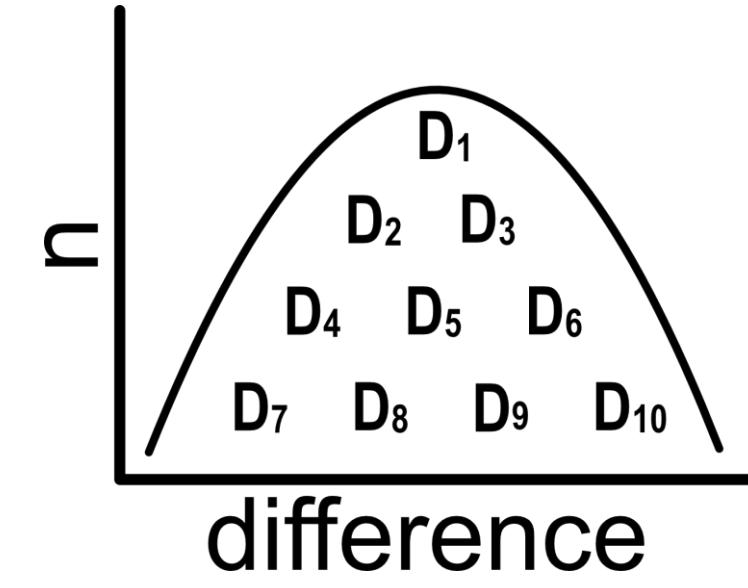
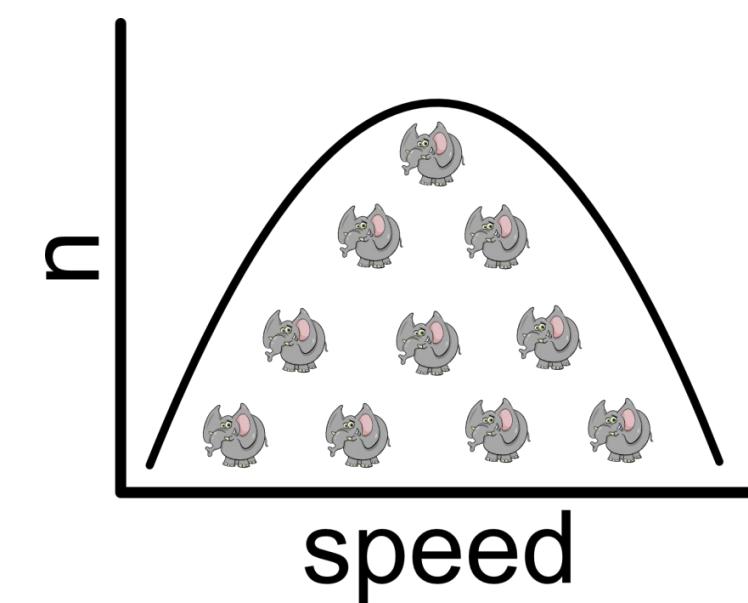
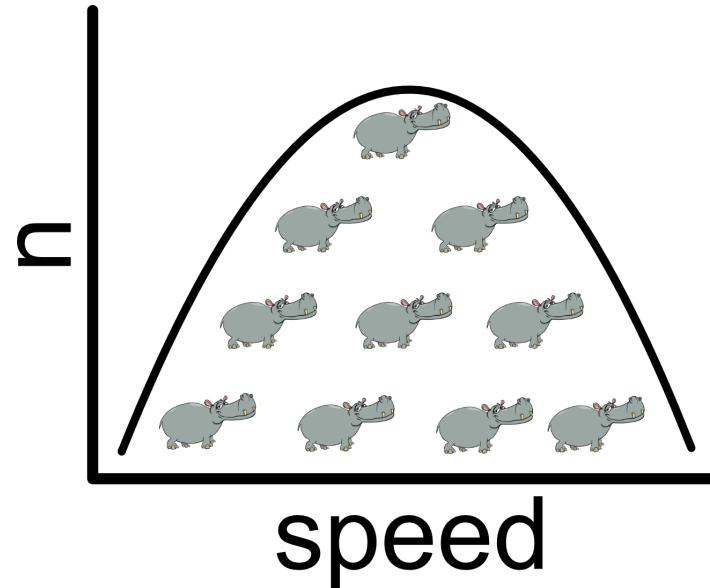
This is hypothesis testing: we have 2% probability that their maximum running speed difference could have just occurred by random chance
(i.e. that the null hypothesis (Hippos = Elephants) is true).

Not very intuitive and informative!

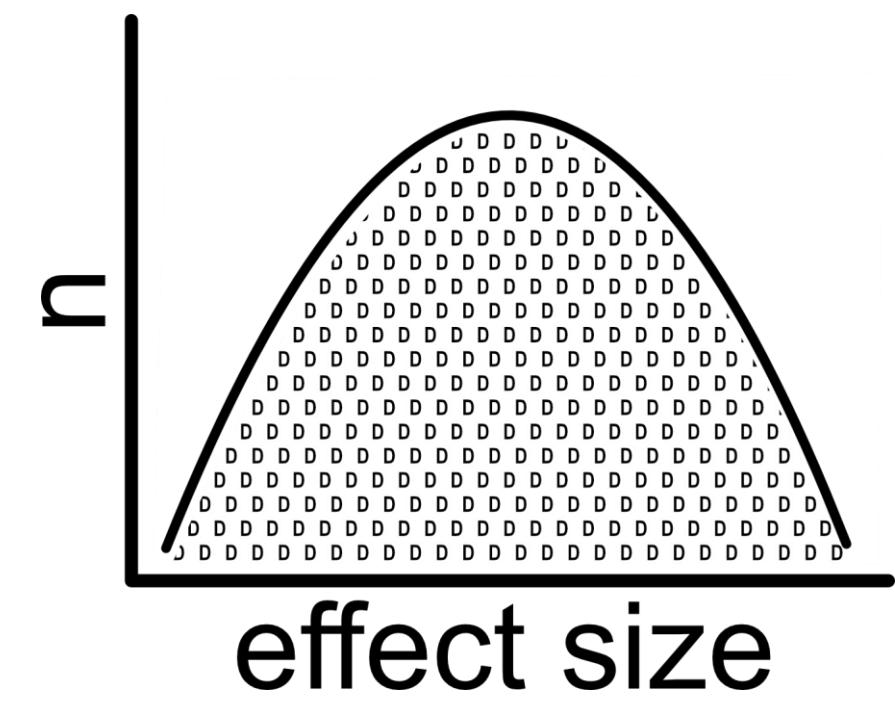


The maximum running speed of group A and B is different ($p = 0.02$).

Sample and population level difference

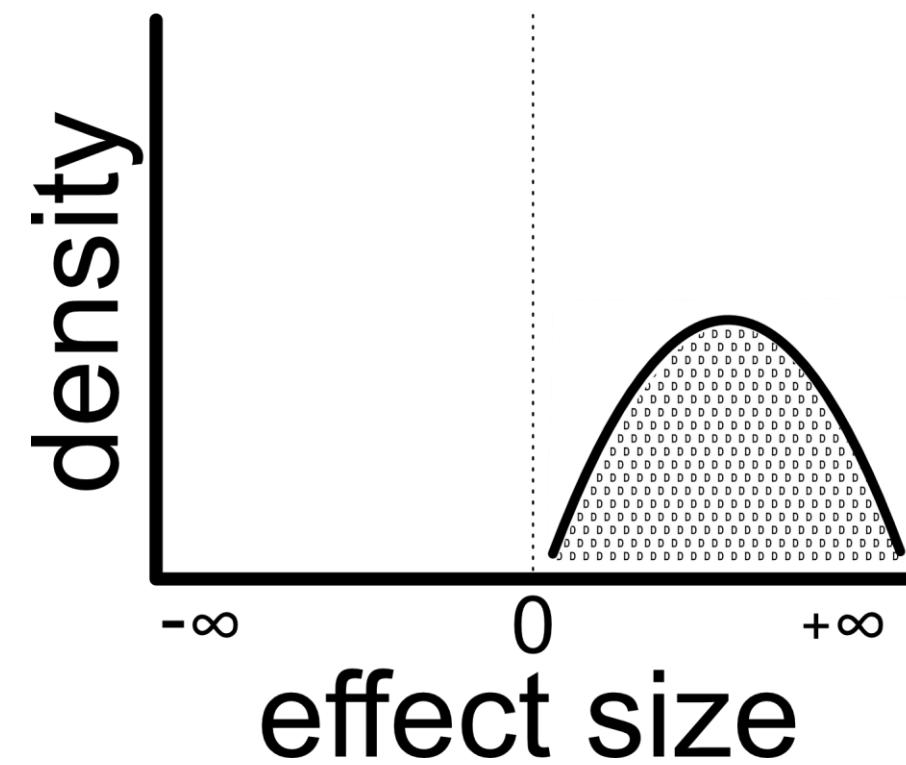


Hippo	Elephant	Difference (D)
44.2	47.6	-3.3
46.2	43.0	3.2
43.1	39.8	3.3
50.0	45.6	4.4
43.0	36.9	6.1
49.1	41.9	7.2
51.5	44.1	7.3
52.4	44.7	7.7
50.9	39.9	11.0
57.6	28.9	28.6



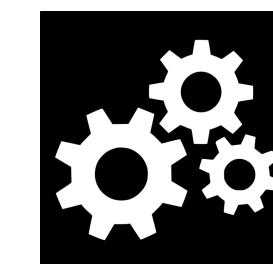
Effect direction

effect size = Hippos - Elephants



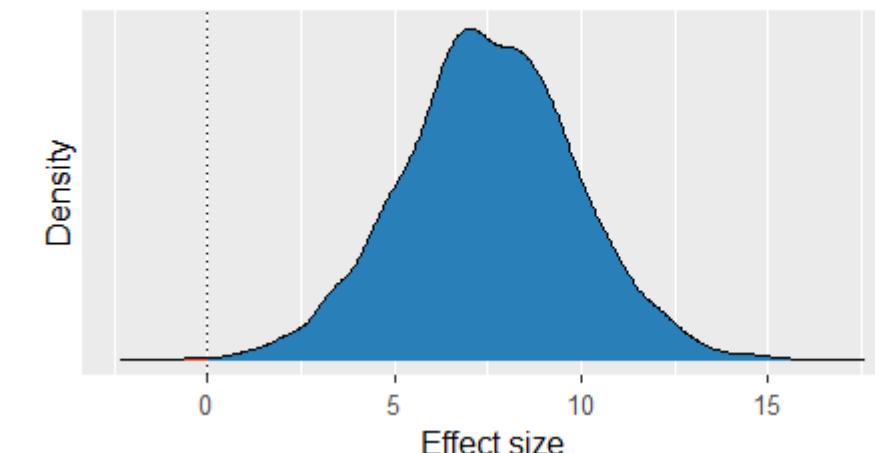
Using effect size for visualising uncertainty

id	y	x
a	0	2
b	1	1
c	2	0



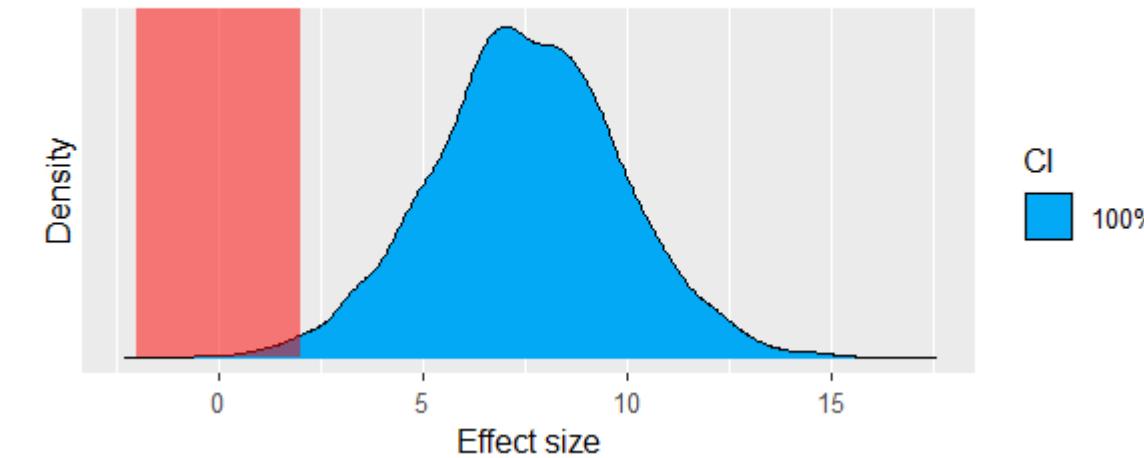
99.93% confidence for hippos being faster

Probability mass in the blue area



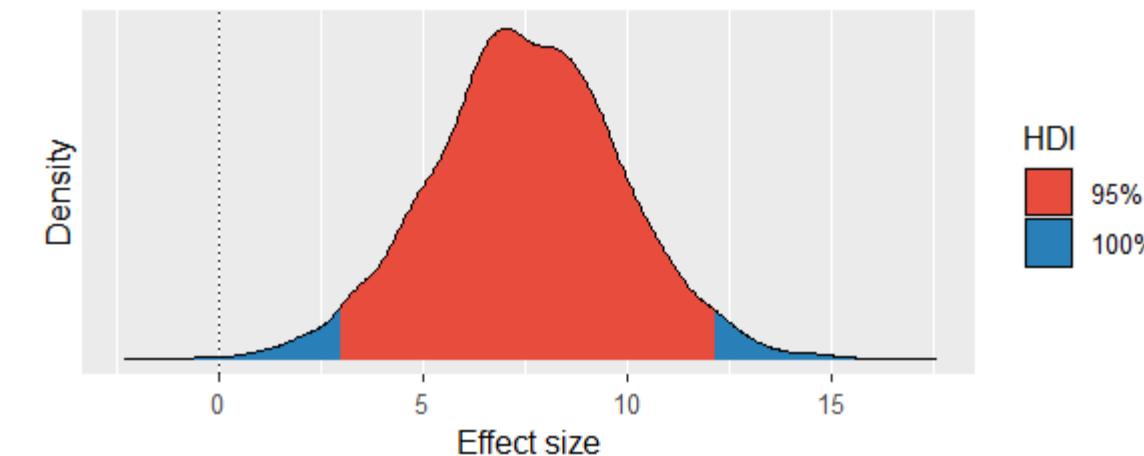
99.1% confidence that the difference is >2 k/h

The probability mass outside of the red rectangle

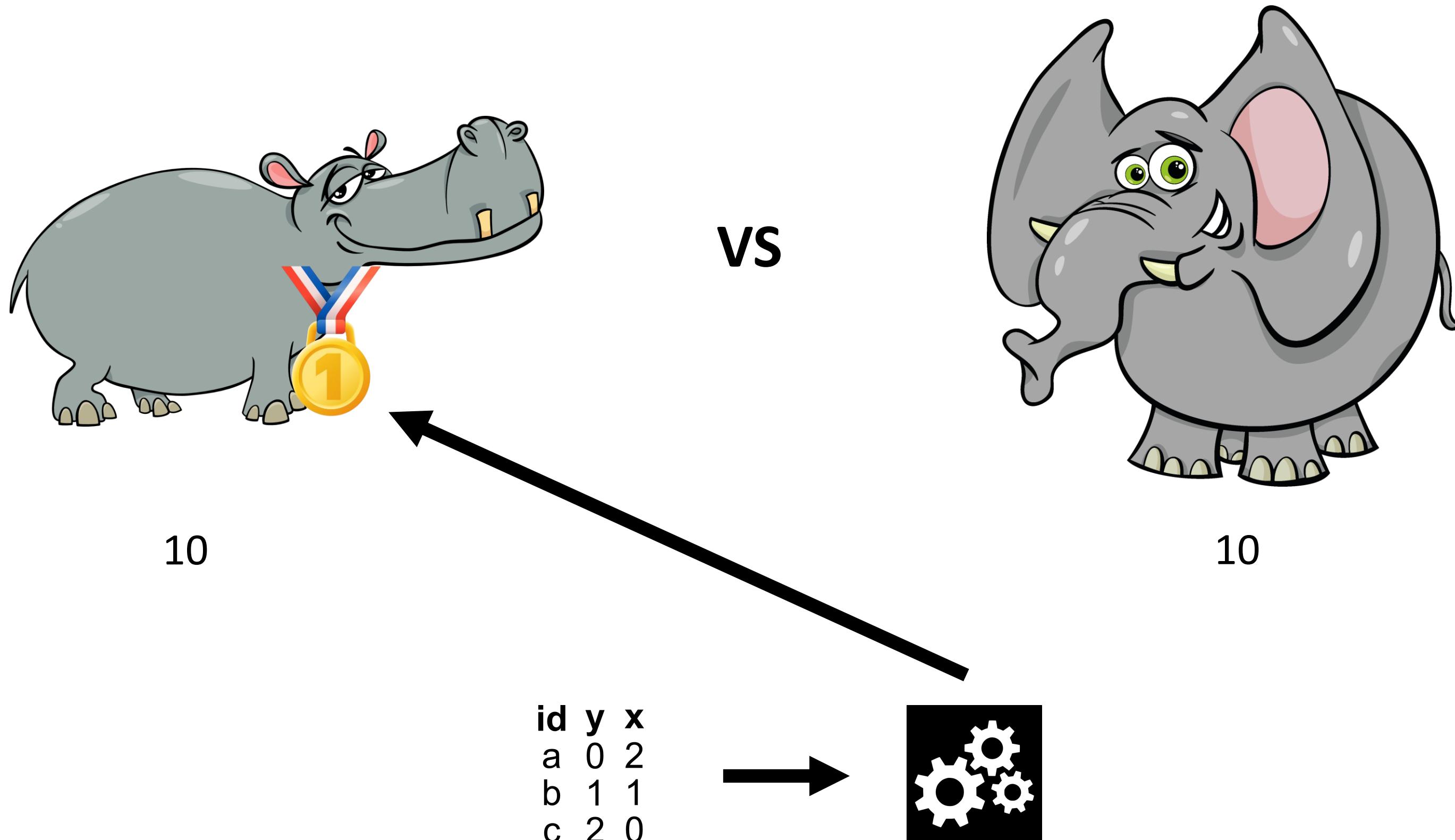


Hippos run 7.2 k/h faster (95% CI: 1.52, 12.40)

