



Data Science & Machine Learning

Dan S. Reznik
www.dat-sci.com
December, 2019

Agenda

- Morning: 8:30~12:00
 - The Business Impact of Data Science and Machine Learning
 - What is Machine Learning
 - What is Data Science
 - Toolbox: R, Python, etc.
- Afternoon: Hands-on w/ R: 13:00~17:30
 - Set up DS environment on the cloud
 - GitHub, Amazon Machine Instance
 - Data Wrangle
 - Option1: Google Forms
 - Option2: Open Data: data.gov
- Takeaways

AI vs Traditional Businesses

- Physical => Digital
- Product, Service => Cloud Platform
- Oil => Data
- Manuf. & Ownership => Sharing economy
- Proprietary Software => Open Frameworks
- Algorithms => Leverage diverse datasets

5,862 views | Apr 5, 2019, 03:59pm

AI Is Destroying Traditional Business Thinking

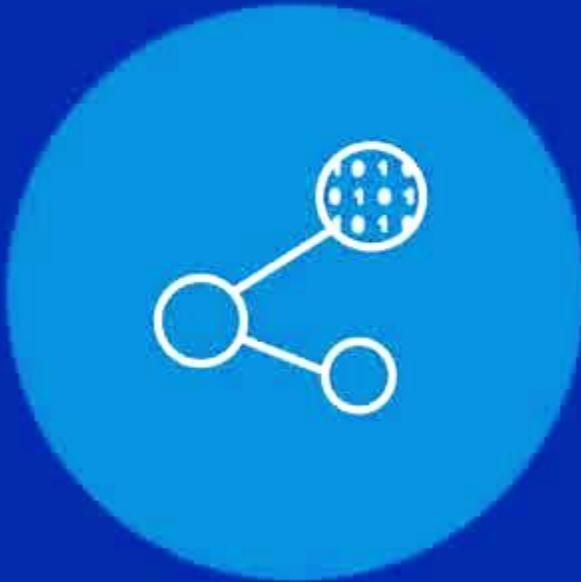


Tom Davenport Contributor ⓘ

[Enterprise & Cloud](#)

Business is being transformed by three trends

Big Data



Cloud



Intelligence



The Jobs Landscape in 2022

emerging
roles,
global
change
by 2022

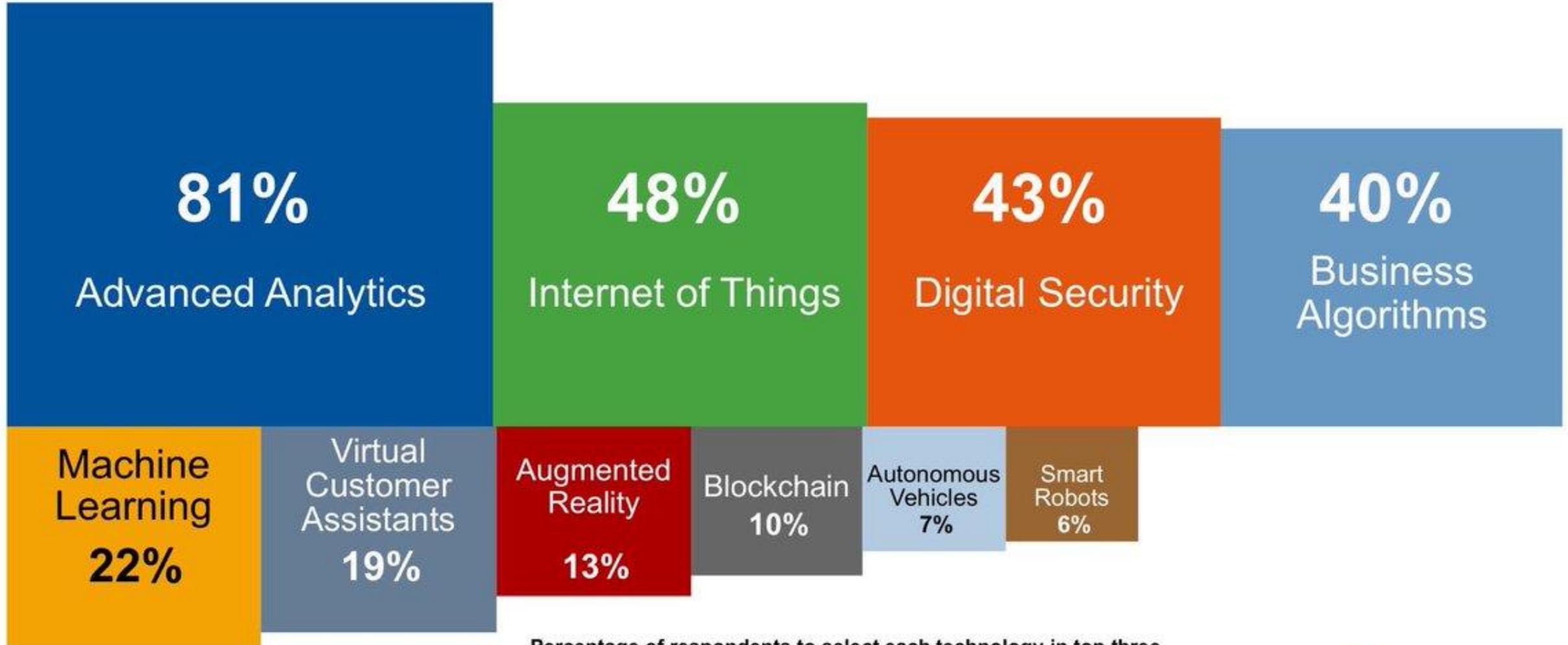
133
Million

Top 10 Emerging

1. Data Analysts and Scientists
2. AI and Machine Learning Specialists
3. General and Operations Managers
4. Software and Applications Developers and Analysts
5. Sales and Marketing Professionals
6. Big Data Specialists
7. Digital Transformation Specialists
8. New Technology Specialists
9. Organisational Development Specialists
10. Information Technology Services

Key Technologies Will Deliver Change

Q. In your opinion, which three of these technologies have the most potential to change your organization over the next five years?



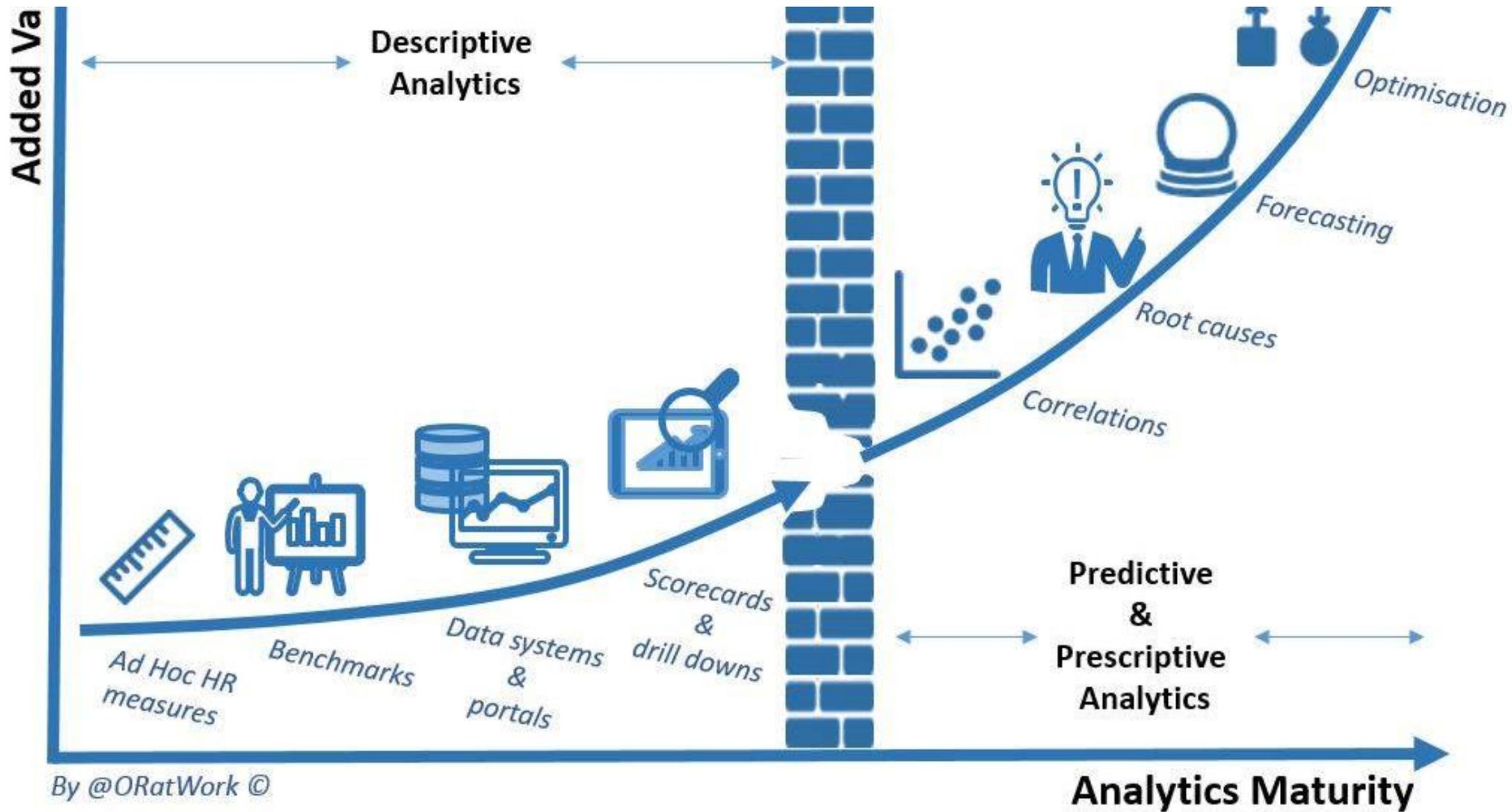
Percentage of respondents to select each technology in top three

#GartnerSYM

24 CONFIDENTIAL AND PROPRIETARY | © 2016 Gartner, Inc. and/or its affiliates. All rights reserved. Gartner and ITXpo are registered trademarks of Gartner, Inc. or its affiliates.

(c) 2019 Dan S. Reznik -- FDC

Gartner



26,472

Likes

13,911

Subscribers

6,524

Followers



5,093

Circled By



45,322

Followers



1,765

Followers



Session Traffic

280,430

29.44% ↑
vs 216,646 (prev.)

11.75% ▲
vs \$15,059 (prev.)

Wistia Video Stats



Metrics Driven
Change
Management

563
Play Count

75%
Engagement

49%
Play Rate

CallRail Today's Average Call Duration

Company: The Klip Factory

10m:8s

Based on 19 calls



Salesforce Accounts by Country



Invalid state information provided for 646 of 13124 accounts.

Quickbooks & Salesforce: CAC (Last 30 Days)



\$15,085

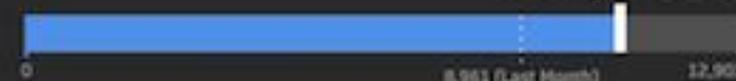
CAC Value

-29.99% ▼

vs \$21,548 prev. 30 days

Leads (This Month)