Probability mass in the red area

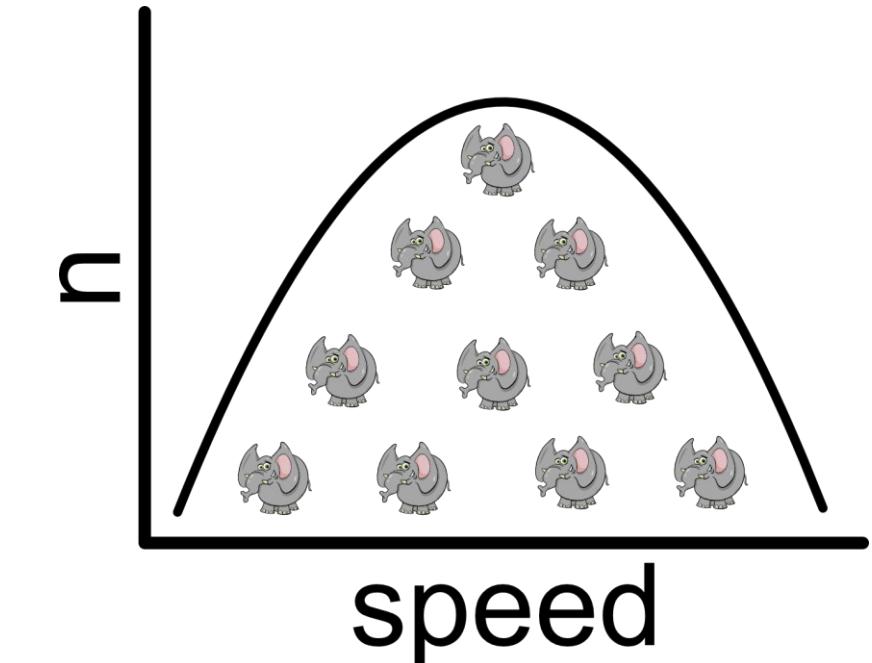
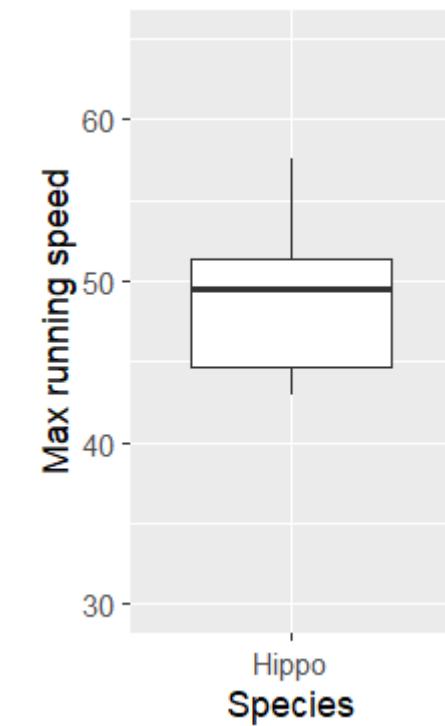
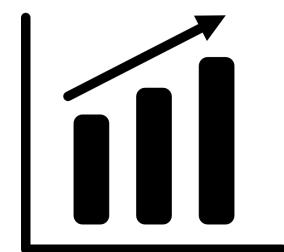


Question > Sample data + Statistics > Population level answer



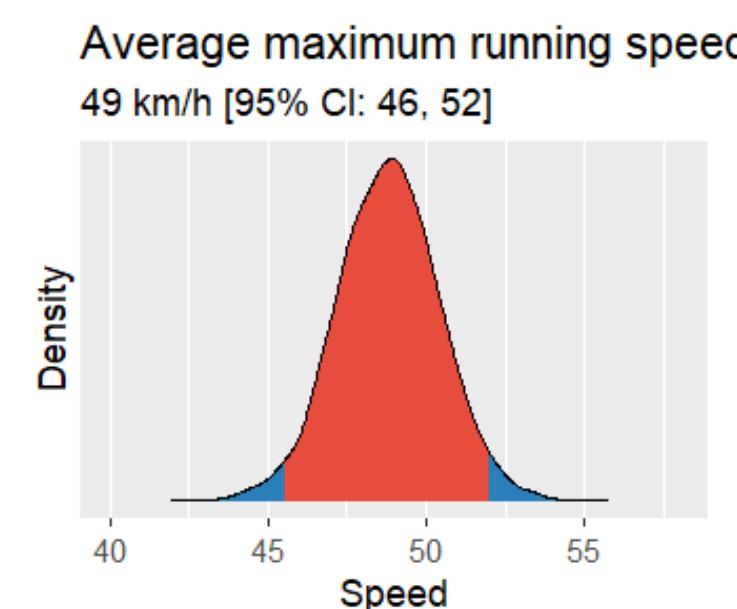
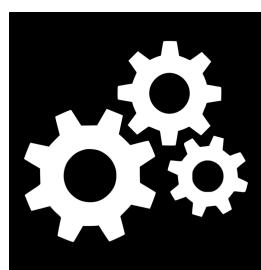
Statistical uncertainty is the uncertainty associated with the use of sample data to make statements about the wider population

id	y	x
a	0	2
b	1	1
c	2	0

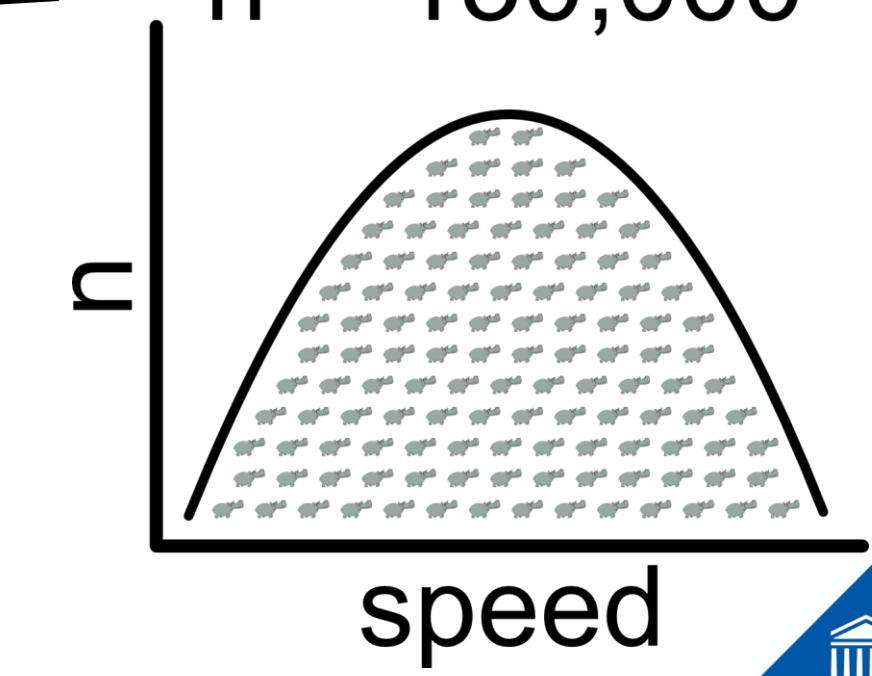


speed

$n = 130,000$

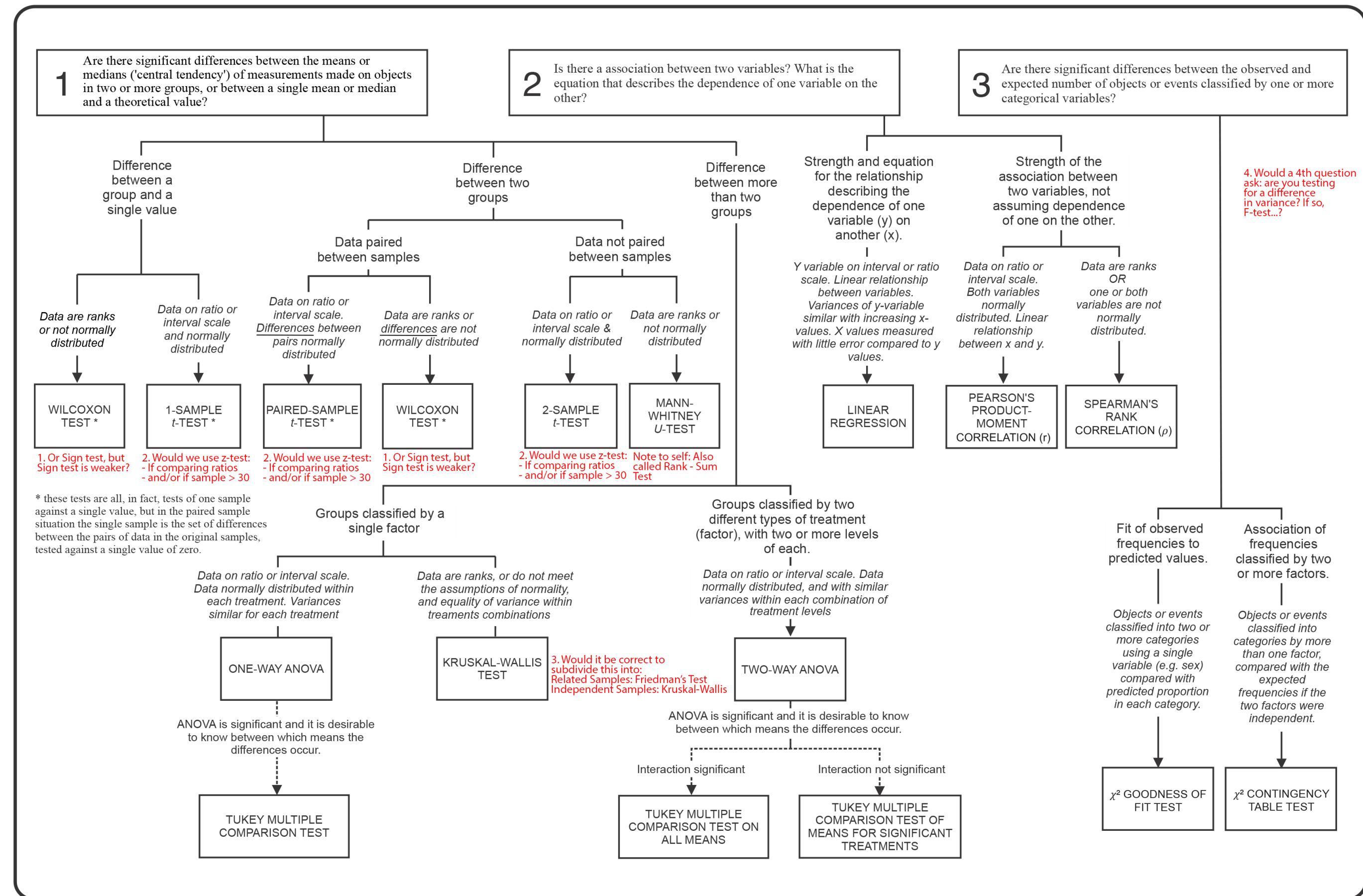


HDI
95%
100%

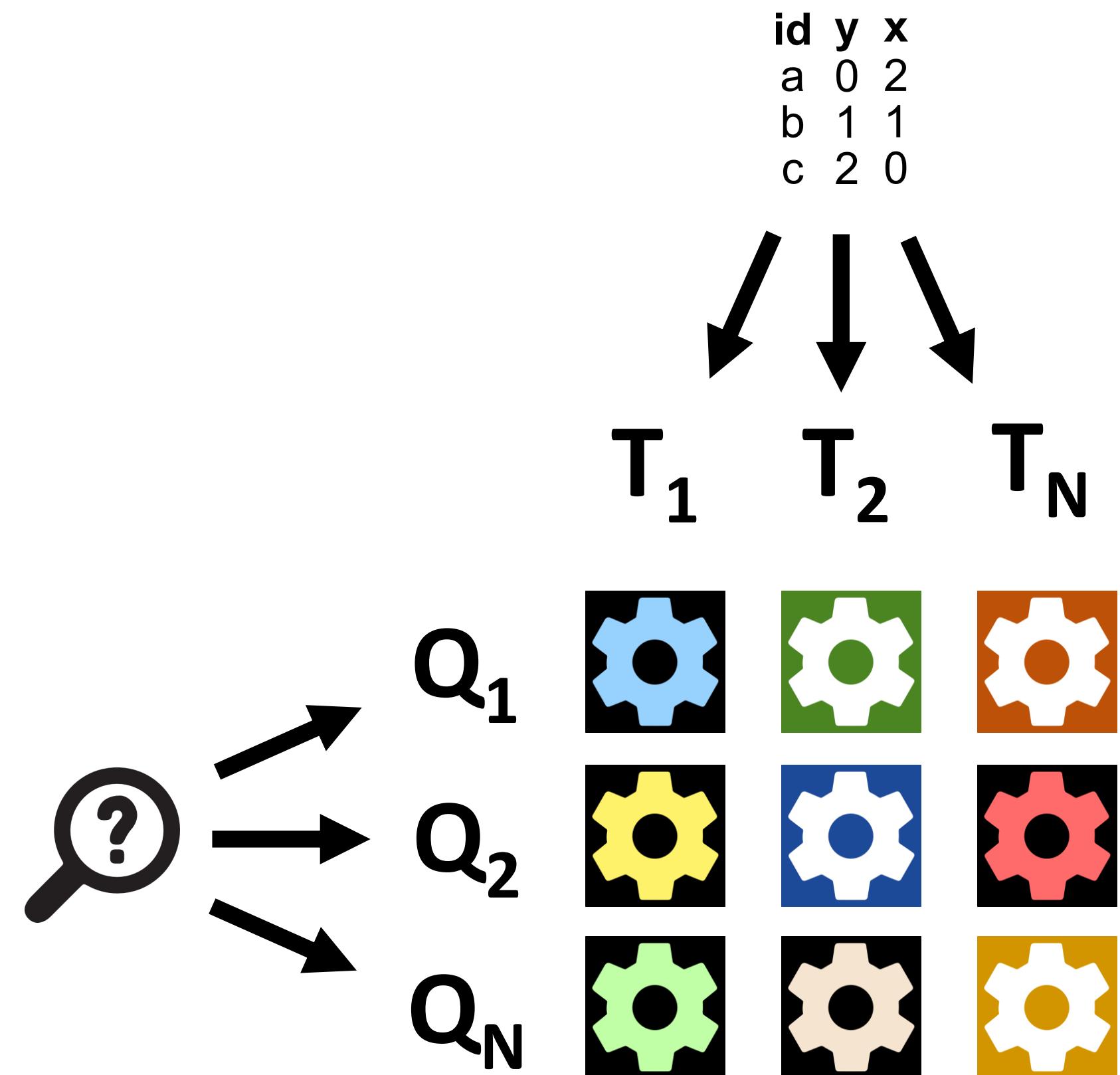


speed

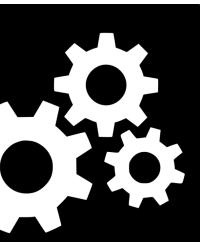
How is the magic (statistics) applied?



How statistical tests are chosen



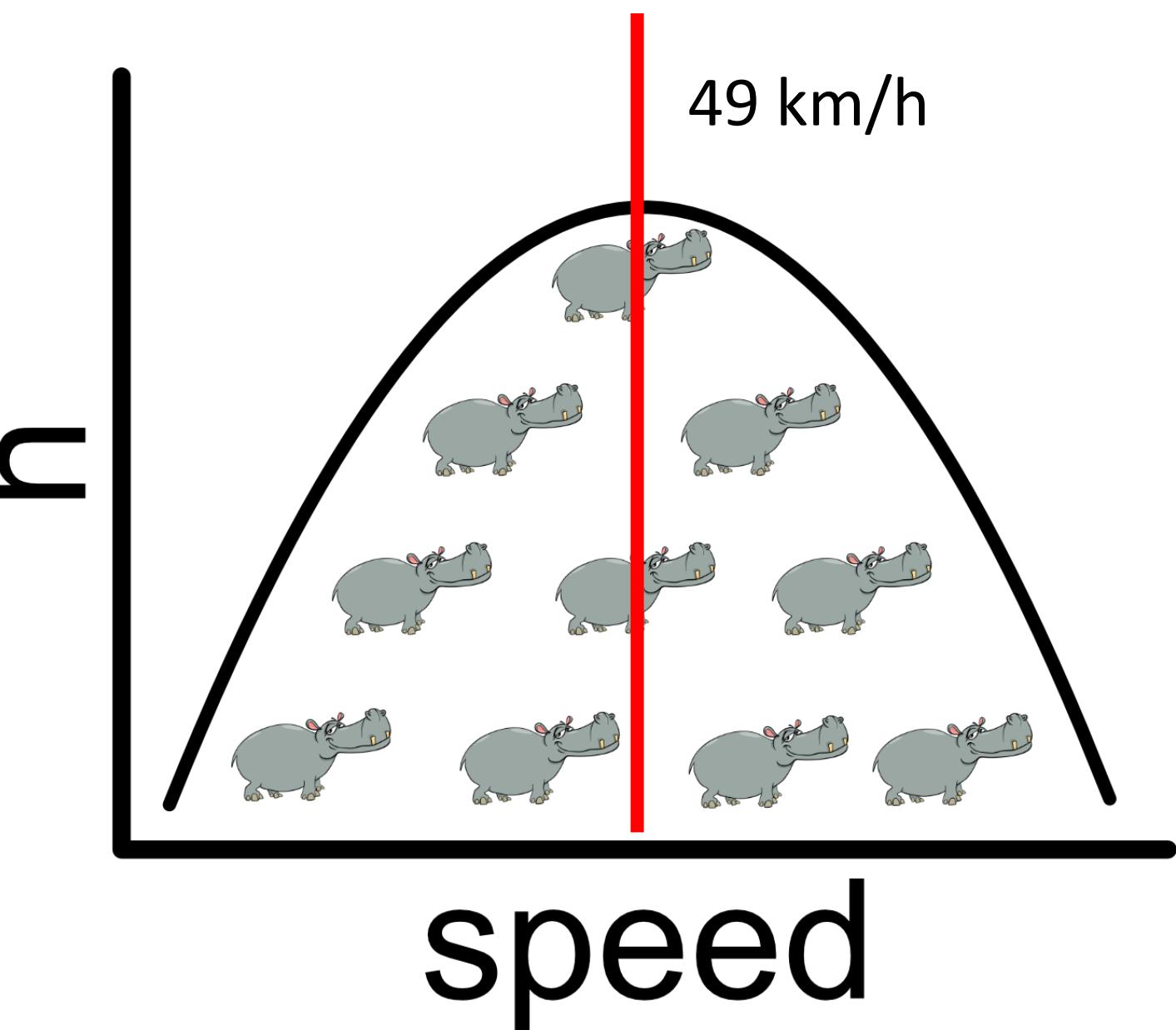
Must know your data (variable types)

		Variable type	Value examples
	x	Integer	1 5 200
	y	Numeric	0.05 1 1.1 “alpha” “gamma” “my text”
id	1	Important since statistical uncertainty is calculated!	
a	0		
b	1		
c	2		
			
		Factor	Female Male Unknown

Different study questions



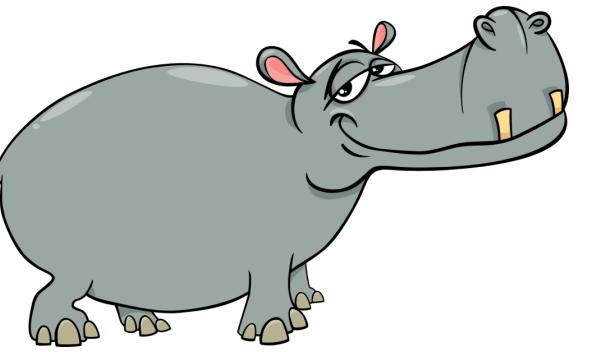
- **Descriptive**
 - What is the maximum running speed of hippos?



Different study questions



- **Descriptive**
 - What is the average running speed of hippos?
- **Comparative**
 - How different is the running speed of hippos from elephants?



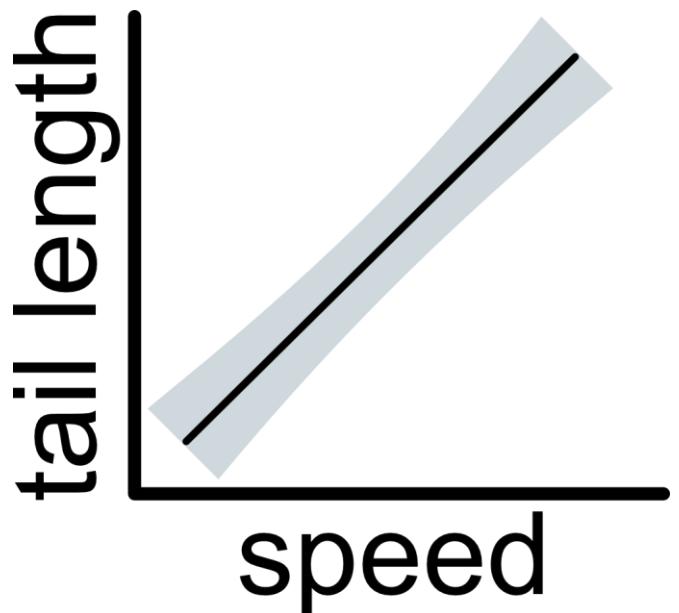
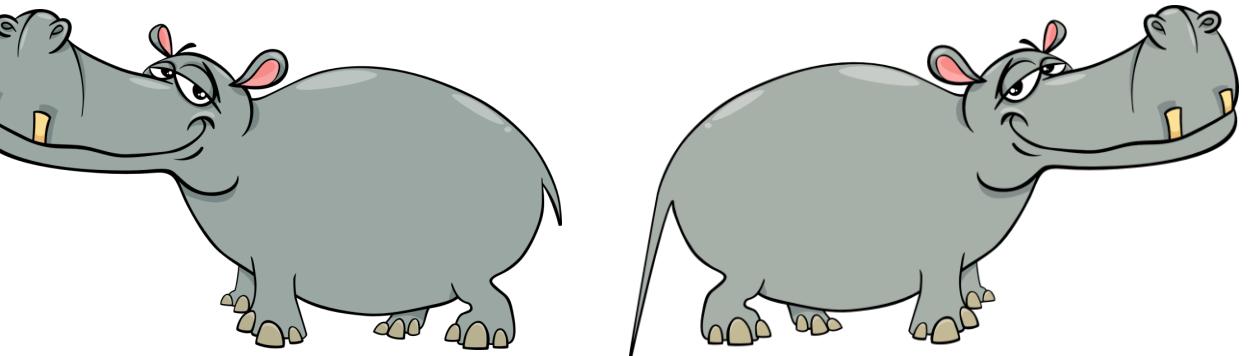
VS



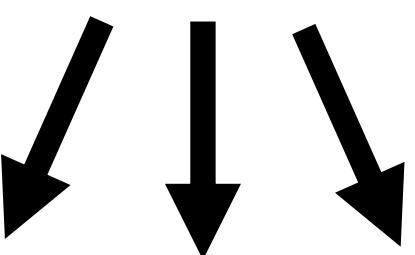
Different study questions



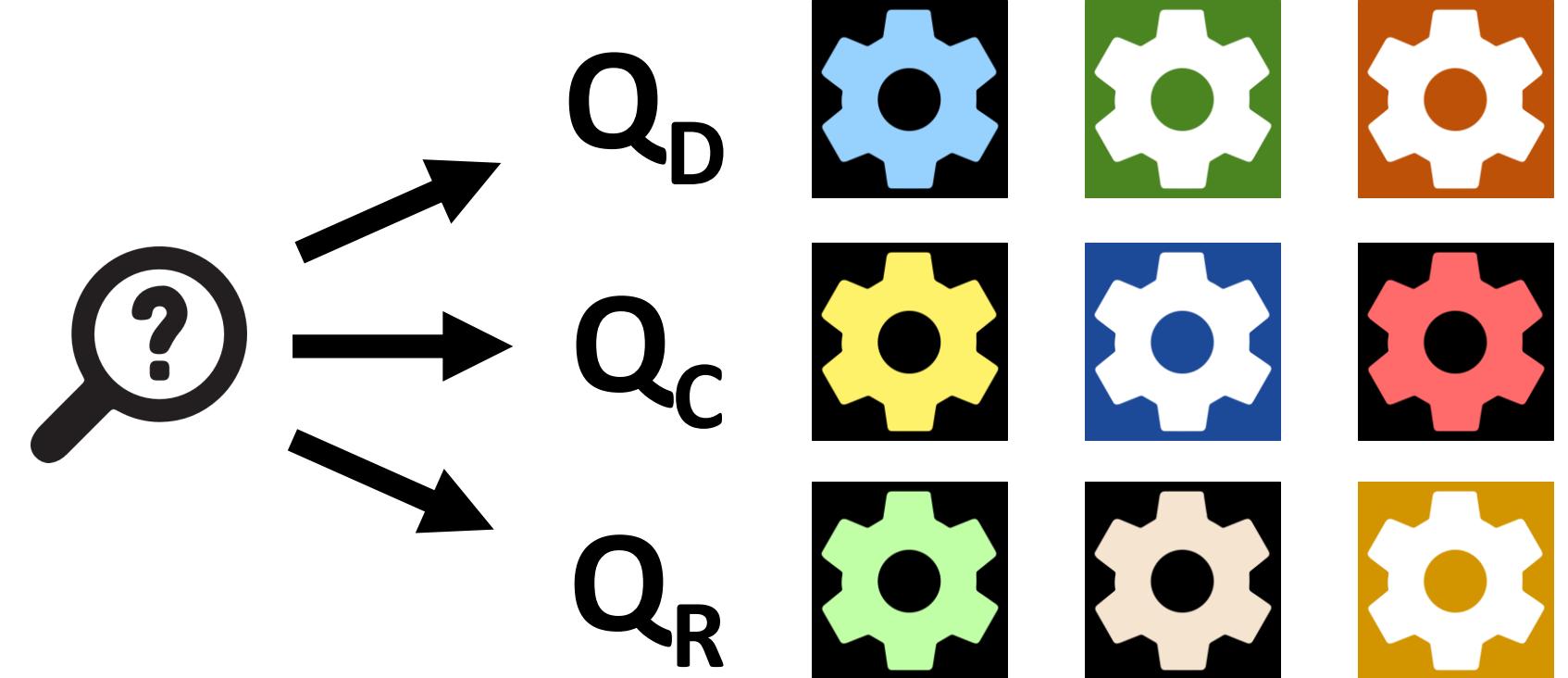
- **Descriptive**
 - What is the average running speed of hippos?
- **Comparative**
 - How different is the running speed of hippos from elephants?
- **Relationship**
 - How does tail length affects running speed of hippos?



id	y	x
a	0	2
b	1	1
c	2	0

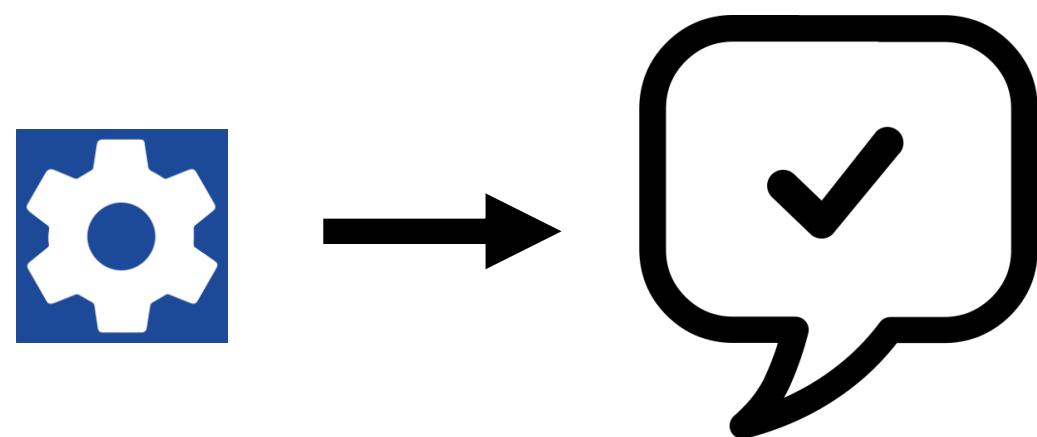
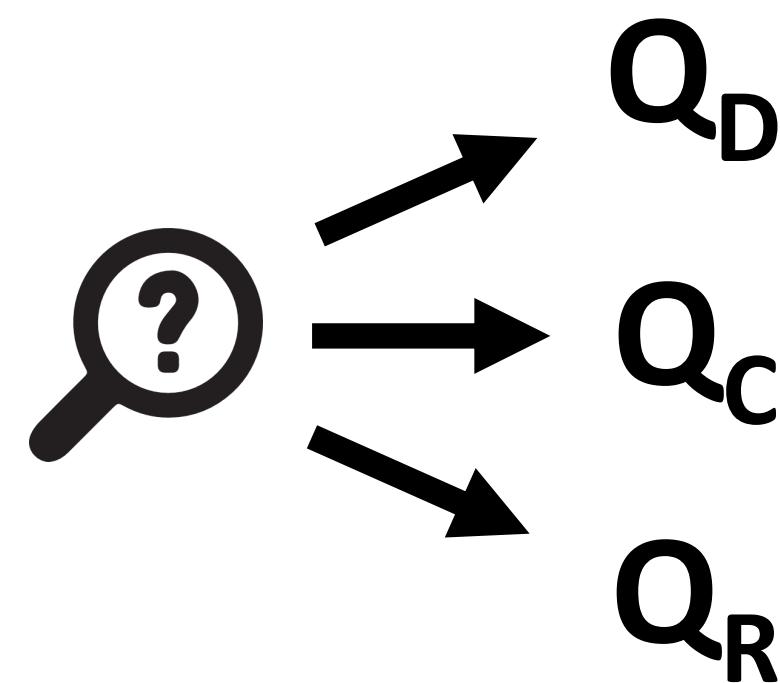


T_{int} T_{num} T_{cat}

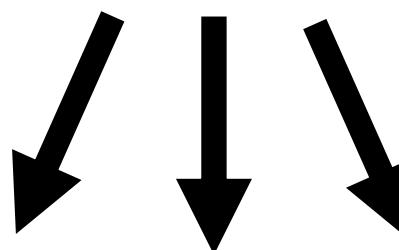


id	y	x
a	0	2
b	1	1
c	2	0

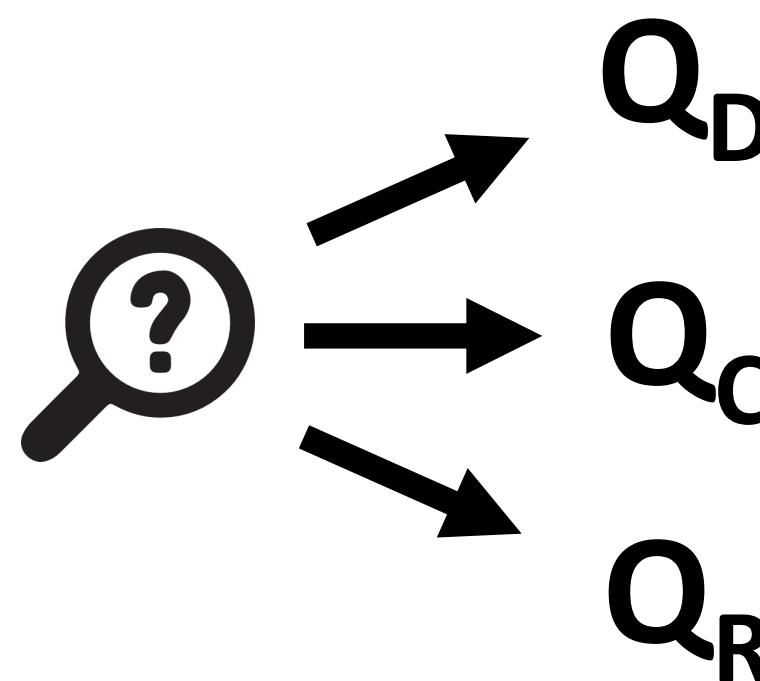
T_{int} T_{num} T_{cat}



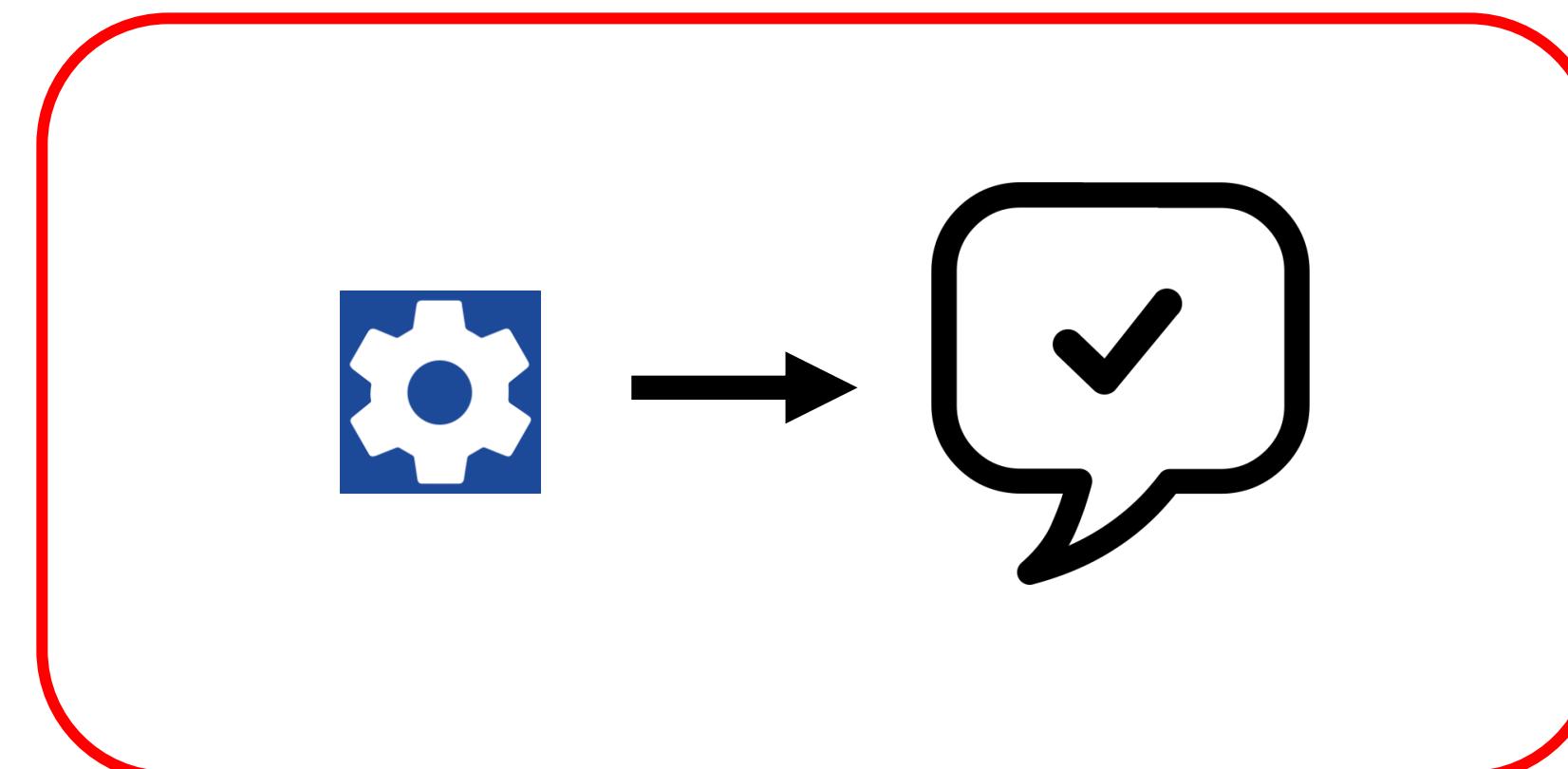
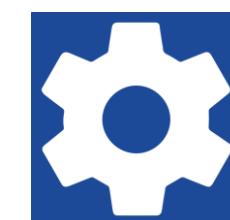
id	y	x
a	0	2
b	1	1
c	2	0