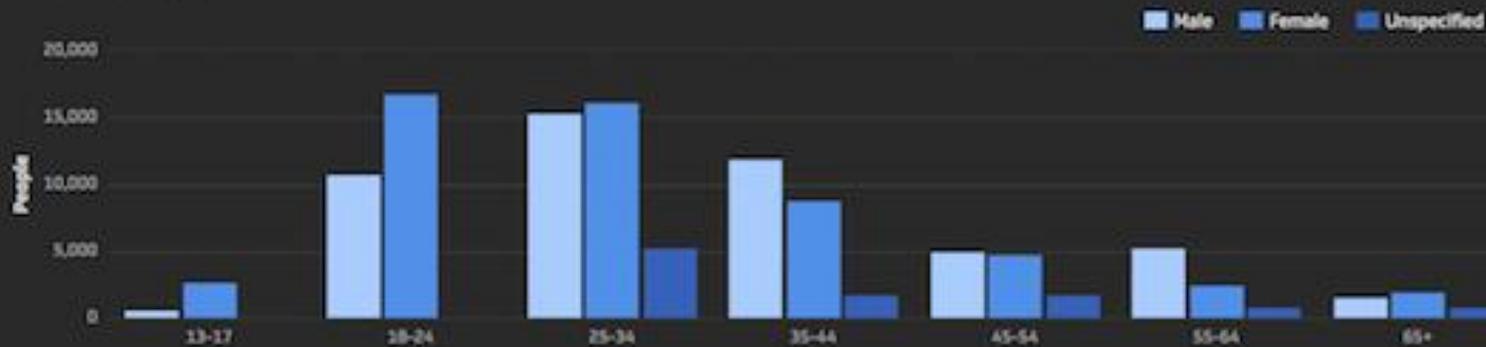
10,753



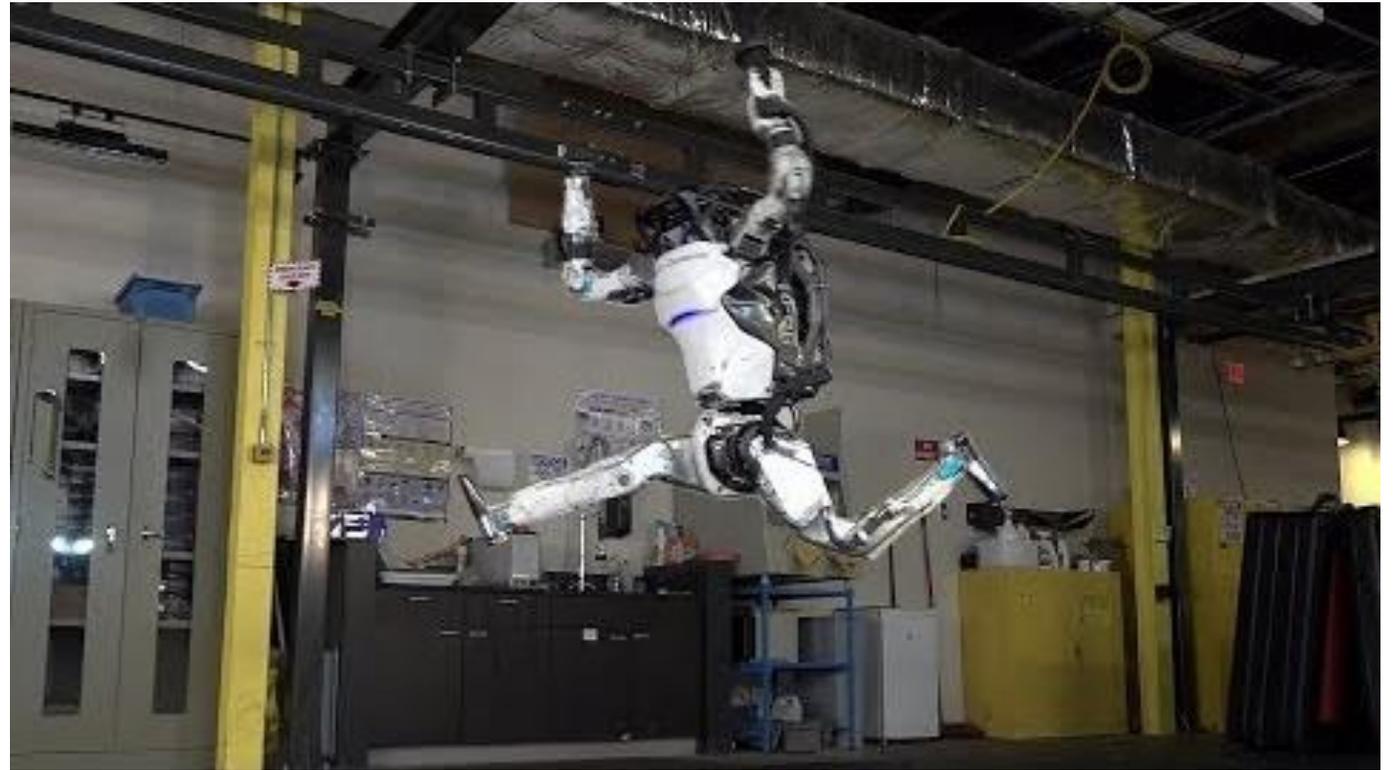
Alexa Page Views per User (Last 30 days)



Facebook Demographics



atlas







١:٢١ ص

Zain JO

"Open app store"

tap to edit

It doesn't look like you have an app named 'app store'. If you'd like, I can help you look for it on the App Store.



In terms of general intelligence,
we're not even close to a rat.

Yann LeCun
Facebook AI Director

Sources of Data-Driven Value



REDUCE COSTS



INCREASE REVENUE

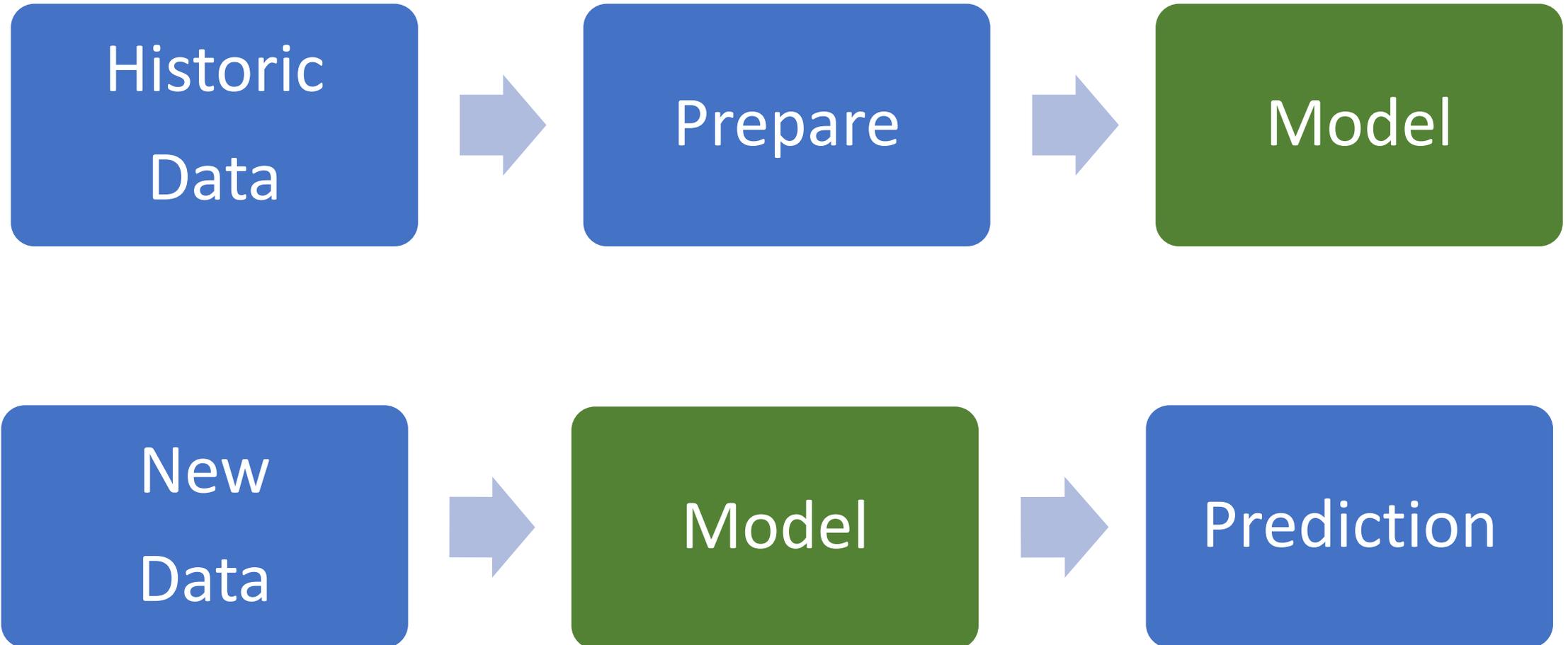


MANAGE RISK

Value from Data

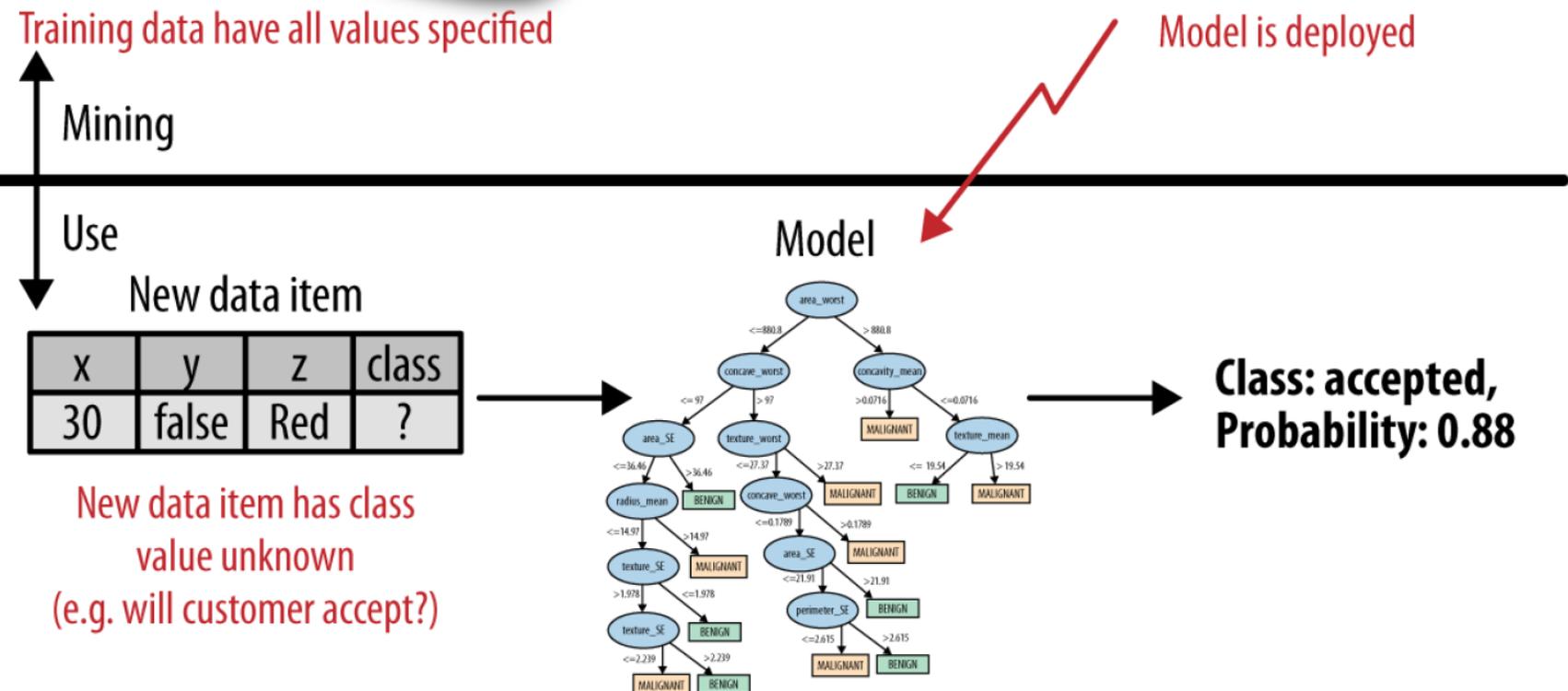
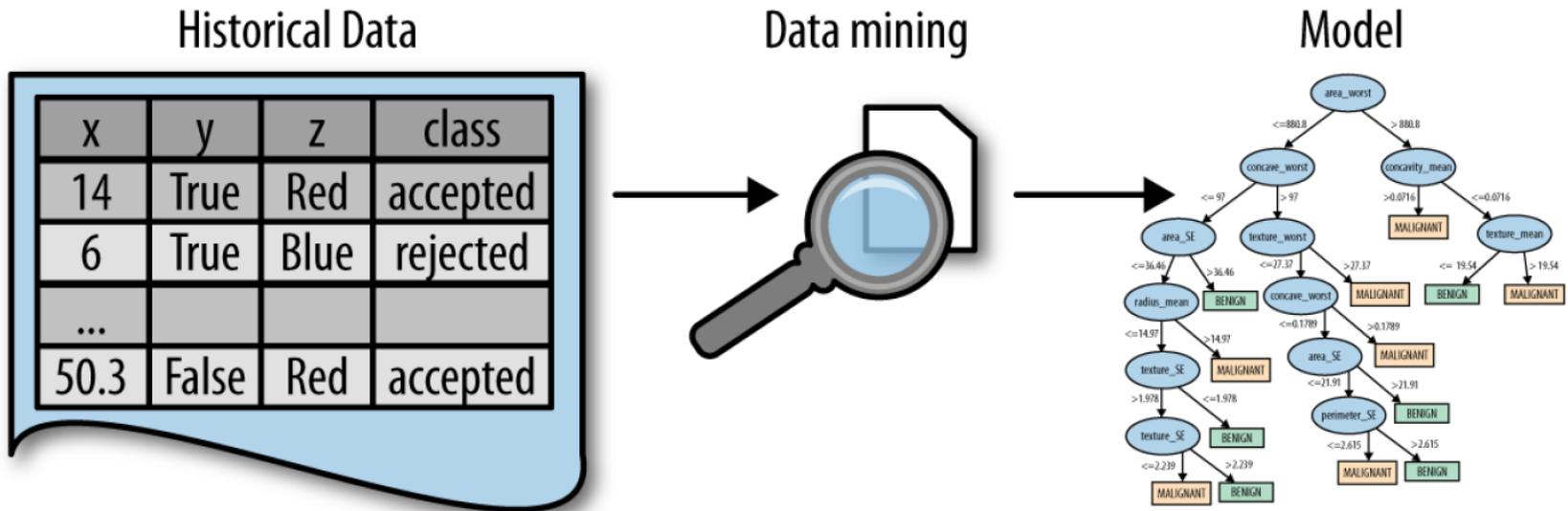