T_{int} T_{num} T_{cat}



Q_D
 Q_C
 Q_R



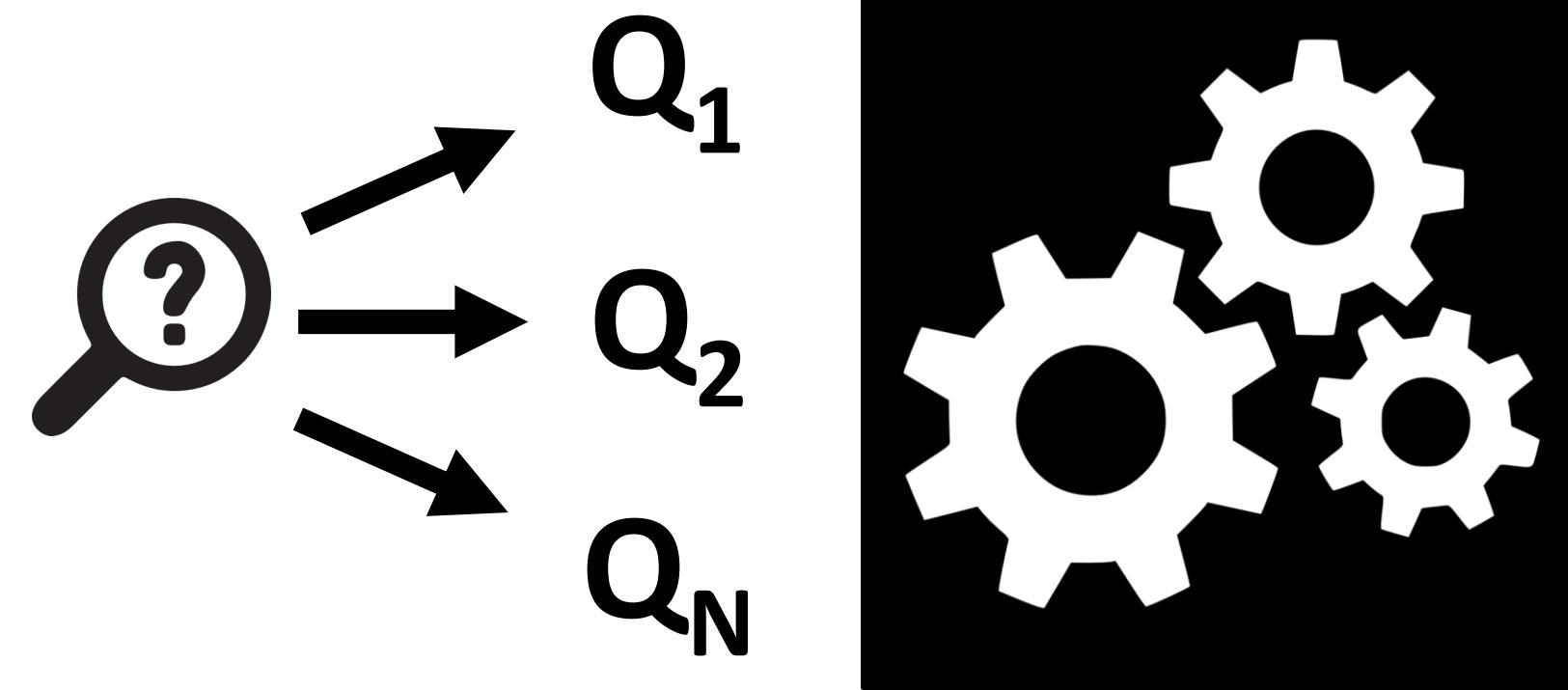
Not a good approach!

A better approach

1. Learn modelling (regression) techniques

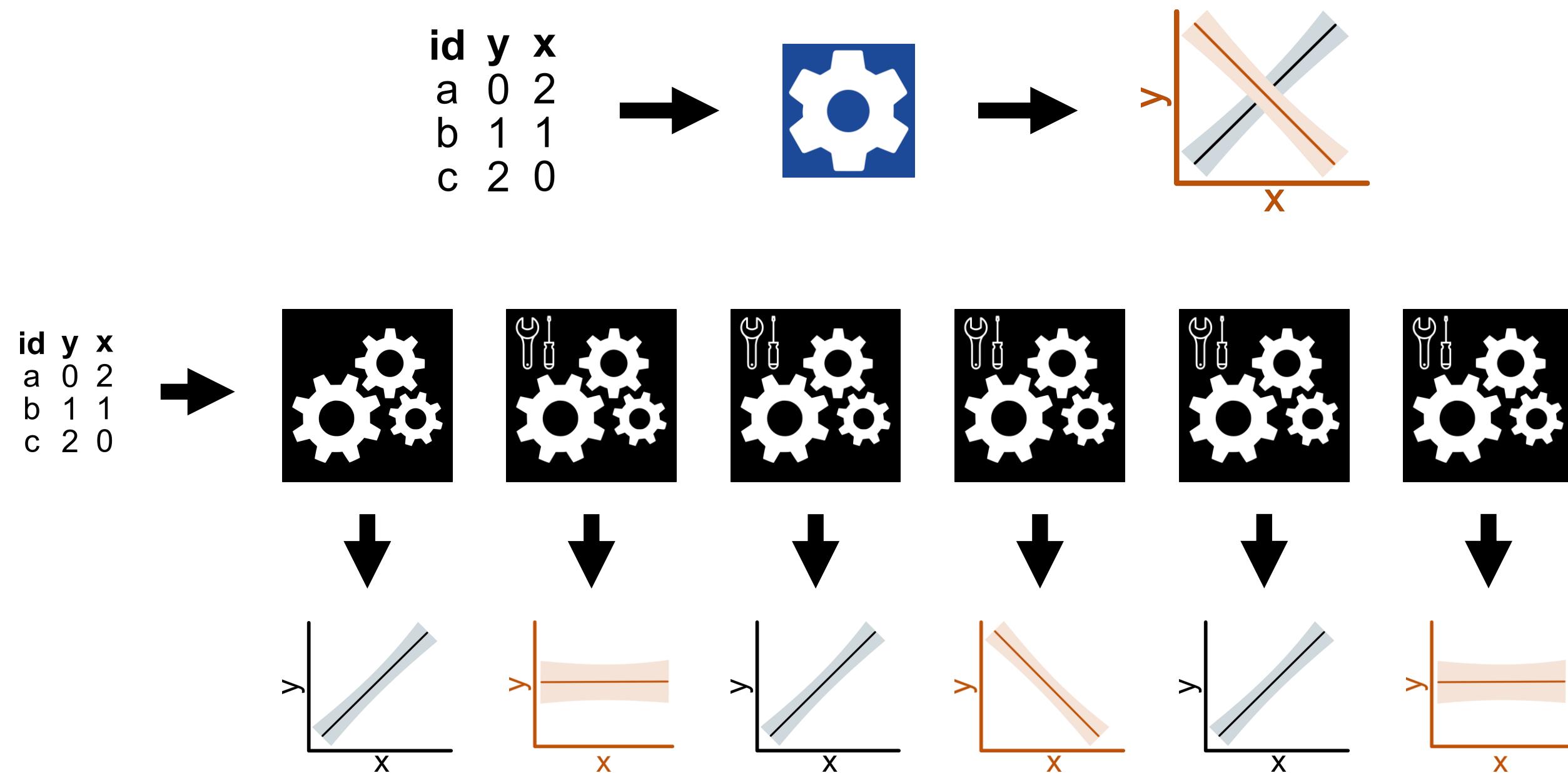
id	y	x
a	0	2
b	1	1
c	2	0

T_1 T_2 T_N



A better approach

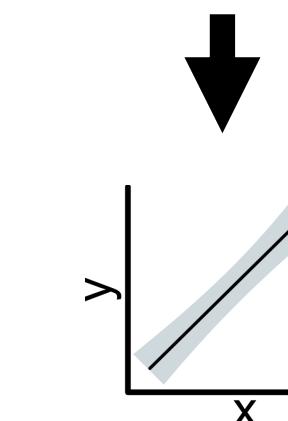
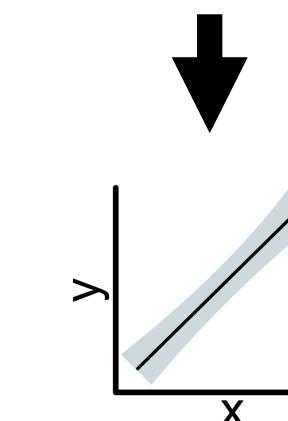
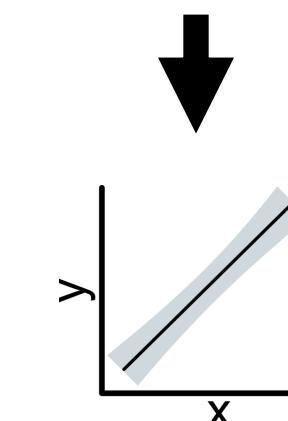
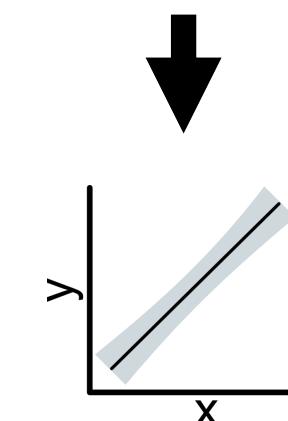
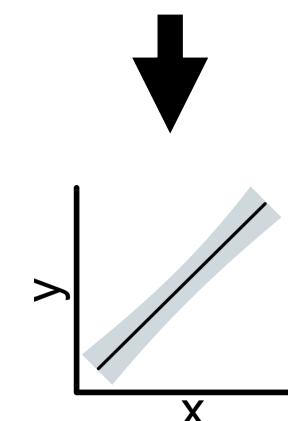
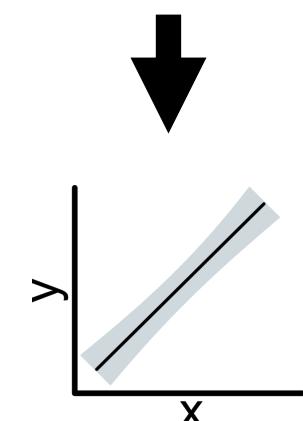
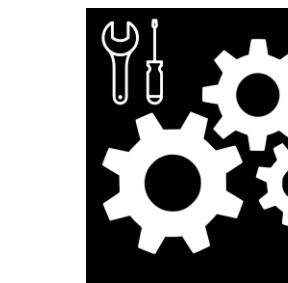
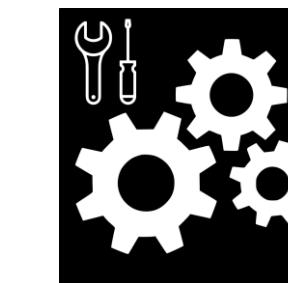
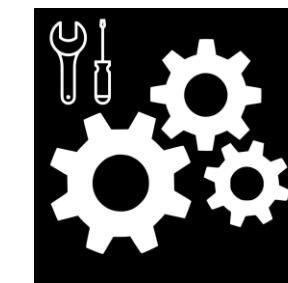
1. Learn modelling (regression) techniques
2. Try to answer your question with multiple models



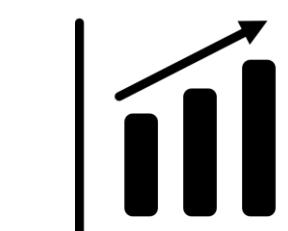
A better approach

1. Learn modelling (regression) techniques
2. Try to answer your question with multiple models
3. Compare their findings with EDA

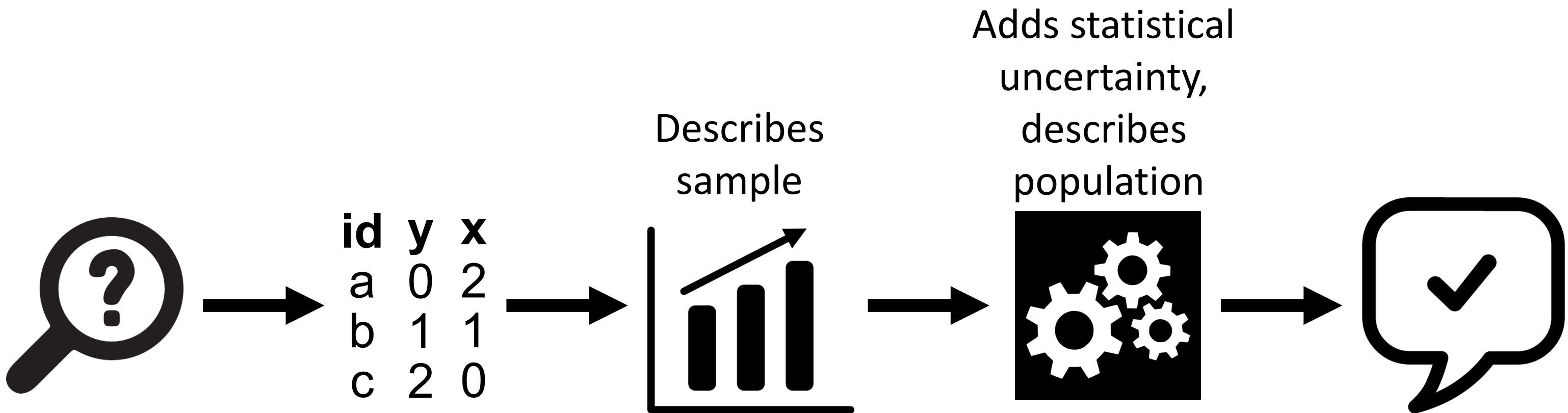
id	y	x
a	0	2
b	1	1
c	2	0

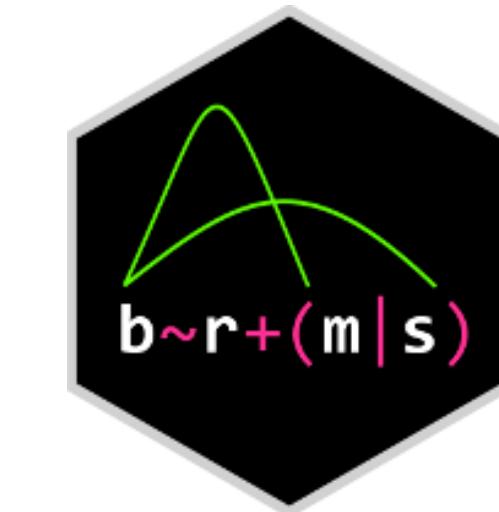
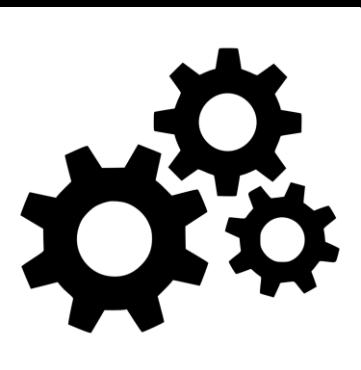


id	y	x
a	0	2
b	1	1
c	2	0



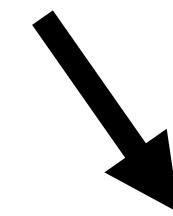
**Analyses give
consistent findings -
a more confident
answer
to your question!**





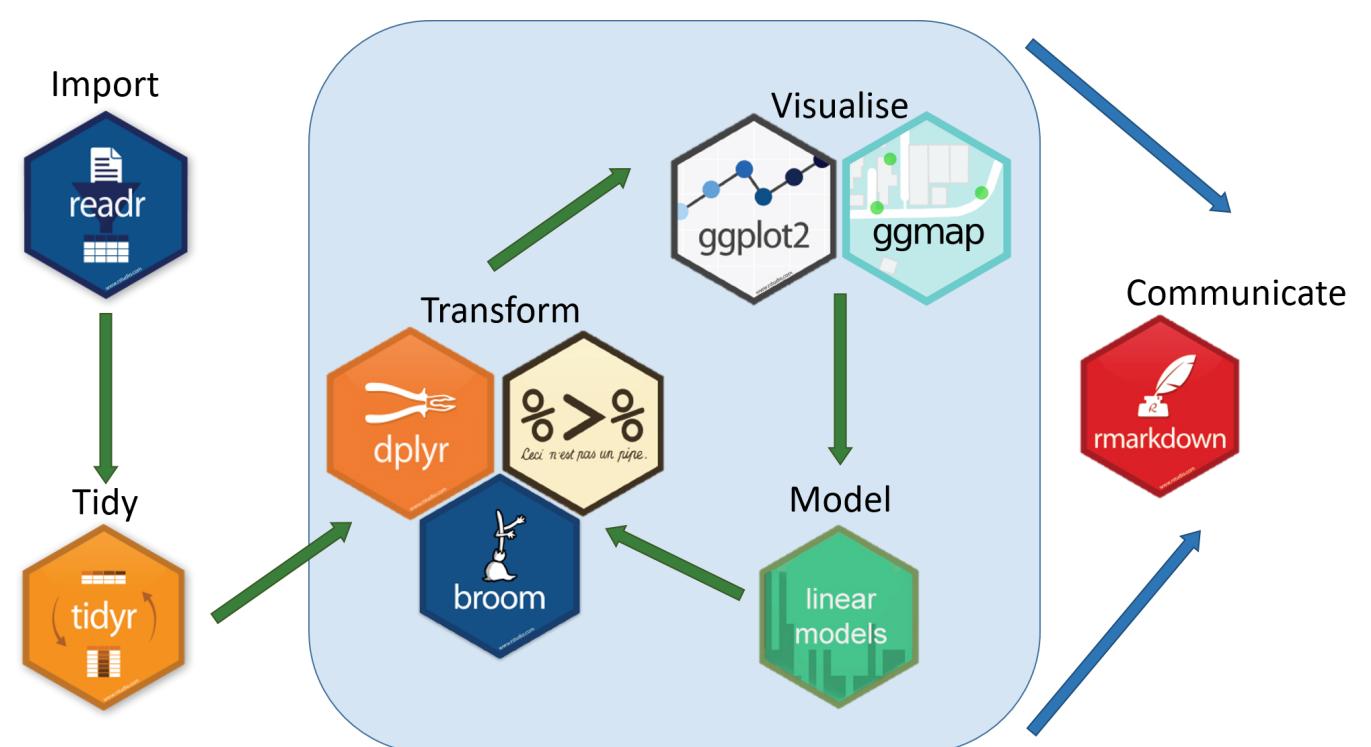


A large, bold blue letter 'R' is centered within a grey oval. The 'R' has a thick, sans-serif font style.

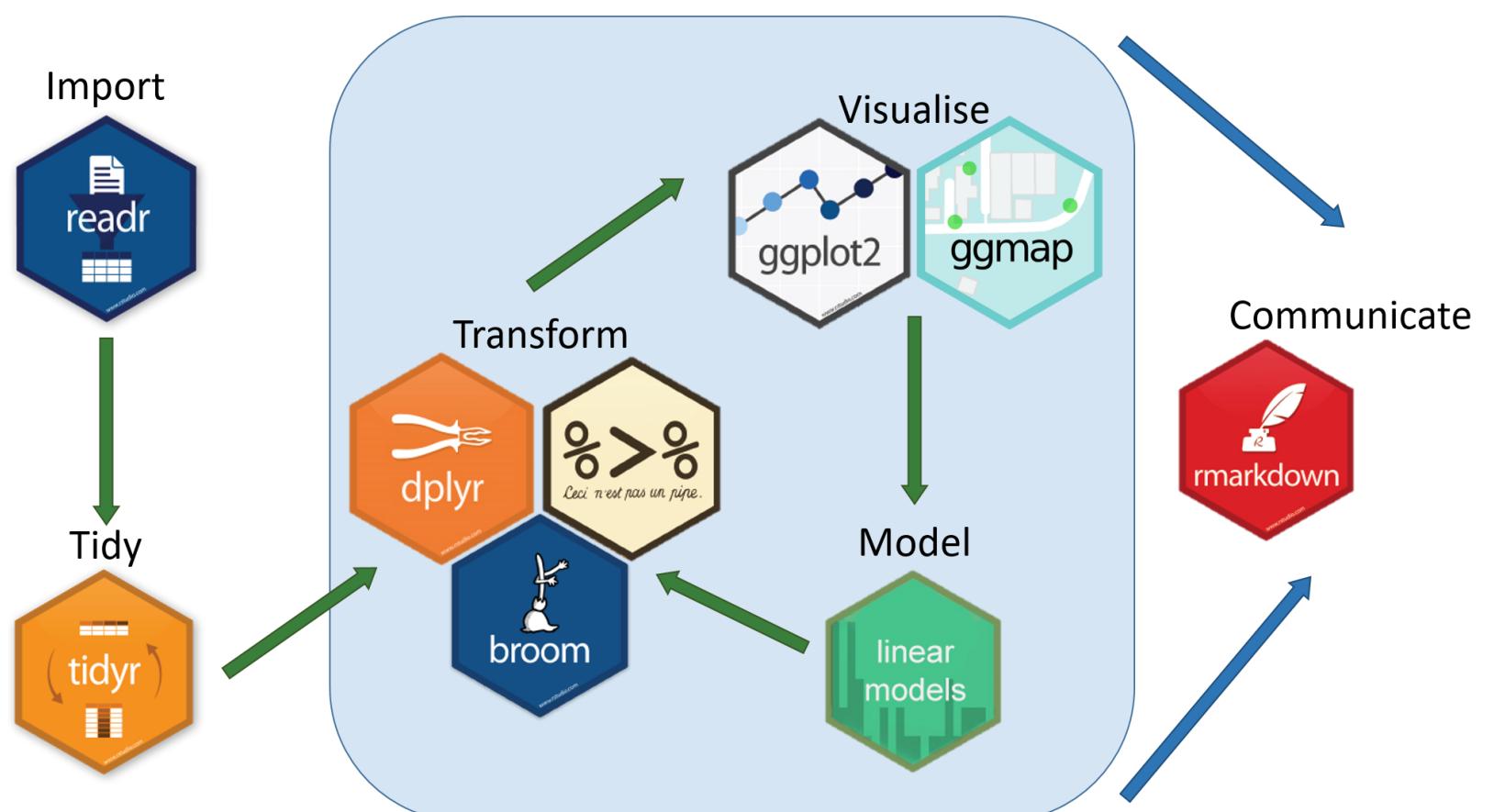
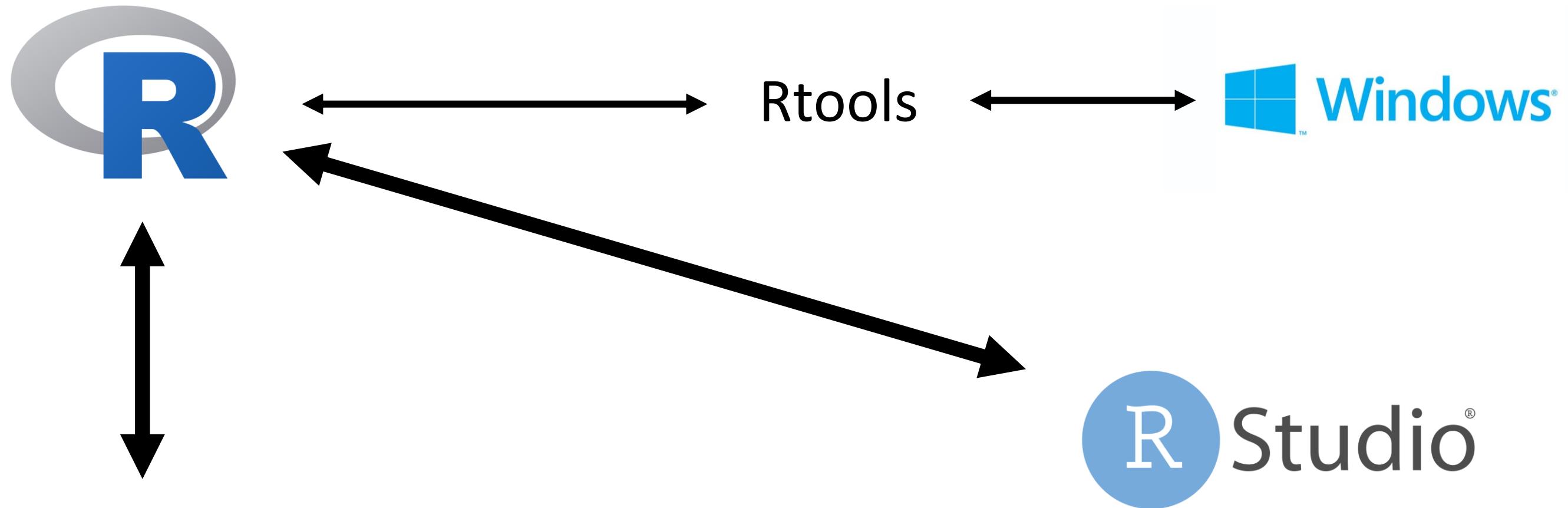


What is R?

- R is a language and environment for statistical computing and graphics.
- Its functionality is divided into modular packages (>18,000).
- Likely, R is the best software for statistical analysis
 - Free
 - Good graphics capabilities
 - Active community (Stack Overflow)



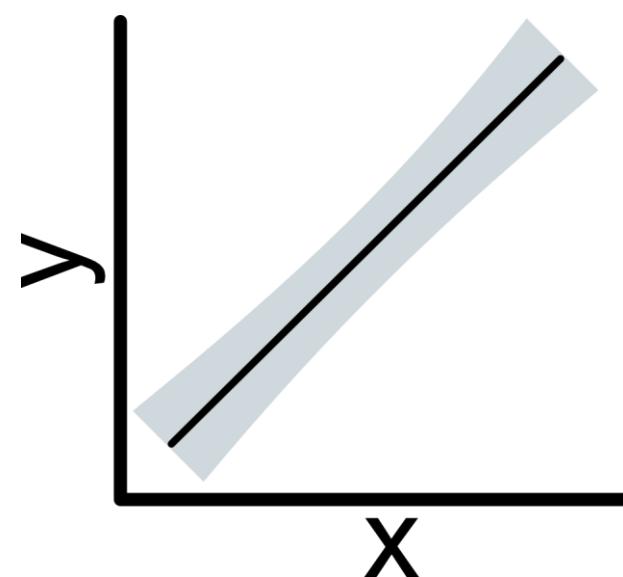
Set-up



Code

```
library(tidyverse)  
print(figure)
```

Output



Not only for statistical analysis

NICK STRAYER

I have made visualizations viewed by hundreds of thousands of people, sped up query times for 25 terabytes of data by an average of 4,800 times, and built packages for R that let you do magic.

Currently searching for a data science position that allows me to building visualization and machine learning to help people explore and understand their data.

EDUCATION

- 2020 | 2015 • PhD. Candidate, Biostatistics Vanderbilt University
 - Working on Bayesian network models & interactive visualization platforms
 - University Graduate Fellow
- 2015 | 2011 • B.S., Mathematics, Statistics (minor C.S.) University of Vermont
 - Thesis: An agent based model of Diel Vertical Migration patterns of Mysis diluviana

RESEARCH EXPERIENCE

- 2020 | 2015 • Graduate Research Assistant TBILab (Yaoxin Xu's Lab) Vanderbilt
 - Primarily working with large EHR and Biobank datasets.
 - Developing network-based methods to investigate and visually relevant patterns in data.
- 2018 | 2017 • Data Science Researcher Data Science Lab Johns Hopkins
 - Building R Shiny applications in the contexts of wearables and education.
 - Work primarily done in R Shiny and Javascript (node and d3.js)
- 2015 | 2013 • Undergraduate Researcher Rubenstein Ecosystems Science Laboratory University
 - Analyzed and visualized data for CATOS fish tracking project.
 - Head of data mining project to establish temporal trends in population densities of Mysis diluviana (Mysis).
 - Ran project to mathematically model the migration patterns of Mysis (honors thesis project.)

Made with the R package [pagedown](#).
The source code is available at github.com/nstrayer/cv.
Last updated on 2019-10-02.

What and why

My assumptions about you

How to use and understand this project

R setup

We have updates

Statistical rethinking with brms, → tidyverse: Second

on of Richard McElreath's (2020a) text, *Statistical Rethinking: A Bayesian course with examples in R and Stan*, the models he covered with Paul Bürkner's *brms* package (Bürkner, 2017, 2018, 2020), which makes it easy to fit Bayesian regression models in R (R Core Team, 2020) using Hamiltonian Monte Carlo. I also prefer plotting and data wrangling with the packages from the *tidyverse* (Wickham, 2019; Wickham et al., 2019). So we'll be using those methods, too.



English vs Base R vs Tidyverse

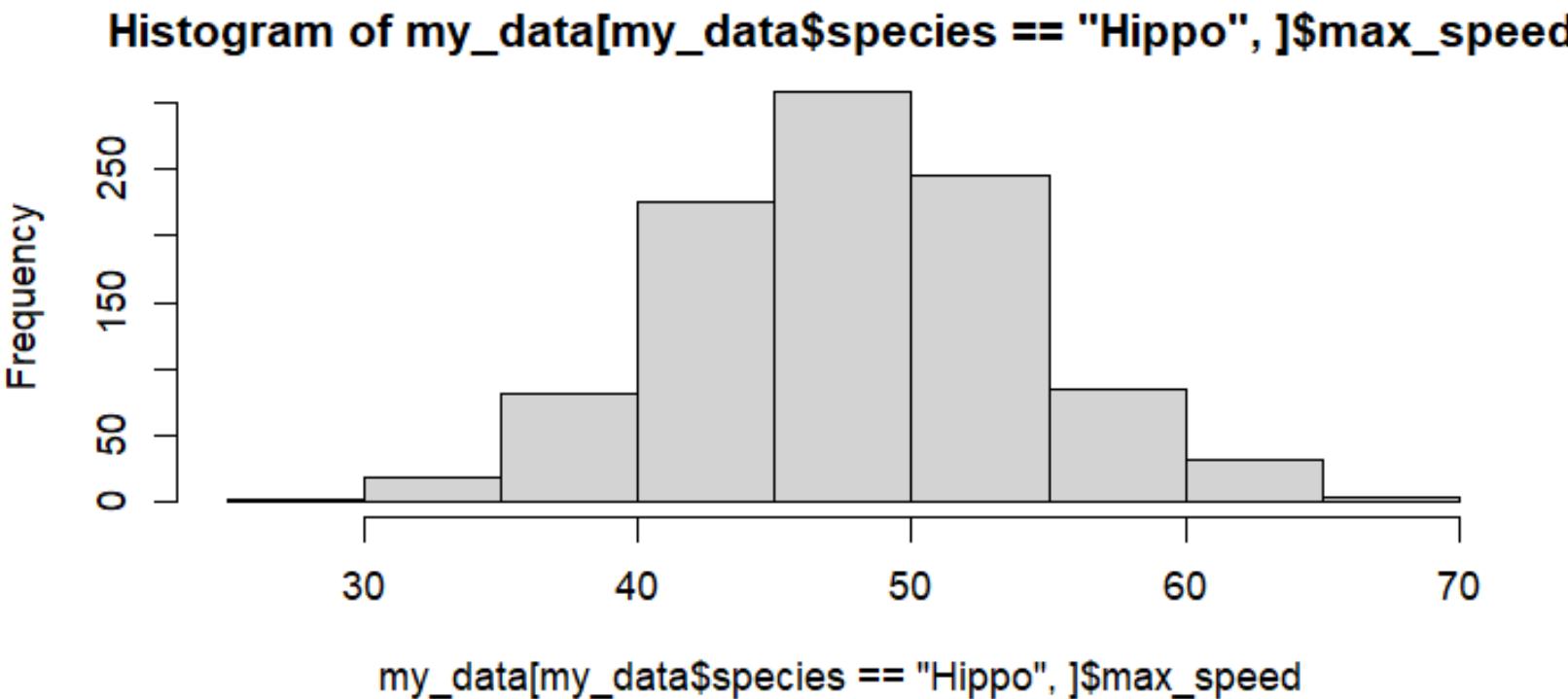


English

Use my hippos-and-elephants data, filter out hippos and plot their maximum running speed values on a histogram.

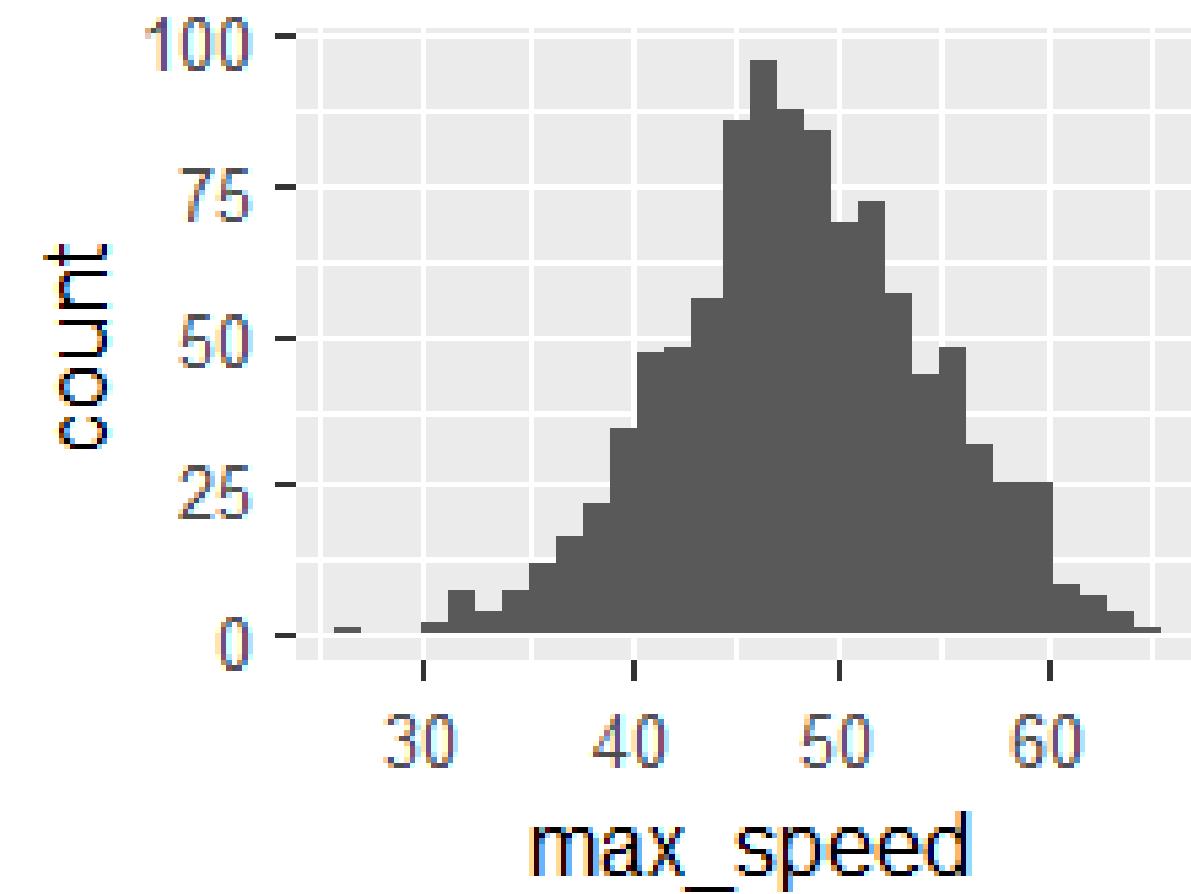
Base R

```
hist(x = my_data[my_data$species ==
  "Hippo", ]$max_speed)
```



Tidyverse

```
my_data %>%
  filter(species == "Hippo") %>%
  ggplot() +
  geom_histogram(aes(x = max_speed))
```



It is a programming language

A

```
my_data %>%
  filter(species == "Hippo") %>%
  ggplot() +
  geom_histogram(aes(x = max_speed))
```

B

```
my_data %>%
  filter(species != "Elephant") %>%
  ggplot(aes(x = max_speed))+
  geom_histogram()
```