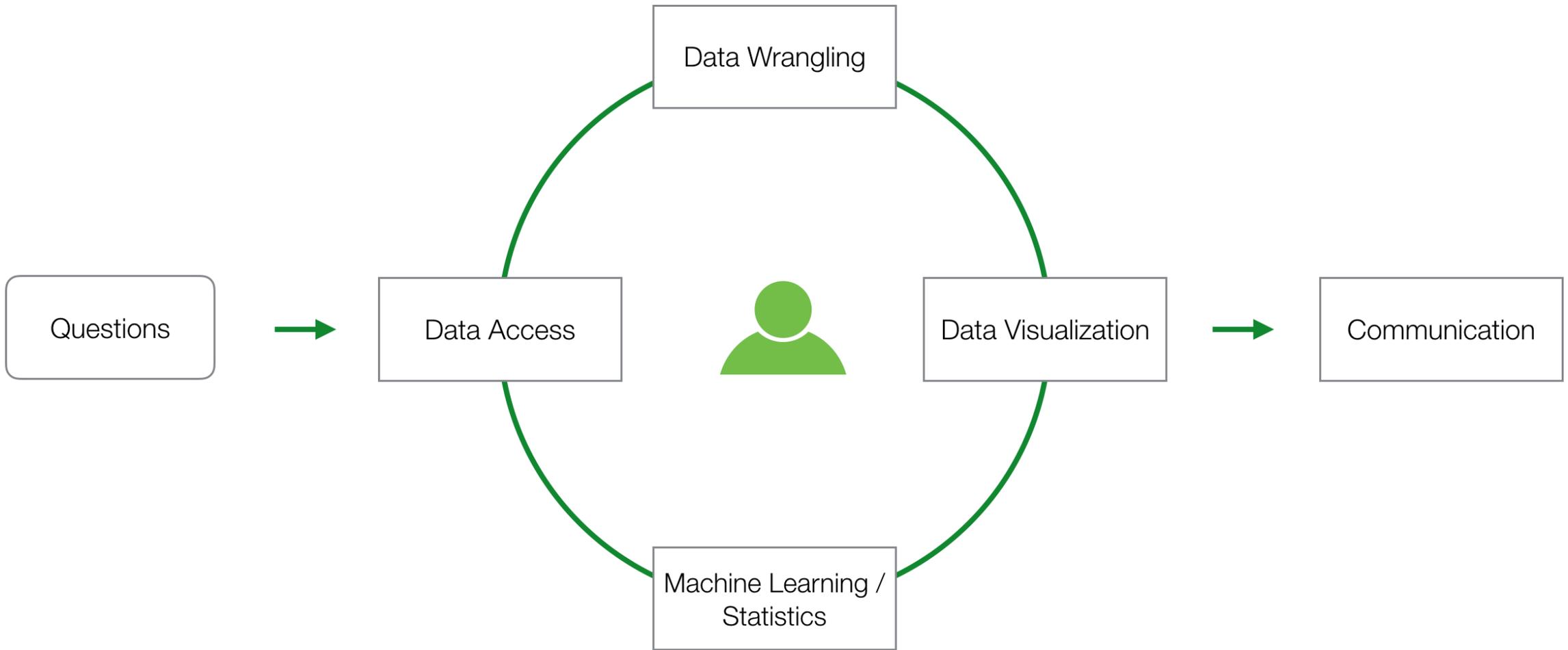
Training Predictive Models



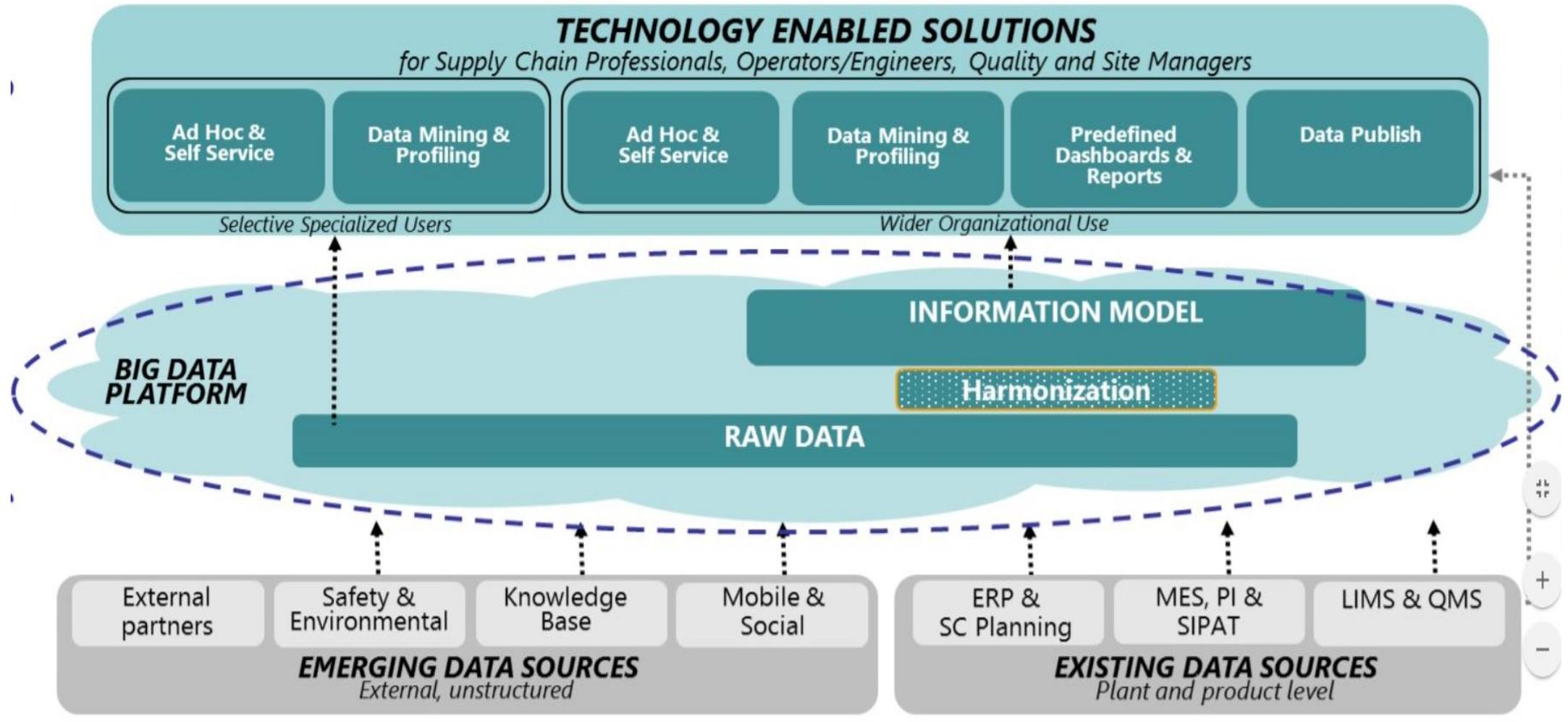
BADIR







Merck's Mantis Data Democratization



Some ML Applications

- Predictive maintenance or condition monitoring
- Warranty reserve estimation
- Propensity to buy
- Demand forecasting
- Process optimization
- Telematics

Manufacturing



- Predictive inventory planning
- Recommendation engines
- Upsell and cross-channel marketing
- Market segmentation and targeting
- Customer ROI and lifetime value

Retail



- Alerts and diagnostics from real-time patient data
- Disease identification and risk stratification
- Patient triage optimization
- Proactive health management
- Healthcare provider sentiment analysis

Healthcare and Life Sciences



- Aircraft scheduling
- Dynamic pricing
- Social media – consumer feedback and interaction analysis
- Customer complaint resolution
- Traffic patterns and congestion management

Travel and Hospitality



- Risk analytics and regulation
- Customer Segmentation
- Cross-selling and up-selling
- Sales and marketing campaign management
- Credit worthiness evaluation

Financial Services



- Power usage analytics
- Seismic data processing
- Carbon emissions and trading
- Customer-specific pricing
- Smart grid management
- Energy demand and supply optimization

Energy, Feedstock, and Utilities



ML in Healthcare

Old World

- Diagnostic care after symptoms arise
- Cancer detected at late stages
- Patients admitted to ICU suffer high mortality rate after sudden deterioration
- Drugs developed without individual considerations

New World

- Real-time preventive care before illness becomes serious
- Patient lives saved due to early ICU admission
- Cancer detected early when still treatable
- Personalized medicine developed for individual needs

ML in Banking

Old World

- Rule-based Anti-Money Laundering (AML) has high false positive rate
- Credit risk scoring relies heavily on past credit history
- Fraud detection mechanisms adapt poorly to ever-changing fraud patterns

New World

- Pattern-based AML improves accuracy and reduces false positives
- Virtual help desk improves customer experience
- Fraudulent behavior is detected faster and with better accuracy
- Product recommendations are personalized and real-time

ML in Insurance

Old World

- Product recommendation from advisor experience
- Risk group classified using static data, e.g. age and medical history
- Labor-intensive claim processes take months to finish
- Customer churn analysis after the fact
- Many loopholes in insurance fraud detection

New World

- Product recommendation from deep analytics of customer data
- Risk group classified using broad spectrum of data, e.g. social media, click streams and web analytics
- Claims processed faster using automated image classification
- Customer churns predicted and intercepted in real time
- Insurance frauds more accurately detected from customer 360° analysis
- Increased revenue

ML in Marketing

Old World

- Product recommendation from experience
- Customer churn analysis after the fact
- Decide on Marketing campaigns based on experience
- Marketing resources allocated based on needs of the past

New World

- Product recommendation from deep analytics of purchase behavior
- Customer churns predicted and intercepted in real time
- Increased revenue
- Marketing resources allocated based on needs of the future

ML in Telecom

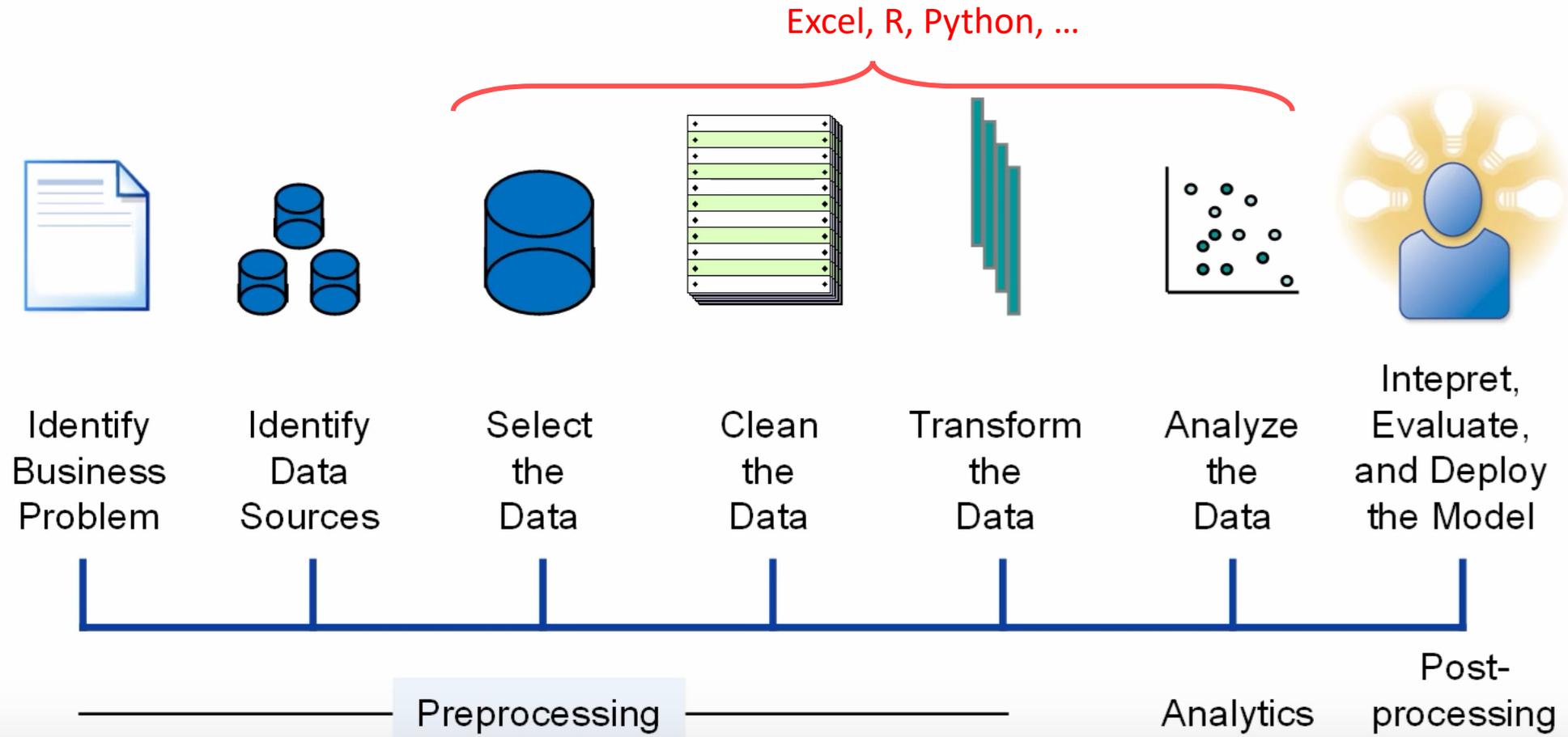
Old World

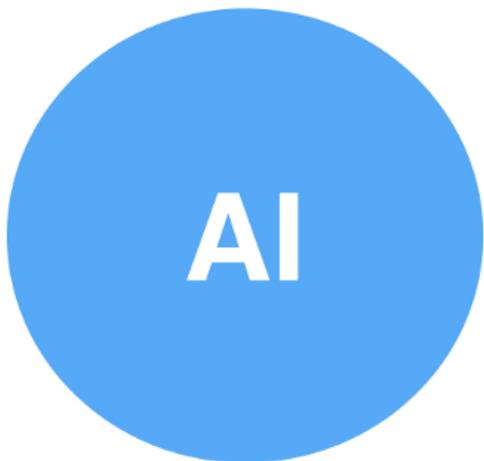
- Reactive Maintenance
- Network optimization with human intervention
- Centralized intelligence
- Security attack repair
- Backlogged customer tickets