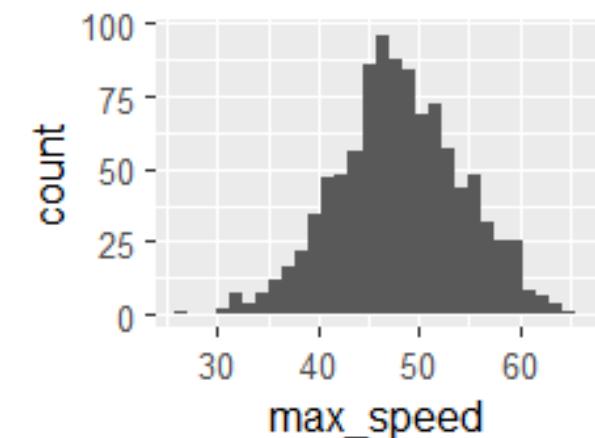
C

```
my_data %>%
  filter(species %>% str_detect("pp")) %>%
  ggplot(aes(x = max_speed))+
  geom_histogram()
```

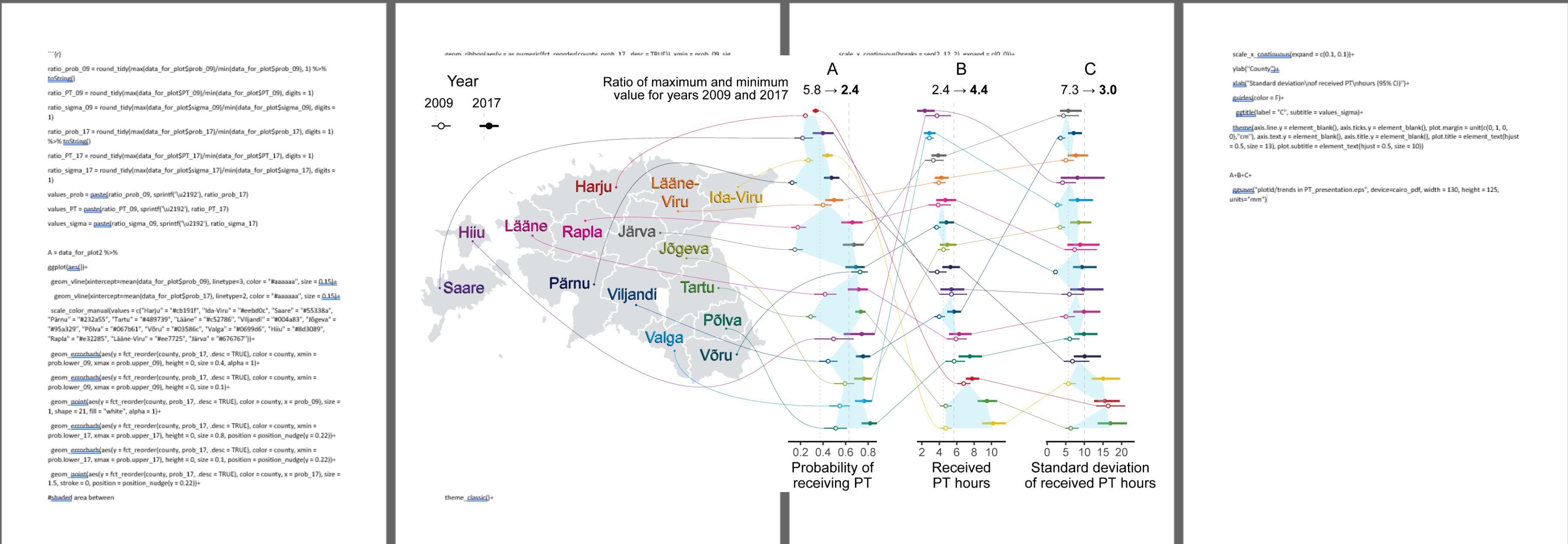
D

```
filtered_data = my_data %>%
  filter(species == "Hippo")

filtered_data %>%
  ggplot(aes(x = max_speed))+
  geom_histogram()
```



Why to choose the hard way - R?



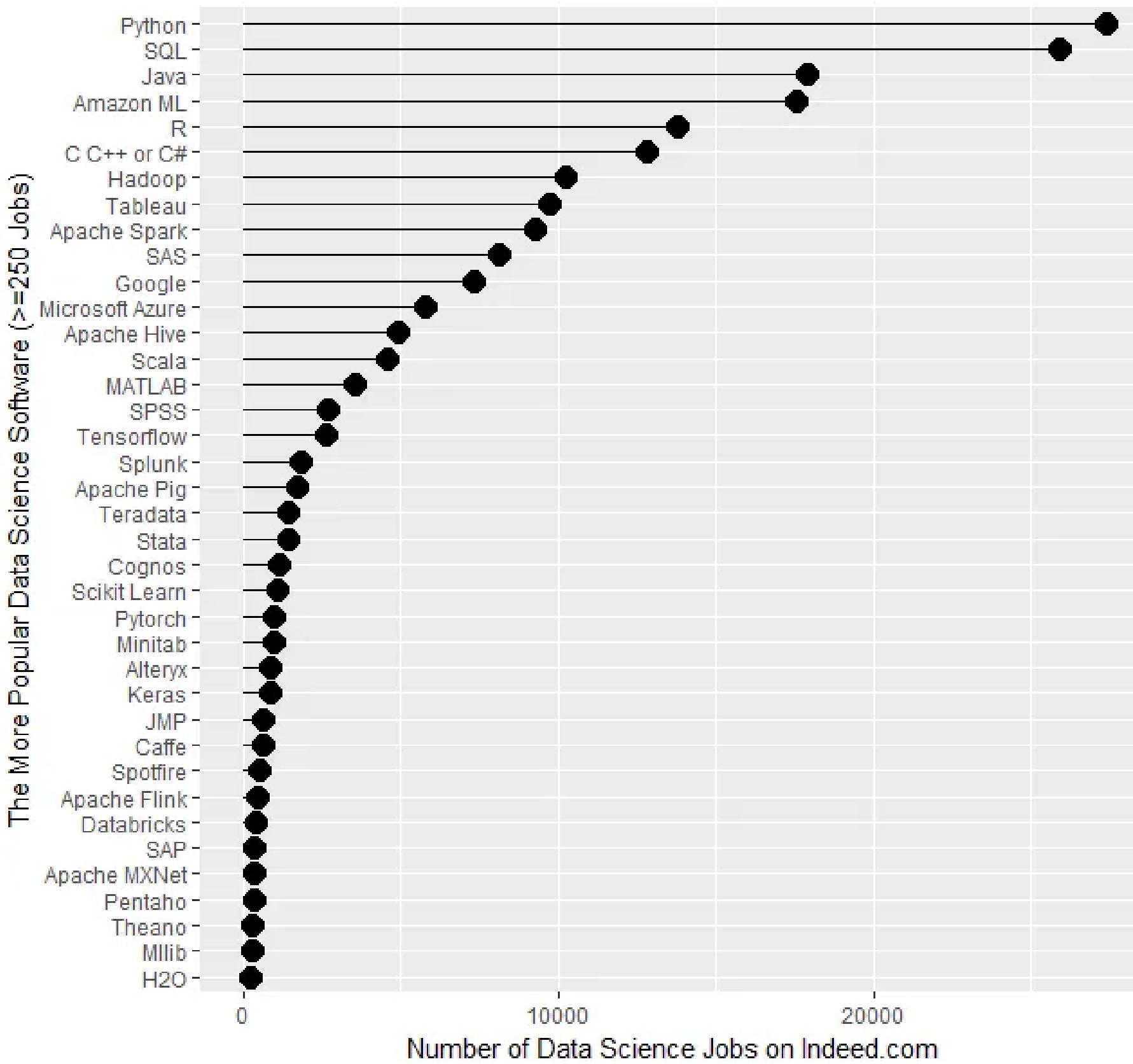
The Popularity of Data Science Software

by Robert A. Muenchen

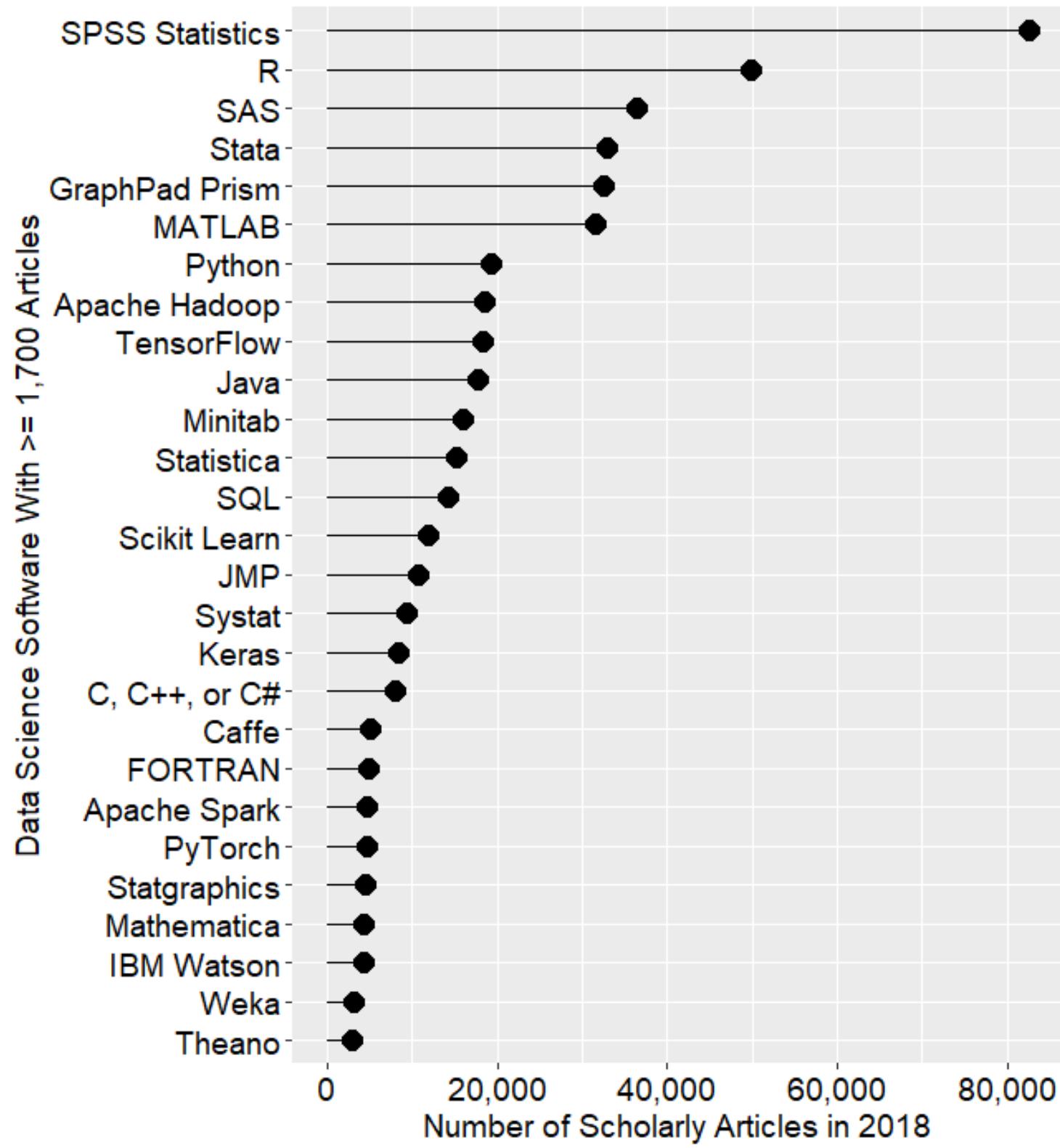
Actuate, Alpine, Alteryx, Angoss, Apache Flink, Apache Hive, Apache Mahout, Apache MXNet, Apache Pig, Apache Spark, BMDP, C, C++ or C#, Caffe, Cognos, DataRobot, Domino Data Labs, Enterprise Miner, FICO, FORTRAN, H2O, Hadoop, Info Centricity or Xeno, Java, JMP, Julia, KNIME, Lavastorm, MATLAB, Megaputer or PolyAnalyst, Microsoft, Minitab, NCSS, Oracle Data Miner, Prognoz, Python, R, RapidMiner, Salford SPM, SAP, SAS, Scala, Spotfire, SPSS, SPSS Modeler, SQL, Stata, Statgraphics, Statistica, Systat, Tableau, Tensorflow, Teradata, Vowpal Wabbit, WEKA/Pentaho, and XGboost.

- Job Advertisements
- Scholarly Articles
- Surveys of Use
- Books
- Blogs
- Discussion Forum Activity
- Growth in Capability

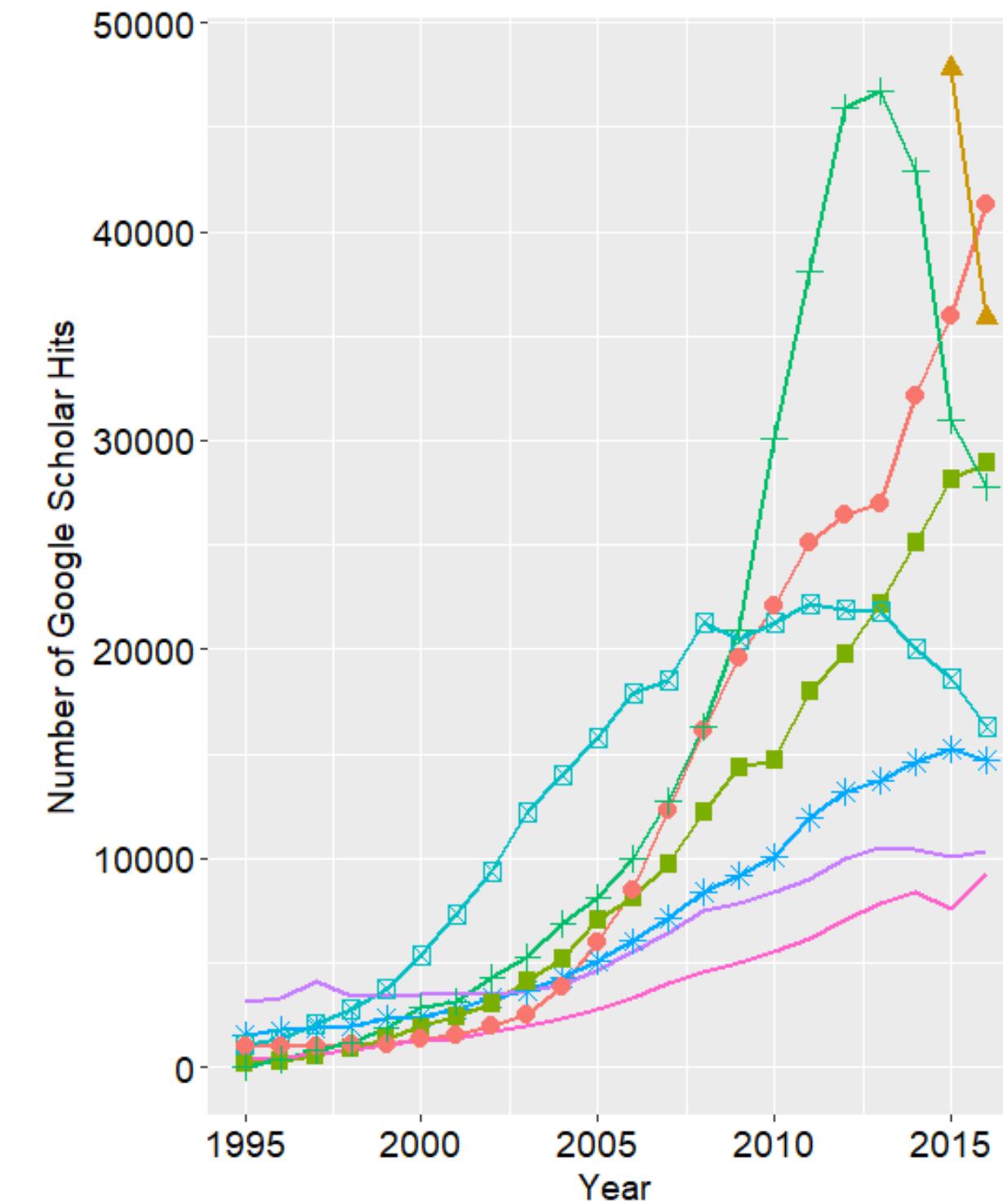
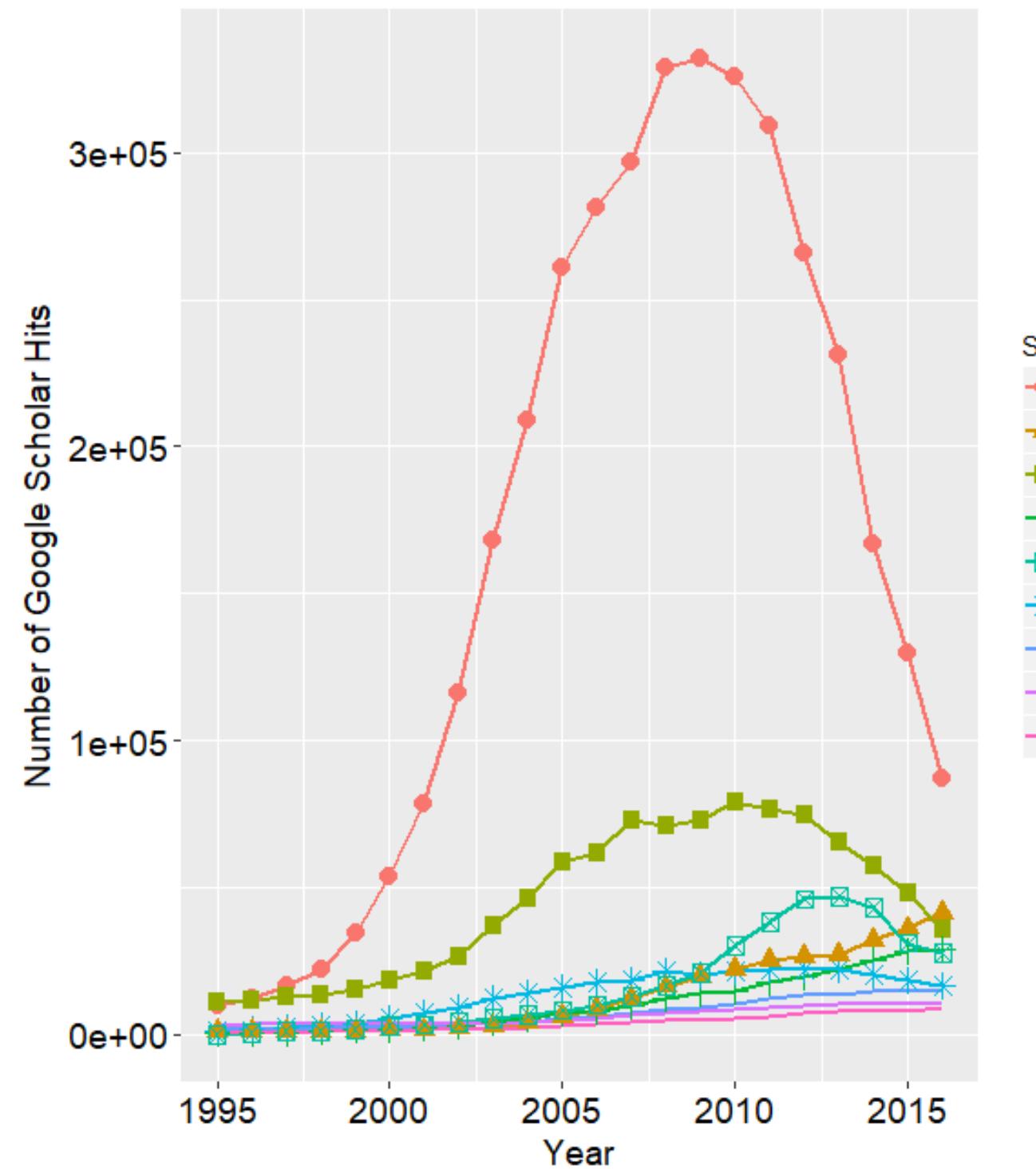
Job Advertisements



Scholarly Articles

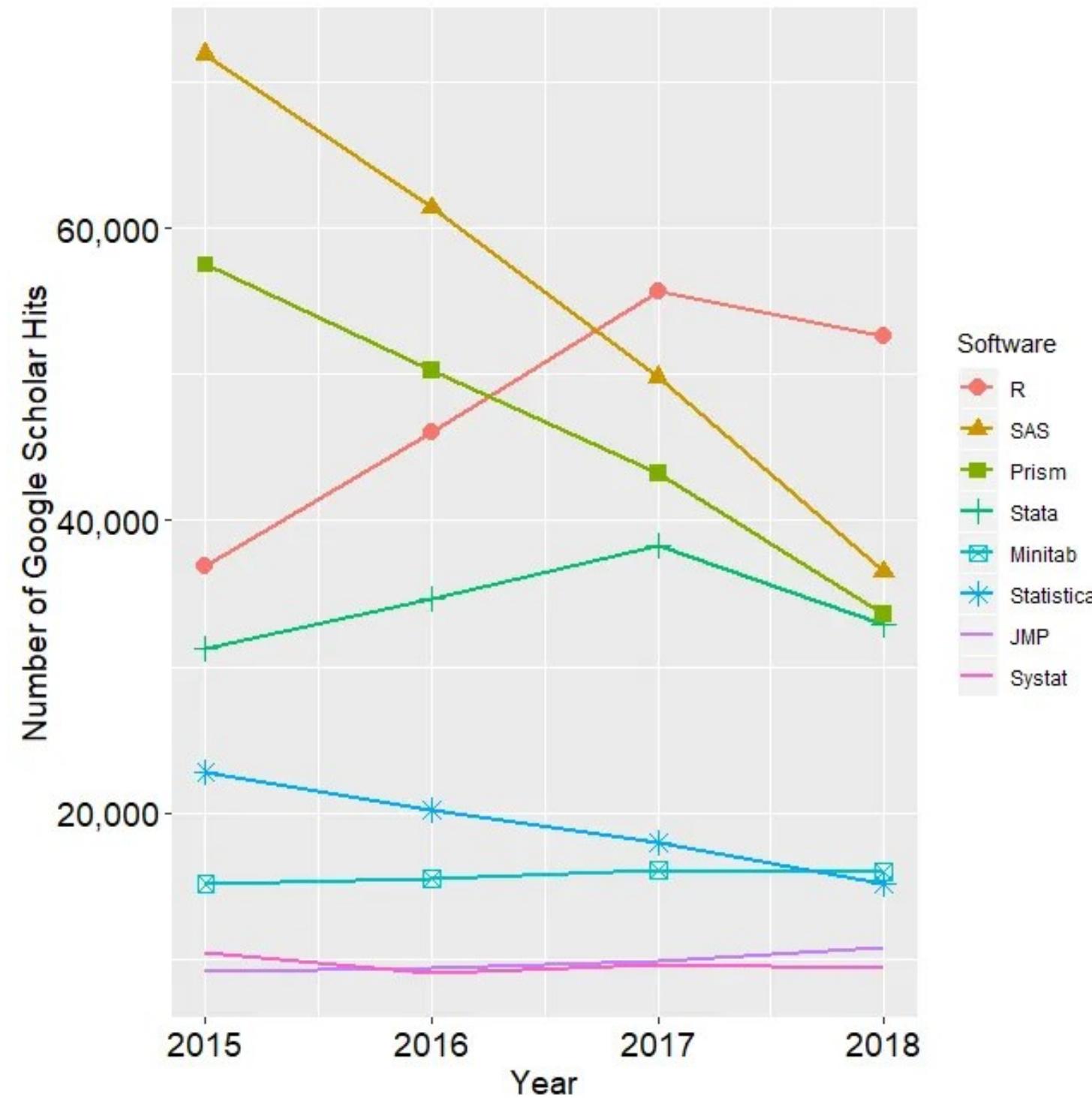


Scholarly Articles: trends for classic statistics packages (1995-2015)



The number of Google Scholar citations for each classic statistics package from 1995 through 2016, this time with SPSS removed and SAS included only in 2014 and 2015. The removal of SPSS and SAS expanded scale makes it easier to see the rapid growth of the less popular packages.

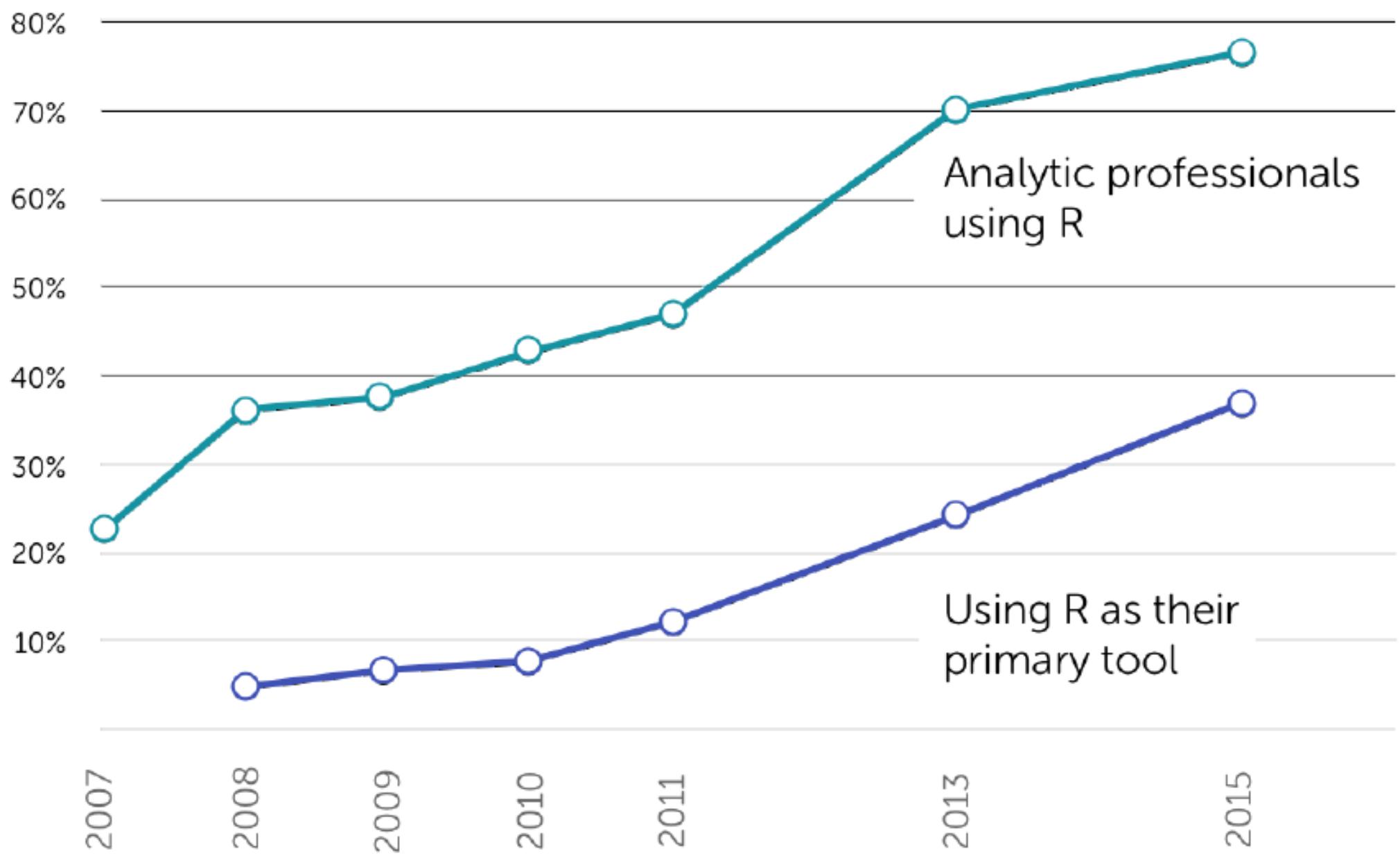
Scholarly Articles: trends for classic statistics packages (2015-2018)



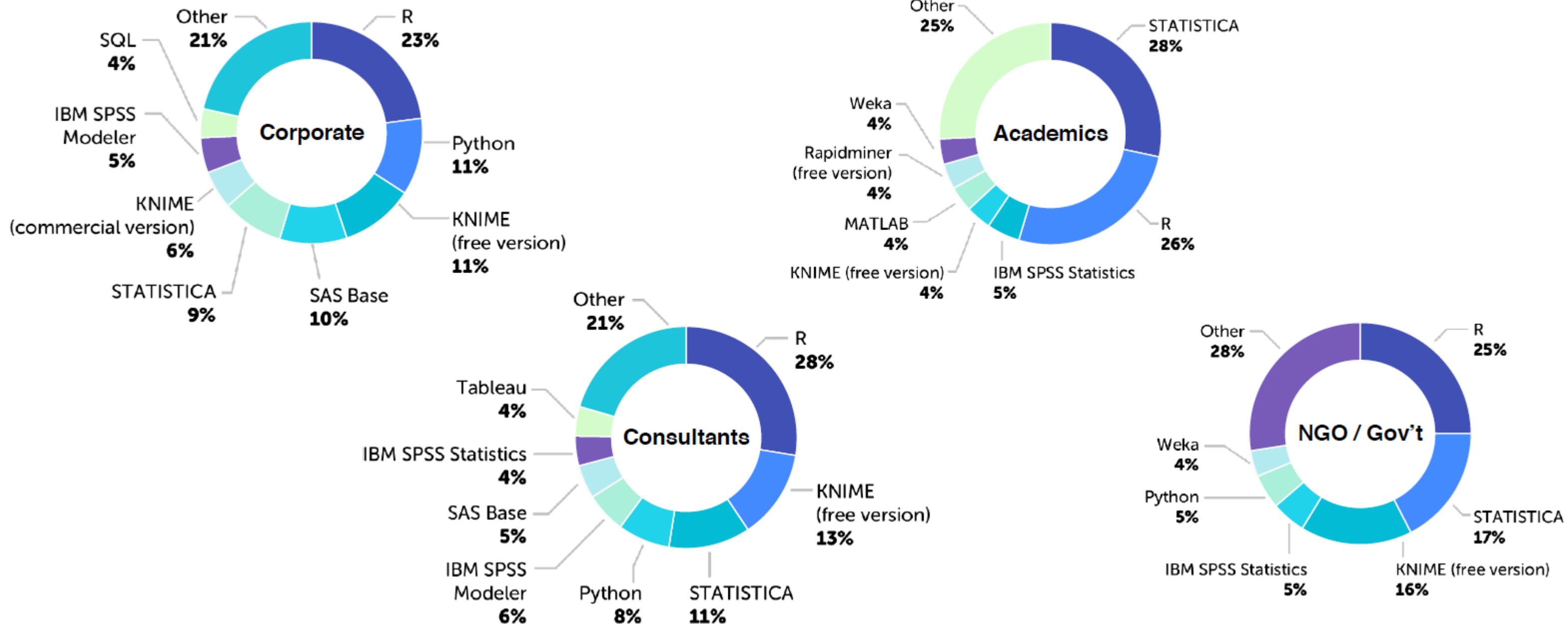
Surveys of Use: Rexter Analytics (2017)

- 67 questions
- 10,000+ invitations emailed & promoted by newsgroups, vendors and bloggers
- Respondents: 1,123 analytic professionals from 91 countries
- Data collected in first half of 2017

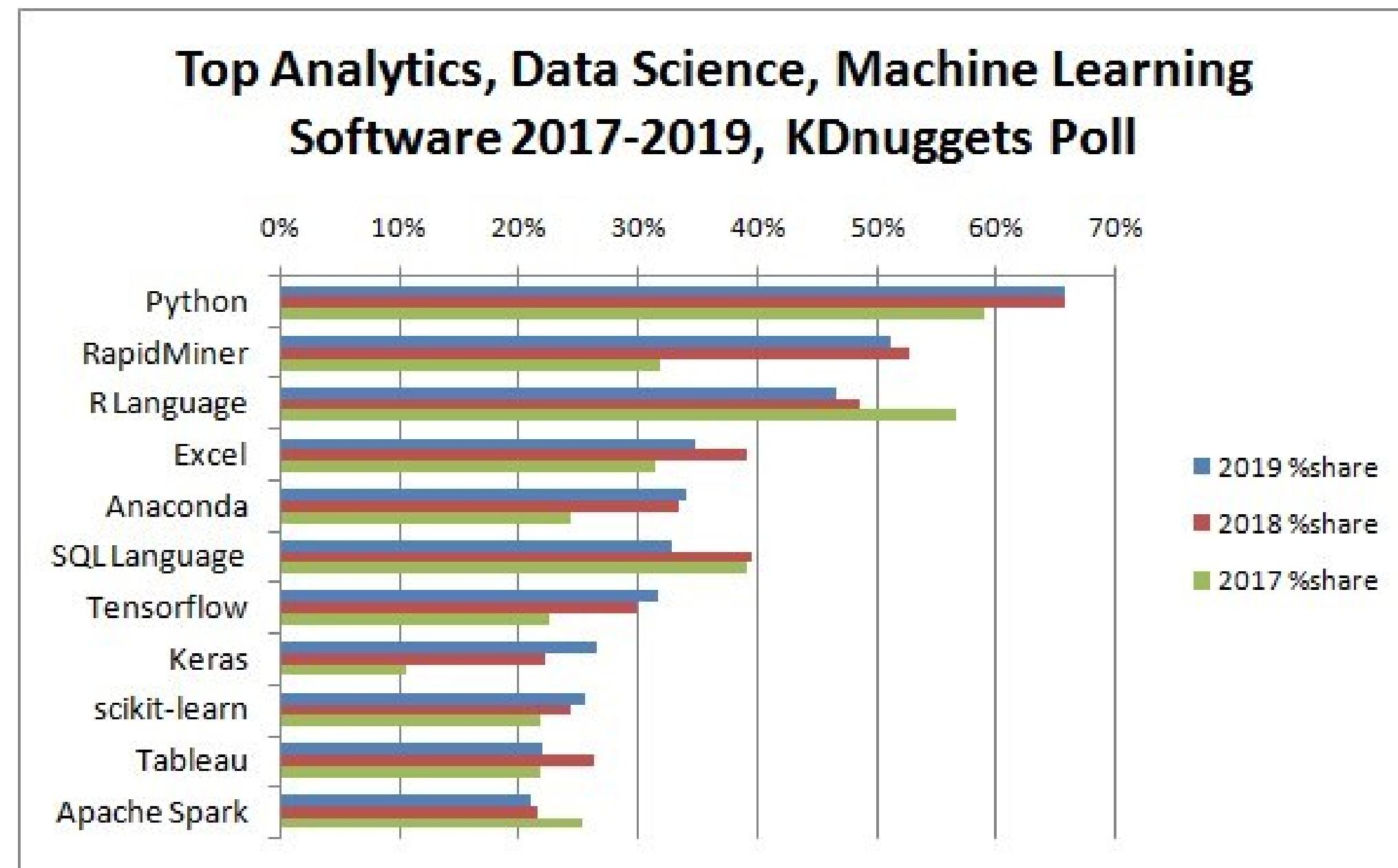
A Decade of Dramatic Growth in the Use of R



Surveys of Use: Rexer Analytics (2017)