New World

- Predictive Maintenance
- Self-optimizing network
- Optimal network quality
- Intelligence at the edge
- Security attack prediction
- Improved customer experience

The Analytics Process Model





- sounds sexy
- gets us money from VCs
- what we all hope is the future



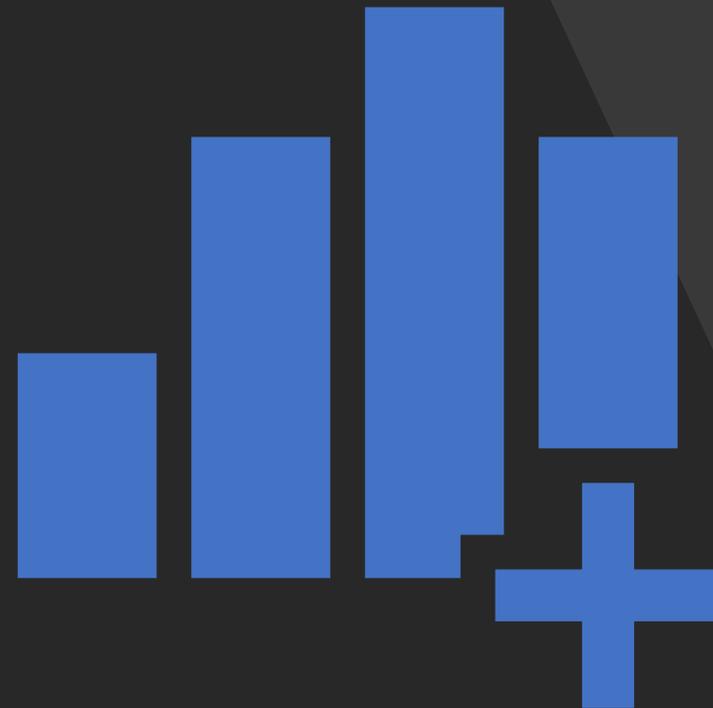
- the only real "AI"
- traditionally an academic discipline
- not concerned with real-world software



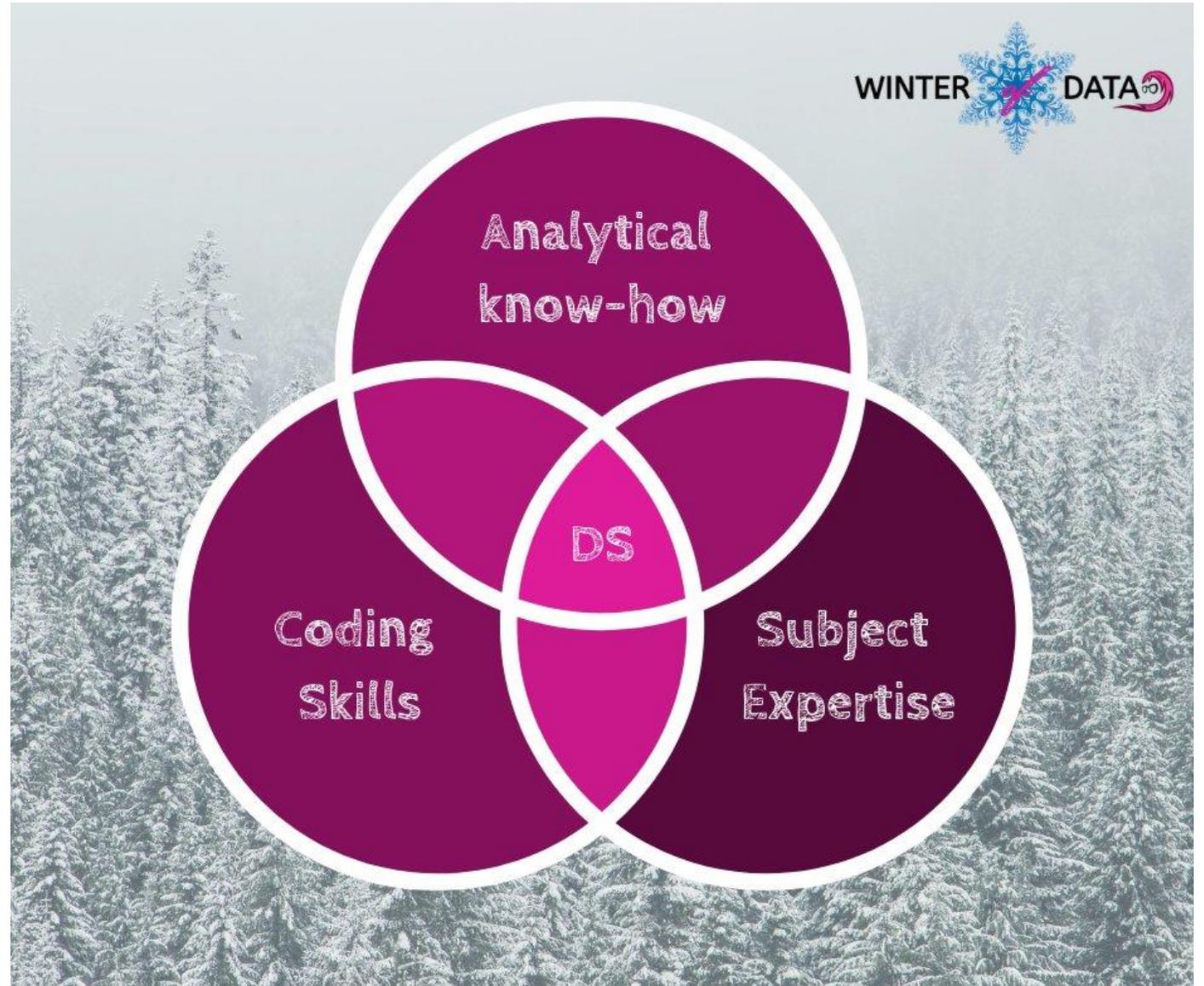
- applies machine learning to create actual products
- deals with real-world complexity

What is Data Science?

- Tools & Methods to produce business value
 - Ingest
 - Understand/Visualize
 - Wrangle/Prepare
 - Model
 - Communicate/Publish
- Support/automate business processes via:
 - Revenue Increase
 - Cost Reduction
 - Risk Management
- Deliver:
 - Dashboards/reports
 - Predictive models/APIs
 - New products



+people
+visualization
+communication







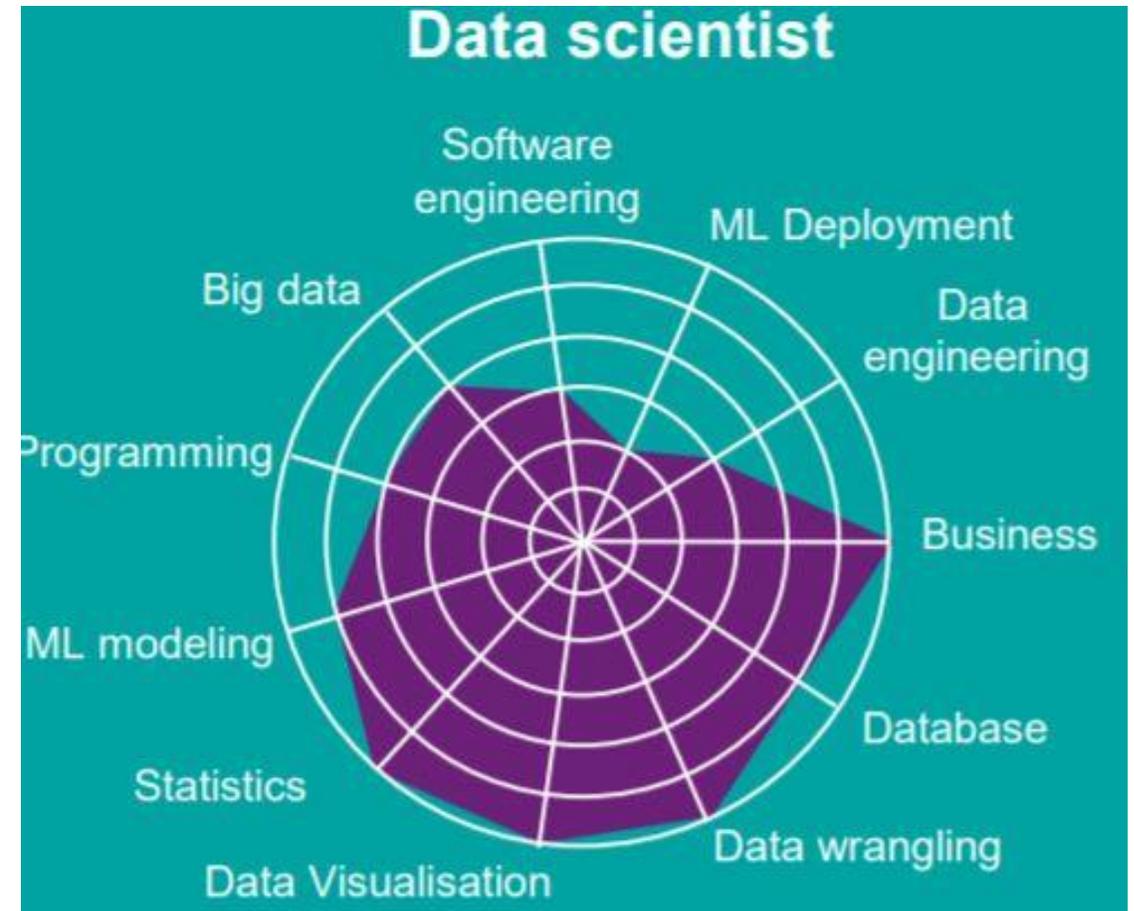
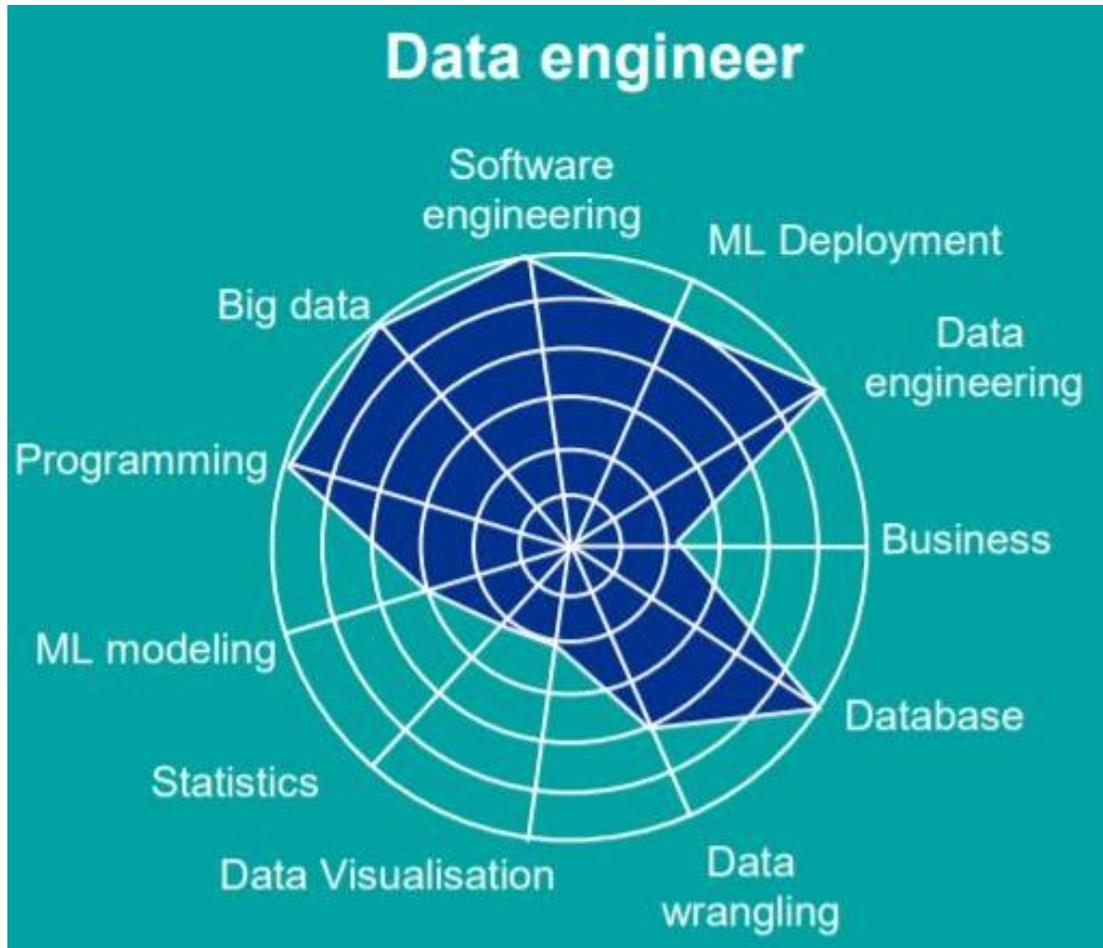
**Data
Wrangler**



**Model
Jockey**

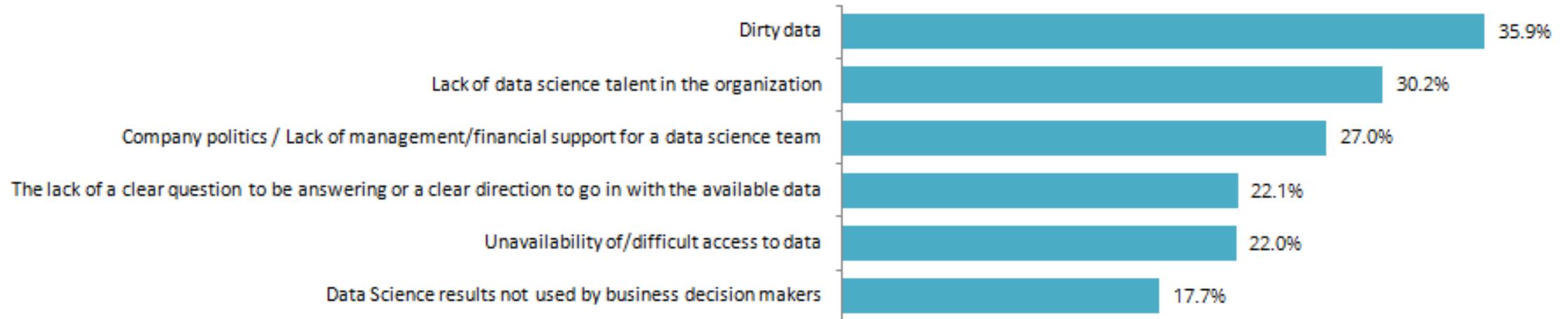


**Data
Scientist**





Challenges that Data Professionals have Faced in the Past Year





Baron Schwartz 

@xaprb

When you're fundraising, it's AI

When you're hiring, it's ML

When you're implementing, it's linear regression

When you're debugging, it's printf()

12:52 AM - Nov 15, 2017

 73

 4,550

 10,284

R Ecosystem



HISTORY AND EVOLUTION OF R



R has developed from the S language

S Version 1

S Version 2

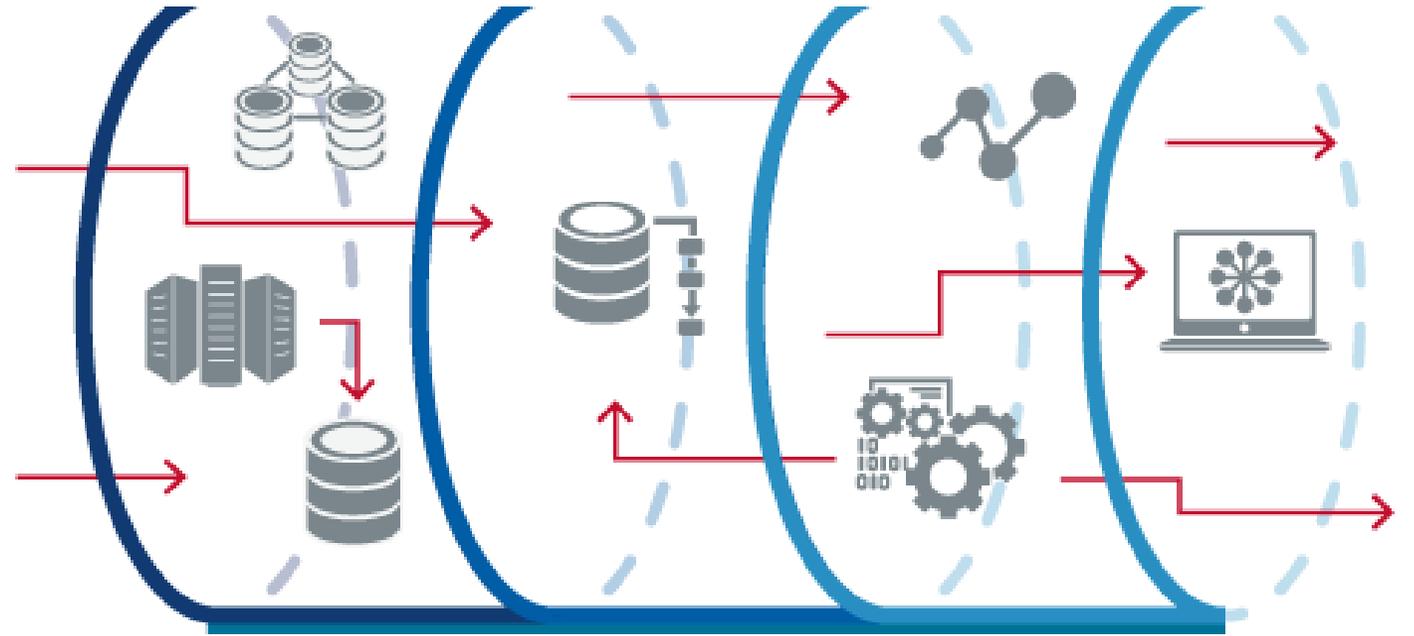
S Version 3

S Version 4

Developed 30 years ago for research
applied to the high-tech industry



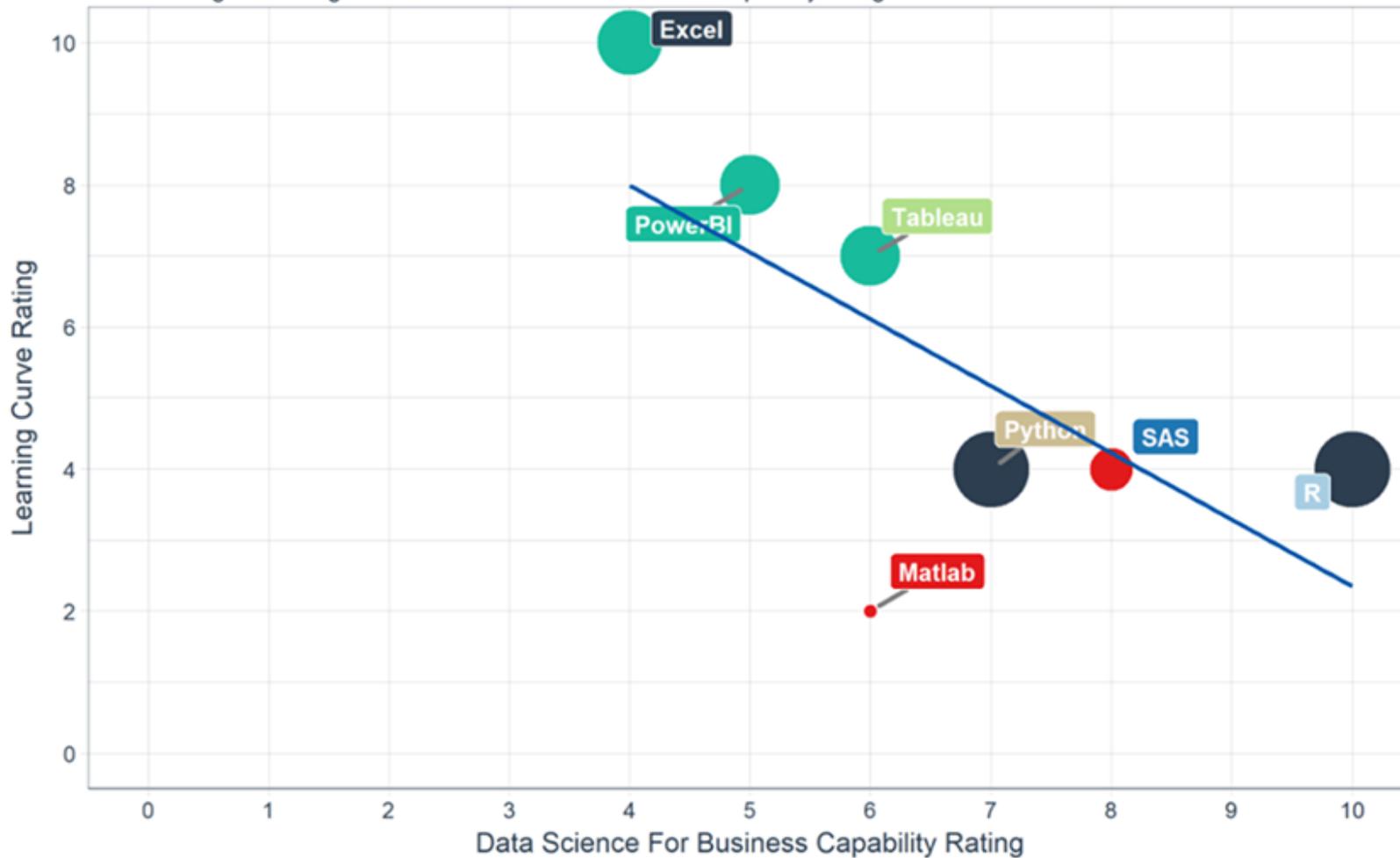




R ecosystem =
Analytics Meat Grinder

DS4B Tools: Capability Vs Learning Curve

R has a longer learning curve but has a massive business capability rating



- Cost
- Free
 - High
 - Low

- Tool
- a Excel
 - a Matlab
 - a PowerBI
 - a Python
 - a R
 - a SAS
 - a Tableau

- Trend
- 2.5
 - 5.0
 - 7.5
 - 10.0



R is the best platform for data munging or visualization

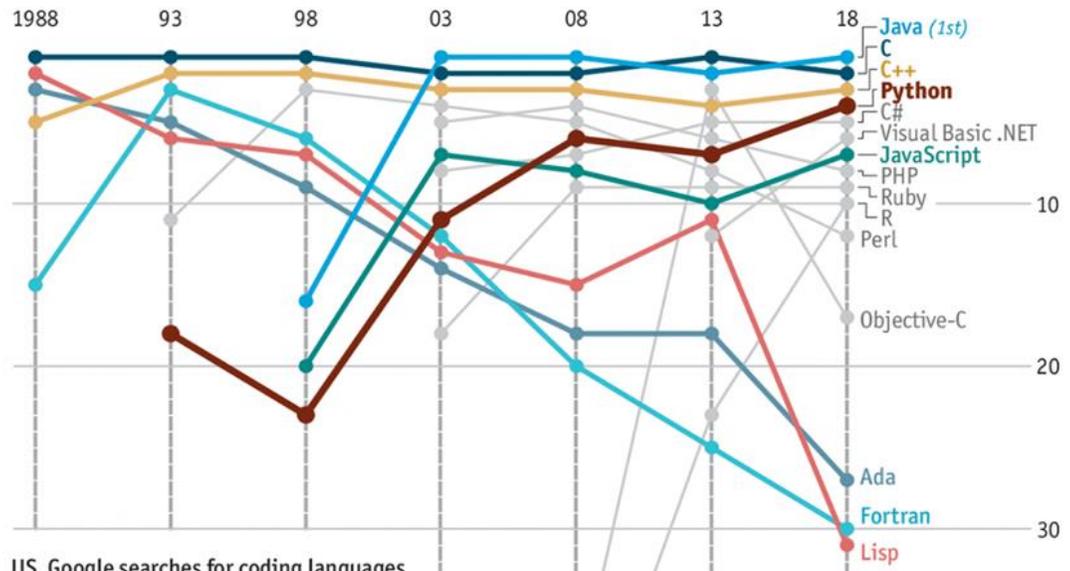
Szilard Pafka
Data Science and R Guru
Sta Monica, California