Surveys of Use: KDnuggets



Books on Amazon.com

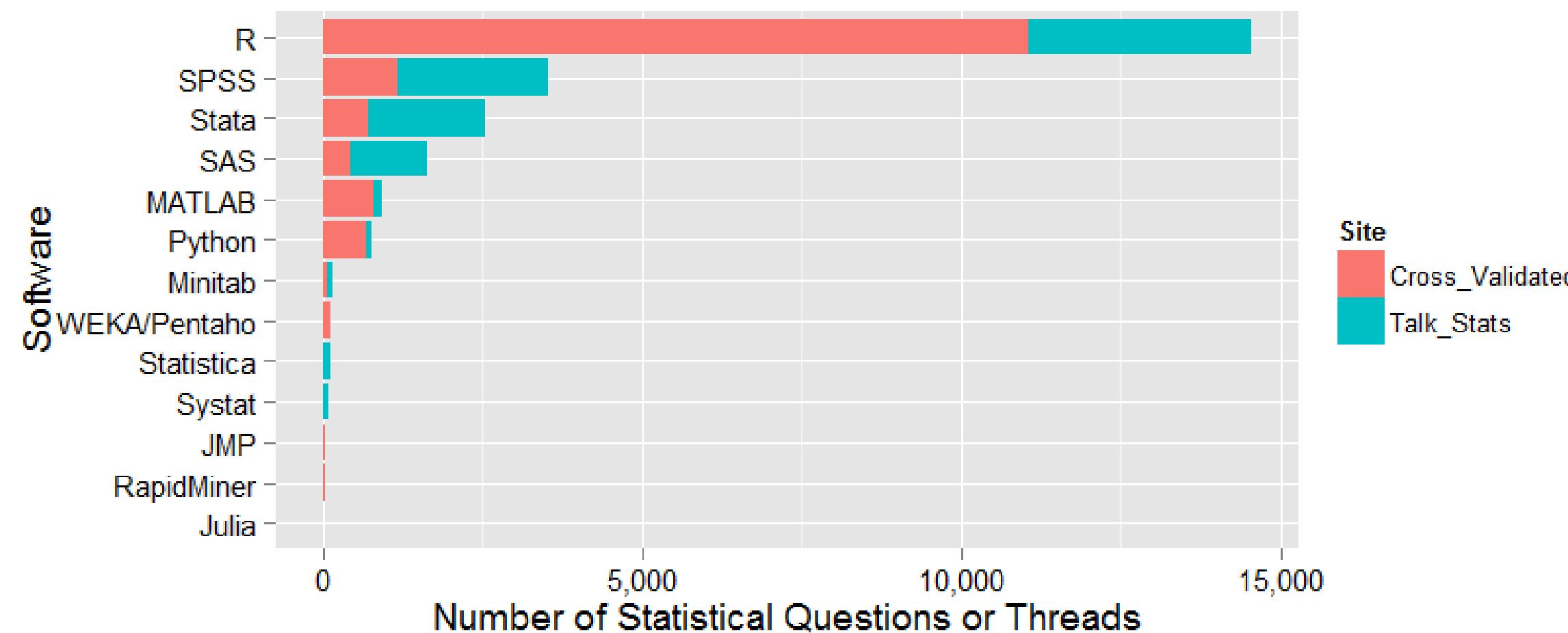
Software	Number of Books
SAS	576
SPSS Statistics	339
R	240 [Corrected from blog post: 172]
JMP	97
Hadoop	89
Stata	62
Minitab	33
Enterprise Miner	32

Blogs

Software	Number	Source
R	550	R-Bloggers.com
Python	60	SciPy.org
SAS	40	PROC-X.com , sasCommunity.org Planet
Stata	11	Stata-Bloggers.com

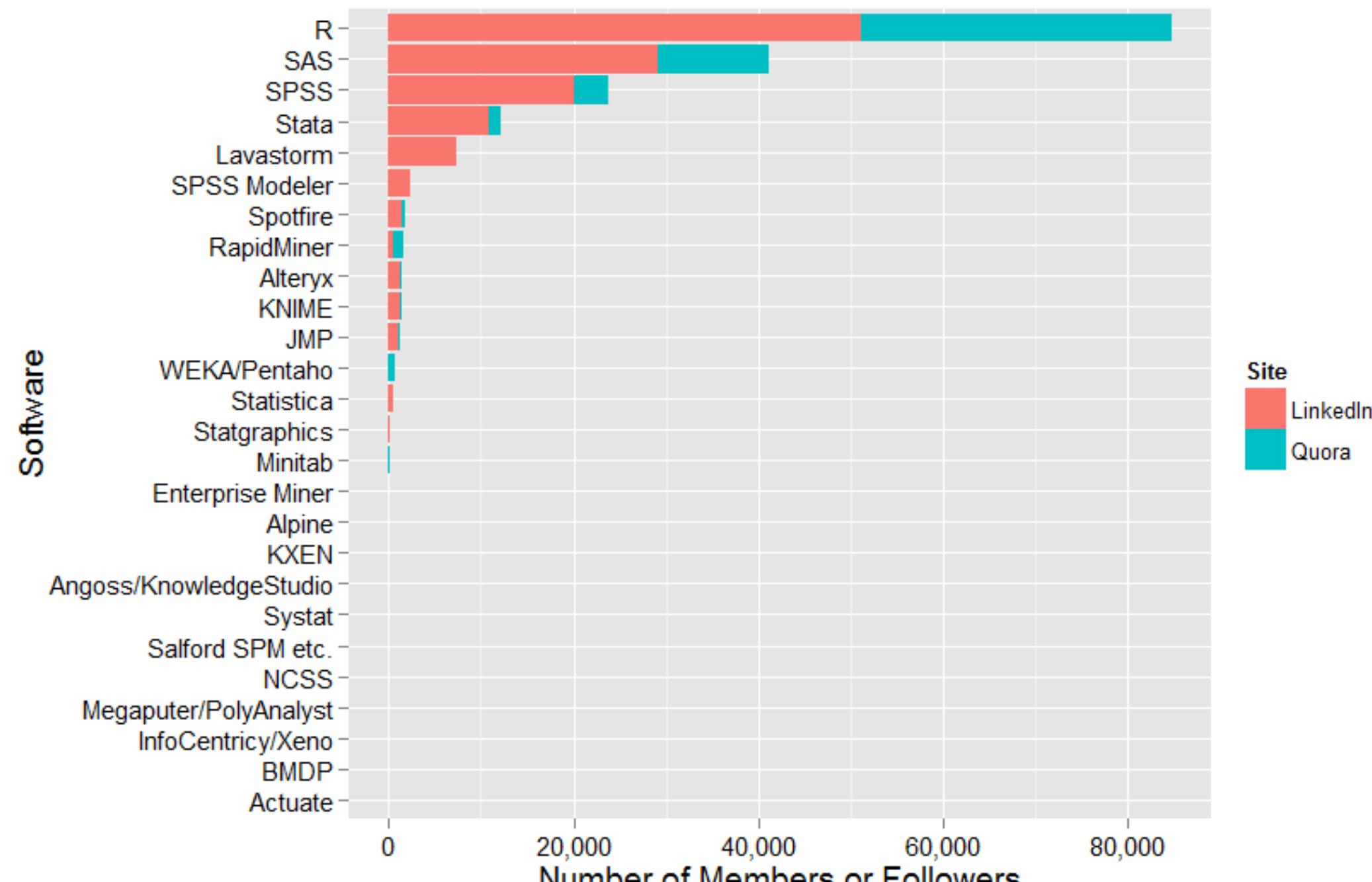
Number of blogs devoted to each software package on April 7, 2014,
and the source of the data.

Discussion Forum Activity: Cross Validated and Talk Stats



Number of statistical questions or threads on Cross Validated and Talk Stats.

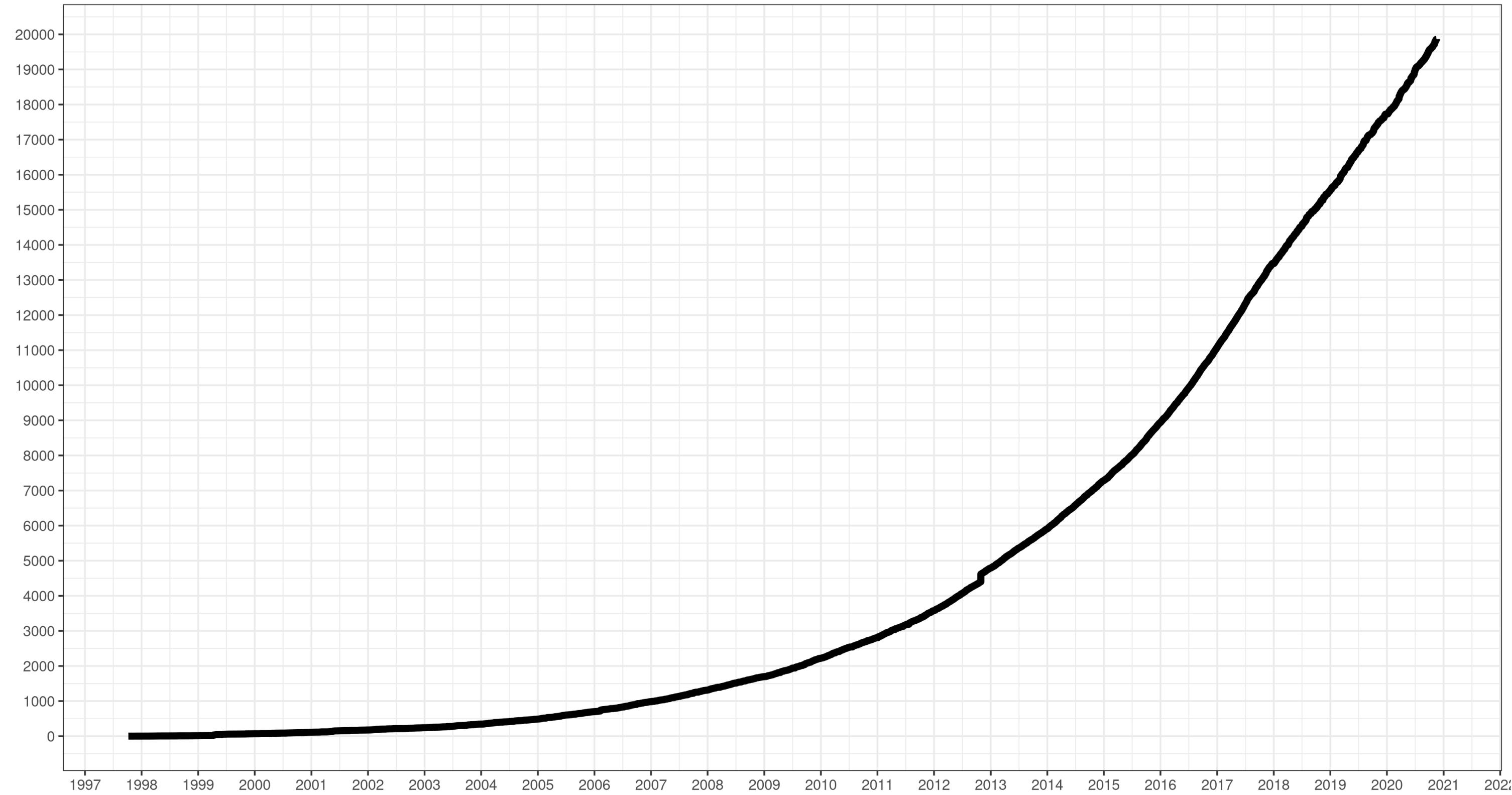
Discussion Forum Activity: LinkedIn and Quora



Number of members or followers on LinkedIn and Quora.

Growth in Capability

Number of R packages ever published on CRAN

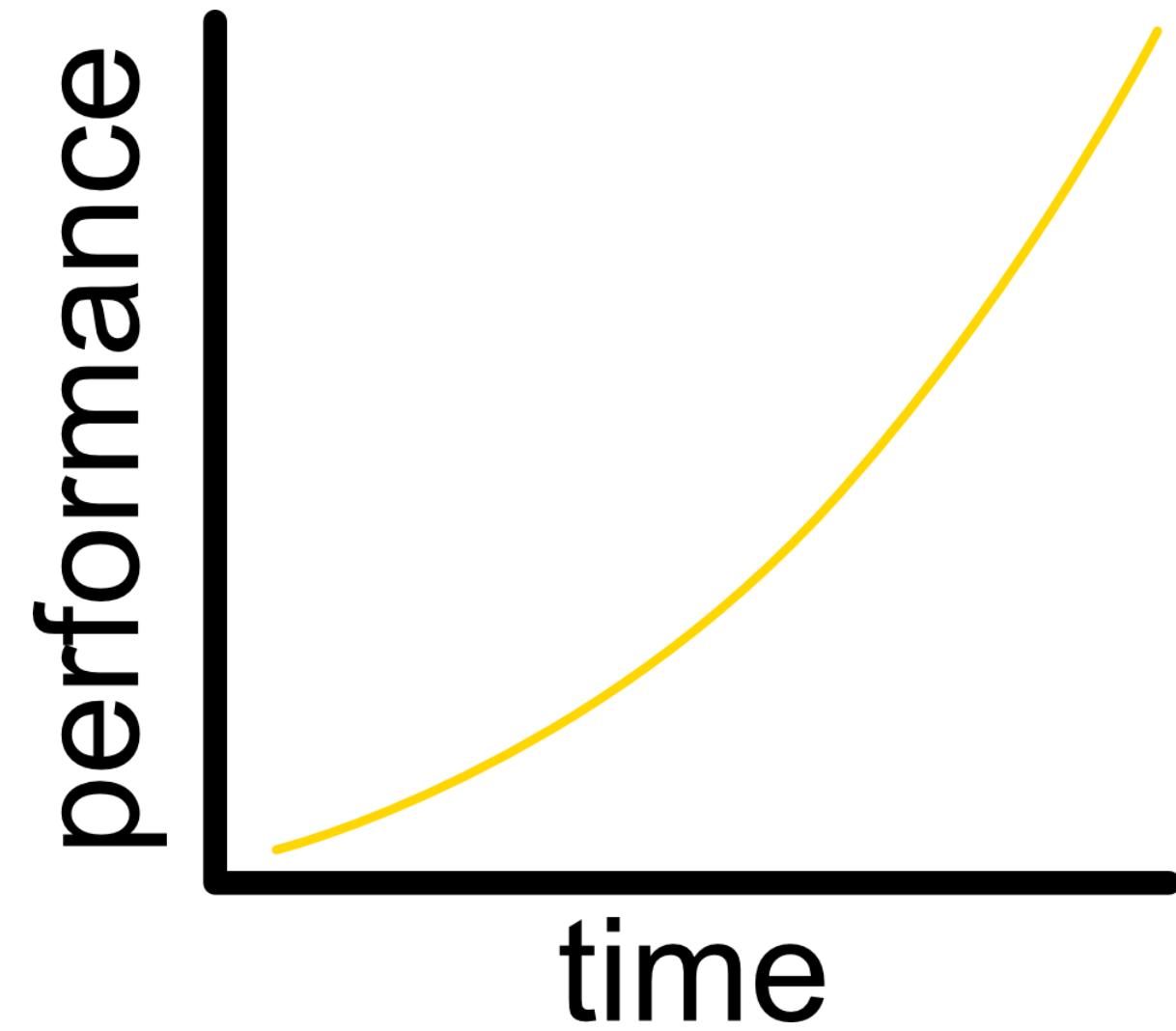


R software: unfriendly but probably the best

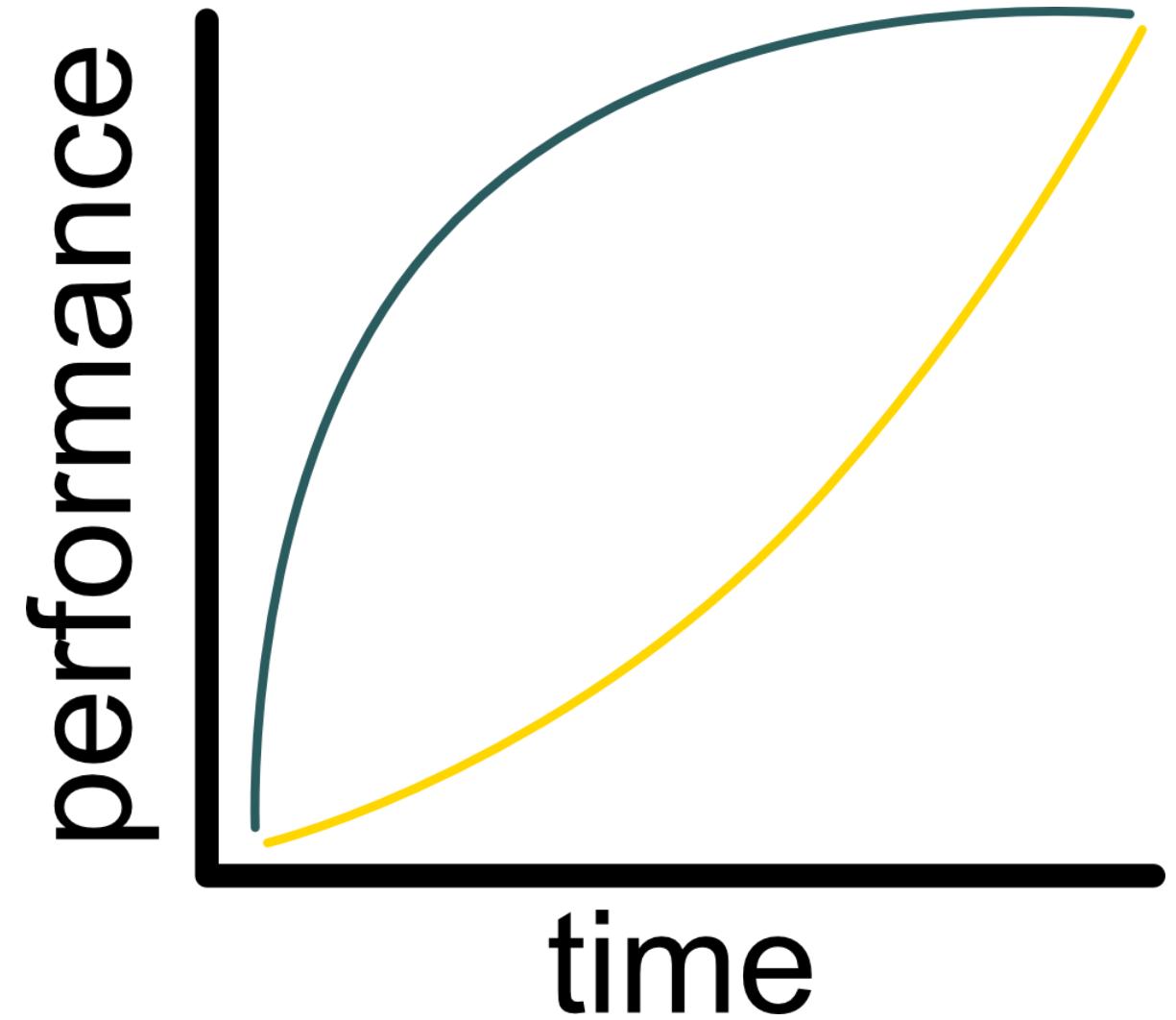
[Branimir K. Hackenberger](#)

R is perhaps the today's best tool for statistical data processing, ranging from frequentist statistics, Bayesian statistics, meta-analysis, machine and deep learning to parallel computing and Big Data processing. For free!

Expect frustration



Expect frustration



Your friends



Google

(R OR tidyverse) AND "your problem"

stack overflow