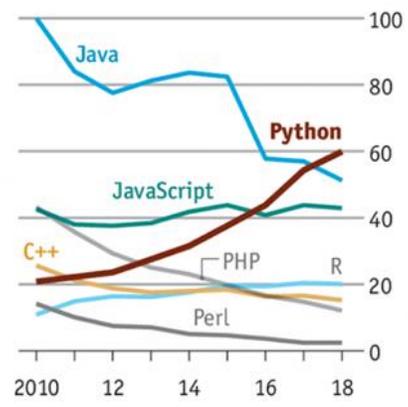
Companies that use R for Analytics



Ranking of programming languages*



US, Google searches for coding languages
100 = highest annual traffic for any language



Source: TIOBE, Google Trends

* Ranked by global search-engine popularity

The Economist

R is a “glue”

```
`` {r}
idoper <- allTablesRowsPct %>%
  filter(col=="ID_OPER") %>%
  mutate(is702=map_int(sprintf("select count(*) cnt from %s where %s=%d", tab, col, 702),
    ~sqlQuery(dbhandle, .x)$cnt),
    is749=map_int(sprintf("select count(*) cnt from %s where %s=%d", tab, col, 749),
    ~sqlQuery(dbhandle, .x)$cnt))
``
```

UNIX (points to `%>%`)

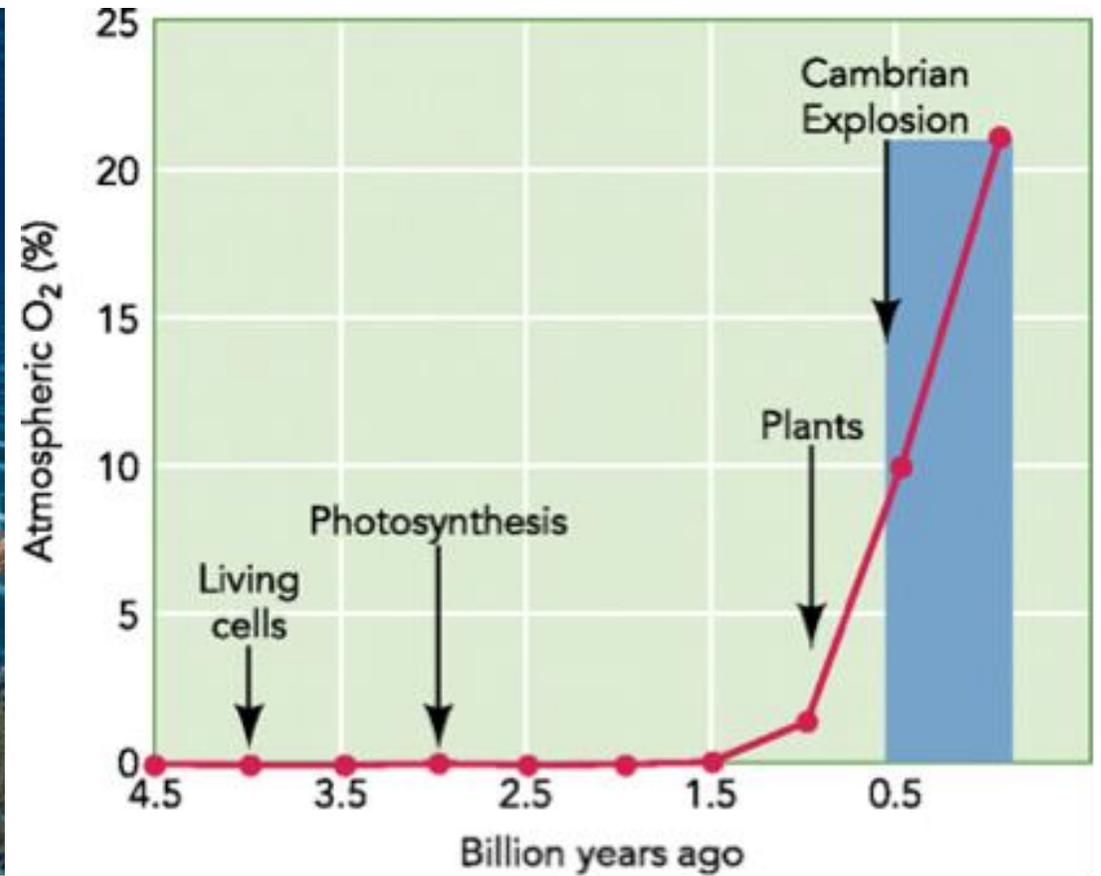
C (points to `map_int`)

SQL (points to the SQL query string)

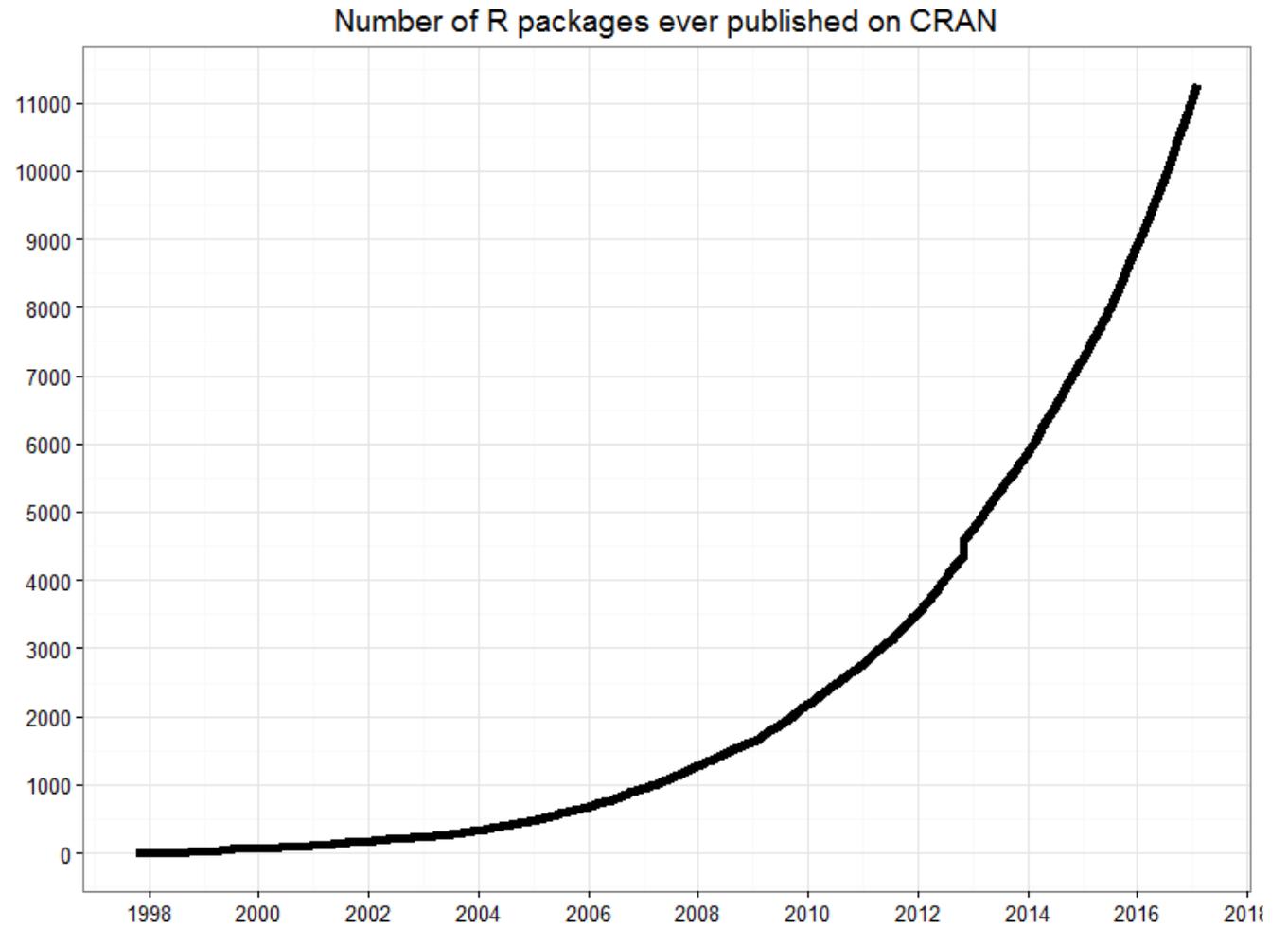
dplyr (points to `mutate`)

lisp (via purrr) (points to `map_int`)

Cambrian Explosion



Package Explosion



Leaderboard

16,989
indexed packages

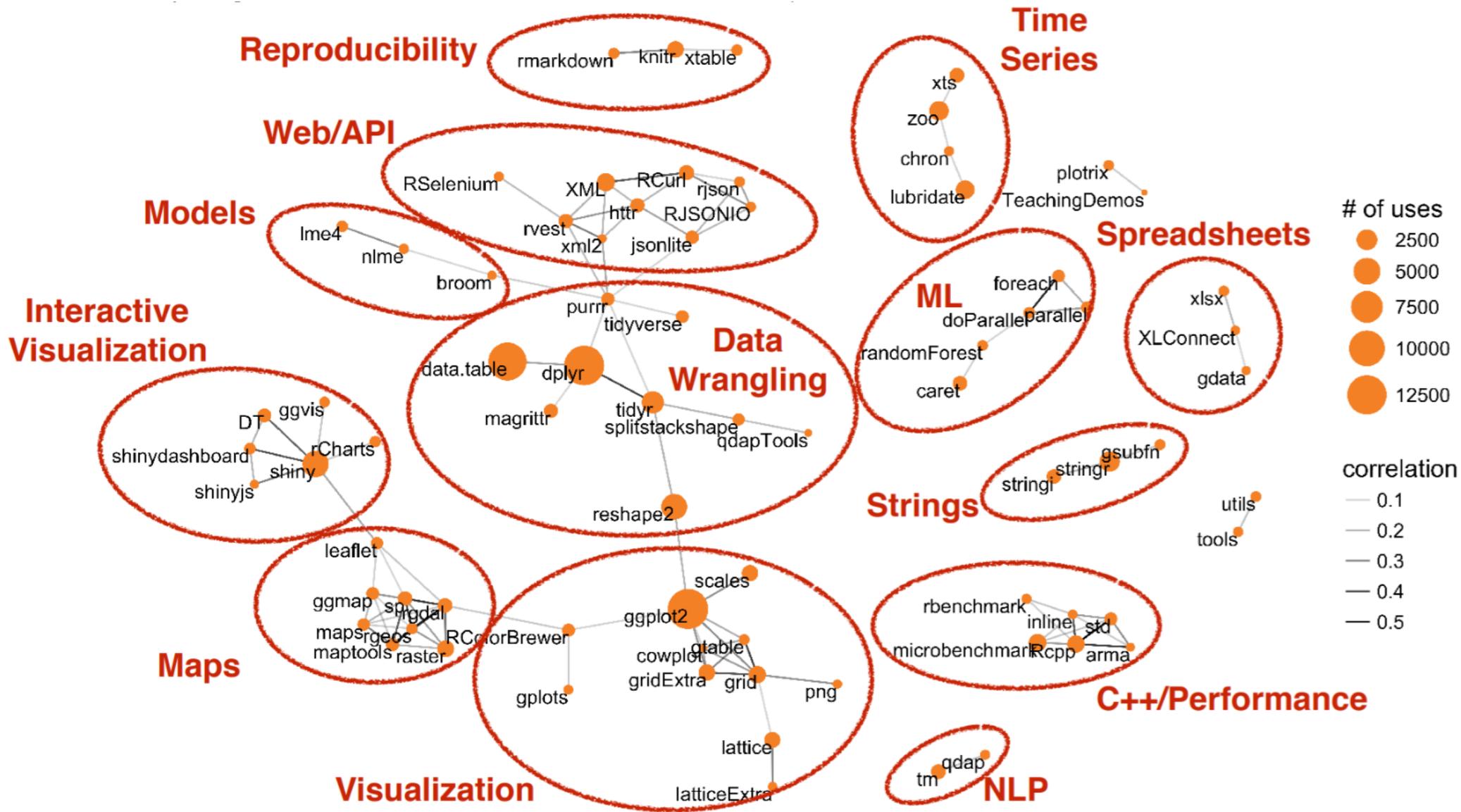
2,450,096
indexed functions

Most downloaded packages

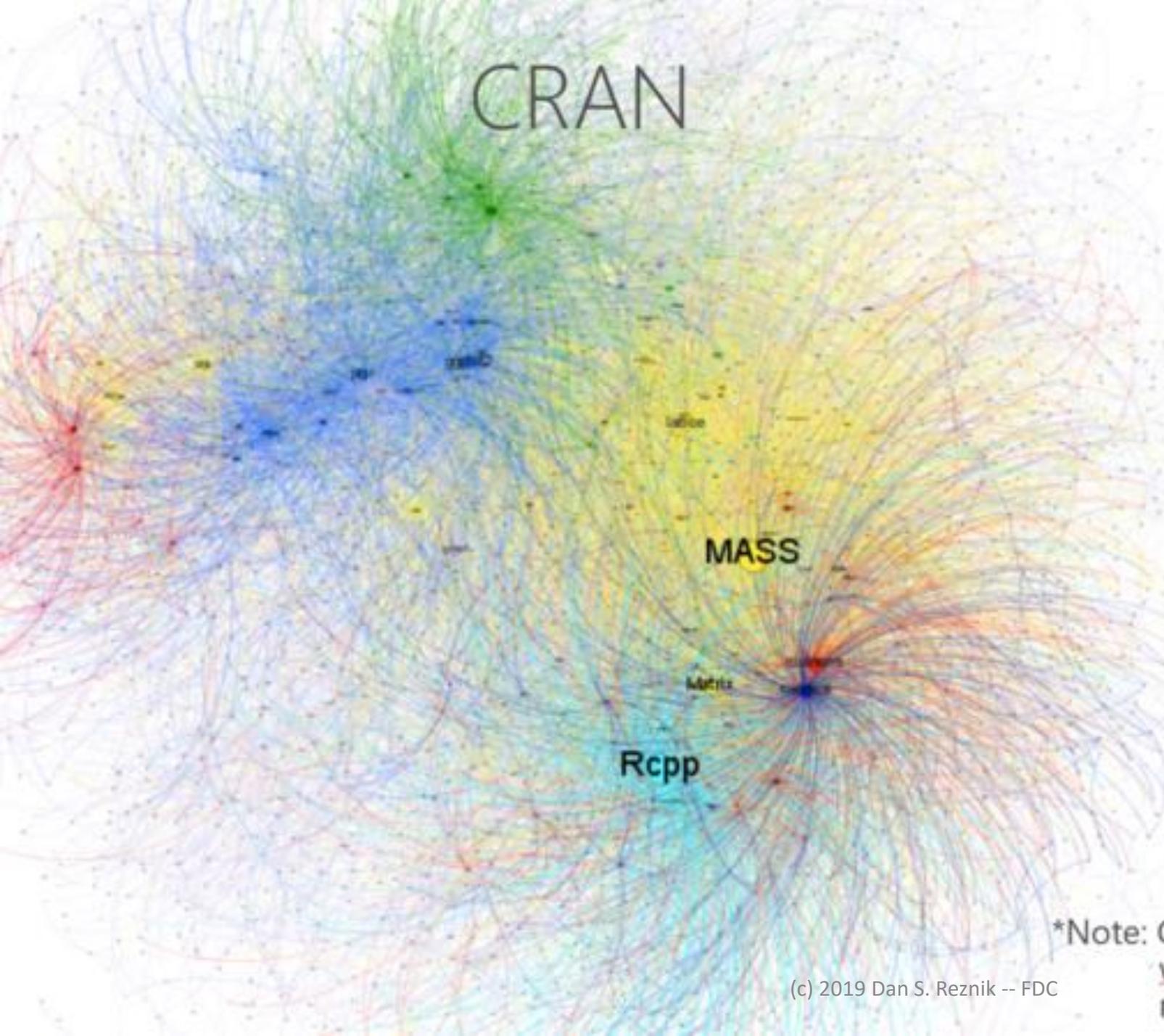
Name	Direct downloads	Indirect downloads	Total
------	------------------	--------------------	-------

Most active maintainers

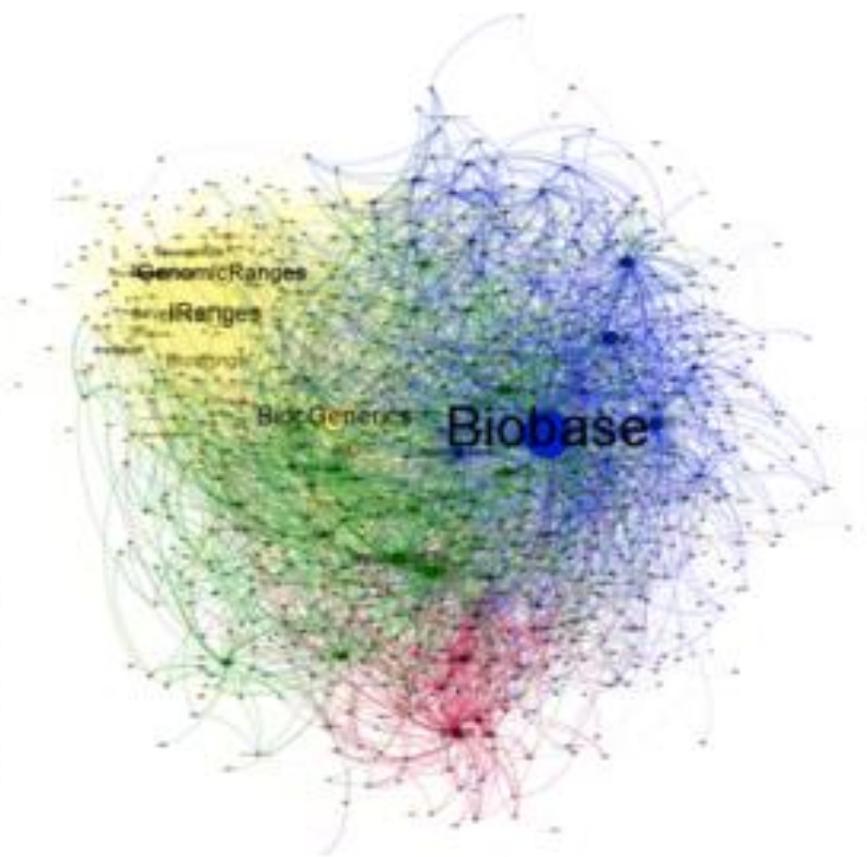
Name	Direct downloads	Indirect downloads	Total
------	------------------	--------------------	-------



CRAN



BioConductor



*Note: Colour indicates communities found by the walktrap algorithm, but has no common meaning in the two networks

Import



Tidy



Transform



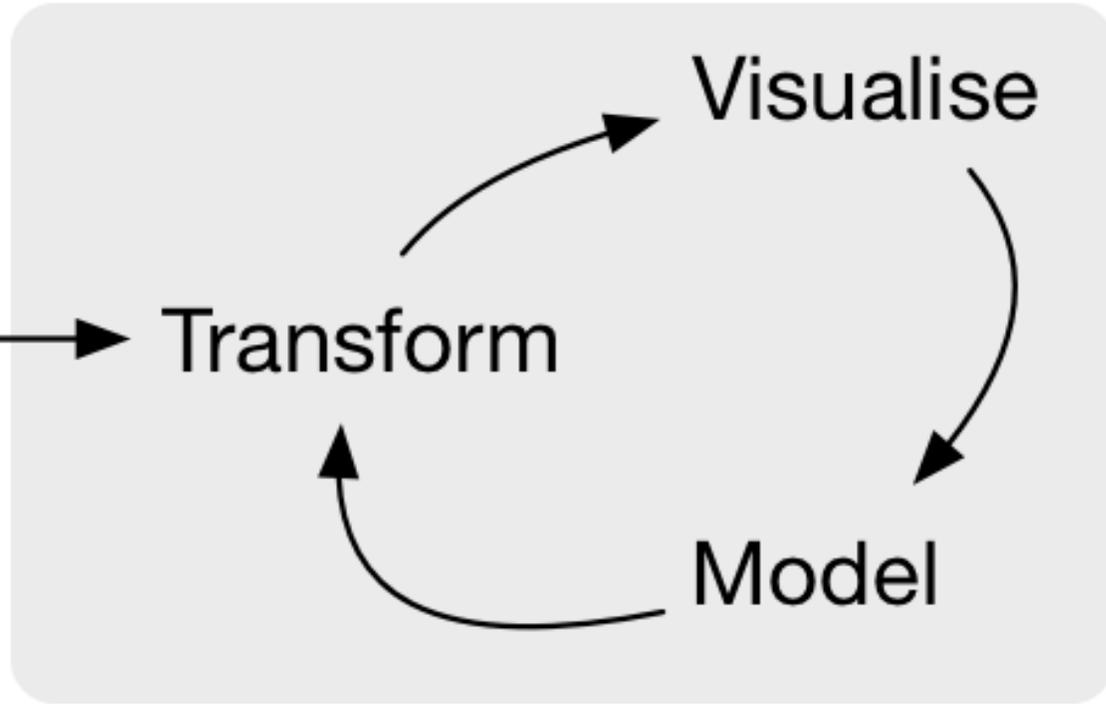
Model



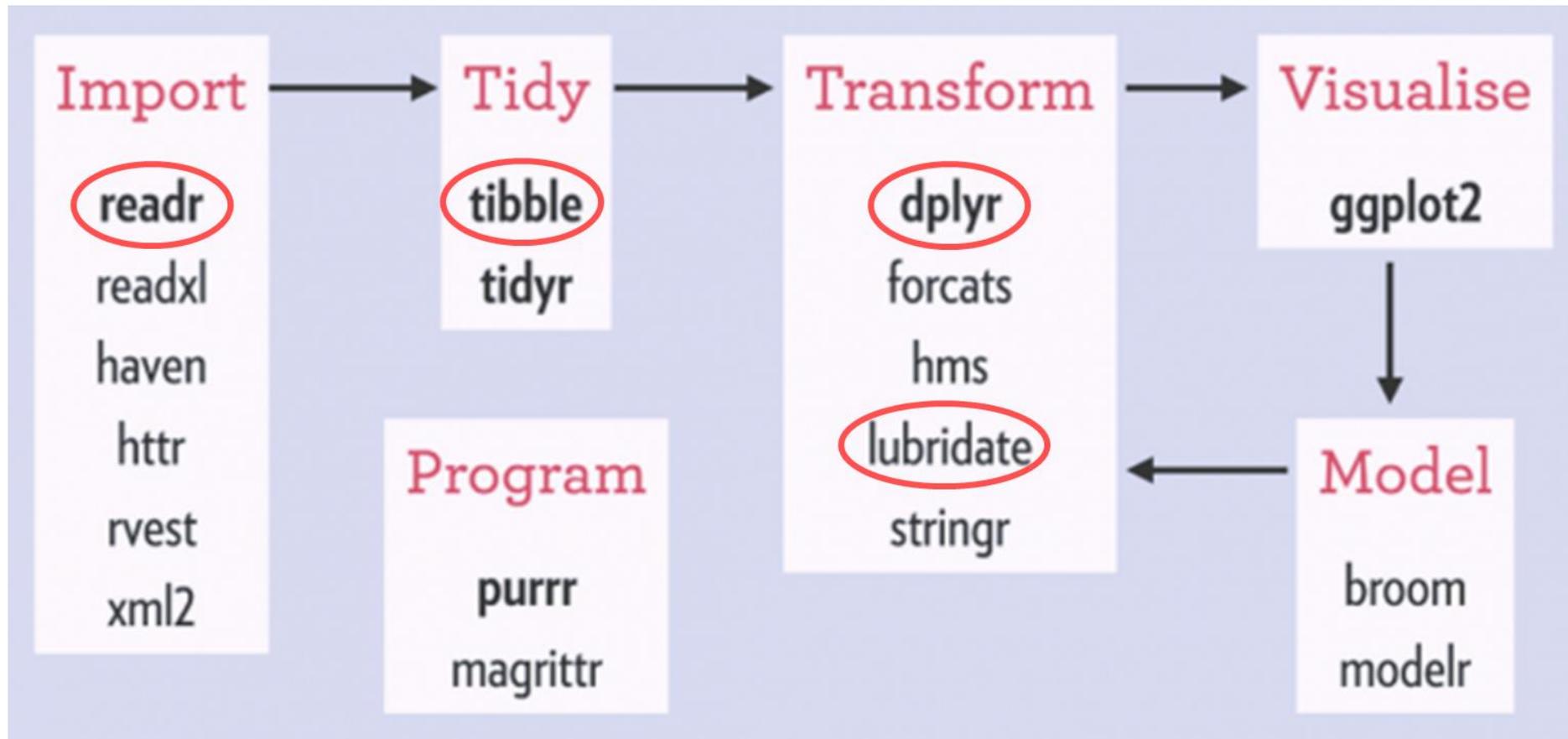
Visualise



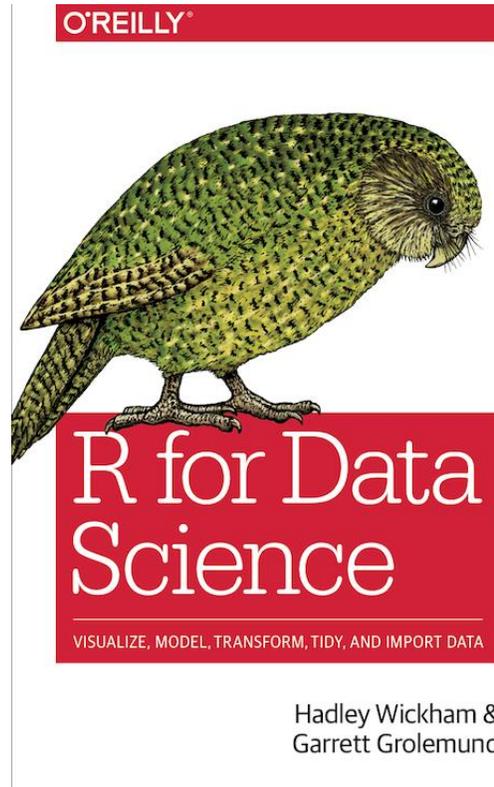
Communicate



Understand



The tidyverse: packages for the data science workflow



<https://r4ds.had.co.nz/>

RStudio

The screenshot displays the RStudio interface with the following components:

- Source Editor:** Shows the R script `arrangements.R` with the following code:

```
6 #' @useDynLib arrangements
7 "_PACKAGE"
8
9 Arrangements <- R6::R6Class(
10   c("Arrangements", "iter", "abstractiter"),
11   private = list(
12     state = NULL
13   ),
14   public = list(
15     nextElem = function() {
16       out <- self$getNext()
17       is.null(out) && stop("StopIteration", call. = FALSE)
18       out
19     }
20   )
21 )
```
- Console:** Shows the R startup message:

```
~/Dropbox/R/arrangements/
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

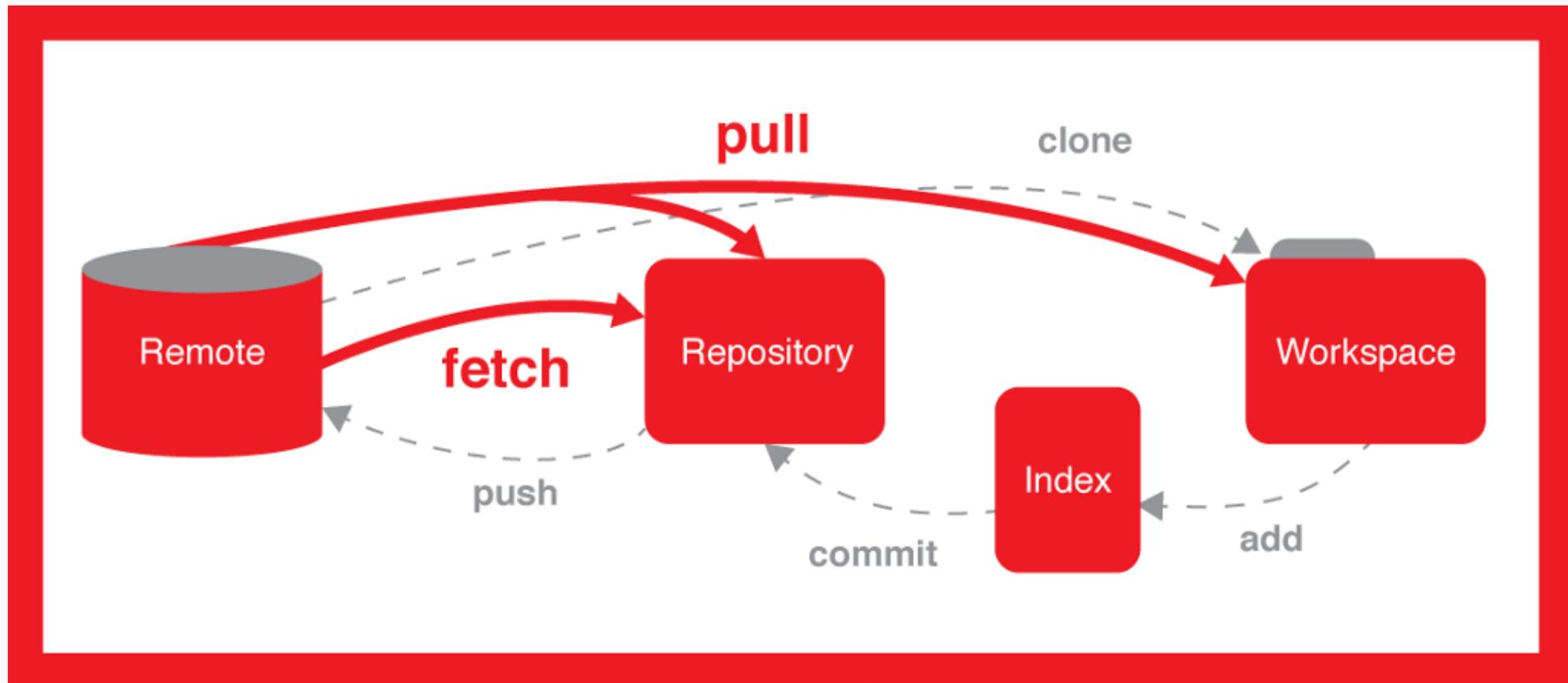
R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> |
```
- Environment Panel:** Shows a table with columns "Connection" and "Status".
- Files Panel:** Shows a file browser view of the directory `~/Dropbox/R/arrangements/R` with the following files:

Name	Size	Modified
..		
arrangements.R	465 B	Aug 9, 2018, 2:32 PM
combinations.R	5.4 KB	Aug 13, 2018, 10:59 PM
partitions.R	4.8 KB	Aug 13, 2018, 11:22 PM
permutations.R	5.7 KB	Aug 13, 2018, 11:13 PM
utils.R	783 B	Aug 23, 2018, 1:48 AM

GitHub



Install R

R-3.6.1 for Windows (32/64 bit)

[Download R 3.6.1 for Windows](#) (81 megabytes, 32/64 bit)

[Installation and other instructions](#)

[New features in this version](#)

If you want to double-check that the package you have downloaded matches the package distributed by CR to the [fingerprint](#) on the master server. You will need a version of md5sum for windows: both [graphical](#) and

Frequently asked questions

- [Does R run under my version of Windows?](#)
- [How do I update packages in my previous version of R?](#)
- [Should I run 32-bit or 64-bit R?](#)

arrangements.R x

```

6 #' @useDynLib arrangements
7 "_PACKAGE"
8
9 Arrangements <- R6::R6Class(
10   c("Arrangements", "iter", "abstractiter"),
11   private = list(
12     state = NULL
13   ),
14   public = list(
15     nextElem = function() {
16       out <- self$getNext()
17       is.null(out) && stop("StopIteration", call. = FALSE)
18       out
19     }

```

16:26 nextElem()

R Script

Console

Terminal x

Jobs x

~/Dropbox/R/arrangements/

Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> |

Environment

History

Connections

Build

Git

New Connection

Connection

Status

Files

Plots

Packages

Help

Viewer

New Folder

Delete

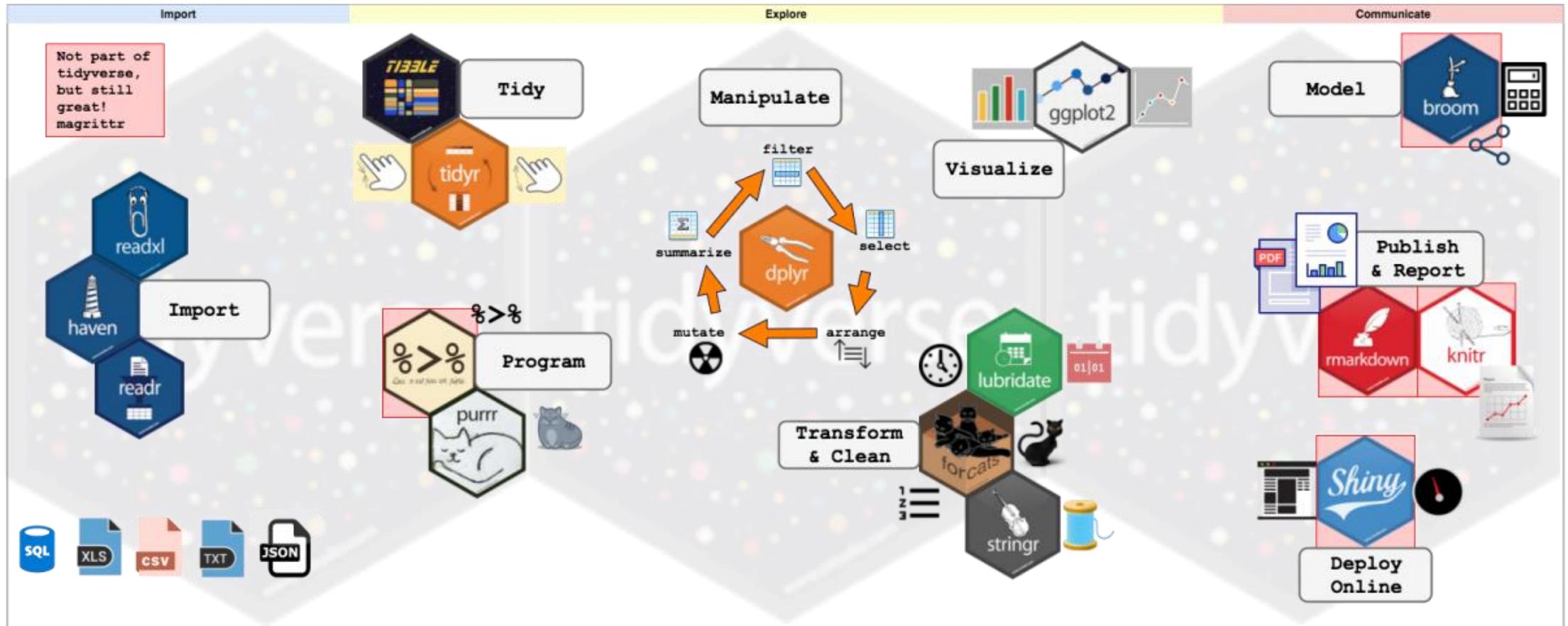
Rename

More

Home > Dropbox > R > arrangements > R

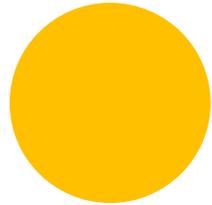
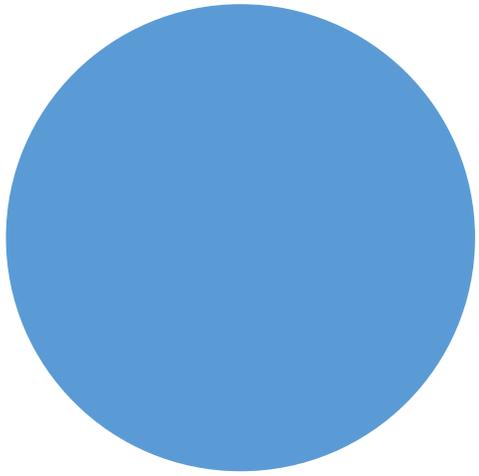
	Name	Size	Modified
	..		
	arrangements.R	465 B	Aug 9, 2018, 2:32 PM
	combinations.R	5.4 KB	Aug 13, 2018, 10:59 PM
	partitions.R	4.8 KB	Aug 13, 2018, 11:22 PM
	permutations.R	5.7 KB	Aug 13, 2018, 11:13 PM
	utils.R	783 B	Aug 23, 2018, 1:48 AM

R DS Pipeline



Takeaways

- Business Impact
 - Data Science
 - Machine Learning
- Tools for Working the Data
 - Cloud
 - GitHub
 - R Ecosystem
- Business Problem First, Technology Later



Thank you!

dan@dat-sci.com

www.dat-sci.com

Web: Data Science

- [book] Strategic Principles of DS <https://lucidmanager.org/strategic-data-science/>
- Towards Data Science: <https://towardsdatascience.com/>
- Automation Tools for DS: <https://www.kdnuggets.com/2018/12/automation-data-science.html>
- Hadley on how to become a DS: https://gist.github.com/hadley/820f09ded347c62c2864?fbclid=IwAR22uF47RlgAmC0ZcErb4MuVRVdHBQZKEn9yxn8S5tOBppxE_Onfjt5iQ
- Intro to GitHub: <https://medium.com/@abhishekj/an-intro-to-git-and-github-1a0e2c7e3a2f>
- Command-Line Data Science: <https://adamdrake.com/command-line-tools-can-be-235x-faster-than-your-hadoop-cluster.html>
- WEF: DS in the new economy: http://www3.weforum.org/docs/WEF_Data_Science_In_the_New_Economy.pdf
- Data Scientist vs Engineer vs Analyst: <https://news.efinancialcareers.com/uk-en/3001517/data-science-careers-finance>
- R for Data Science: <https://r4ds.had.co.nz/>

Web: ML/AI

- Data Science vs ML: <https://www.kdnuggets.com/2018/12/learning-machine-learning-data-science.html>
- Davenport: Shift Toward AI: <https://www.forbes.com/sites/tomdavenport/2019/04/05/ai-is-destroying-traditional-business-thinking>
- Amazing Innovation w/ Self-Driving Cars: <https://www.cbinsights.com/research/startups-drive-auto-industry-disruption>
- All ML Algos in one page: <https://www.r-bloggers.com/101-machine-learning-algorithms-for-data-science-with-cheat-sheets/>
- State of ML 2019: <https://www.analyticsinsight.net/the-state-of-machine-learning-in-2019/>
- ML Process: <https://www.kdnuggets.com/2018/12/essence-machine-learning.html>
- ML with R: <https://bradleyboehmke.github.io/HOML/>
- The AI “Clown Show”: <https://blog.pieknewski.info/2019/05/30/ai-circus-mid-2019-update/>
- How to recognize AI “snakeoil”: <https://www.cs.princeton.edu/~arvindn/talks/MIT-STS-AI-snakeoil.pdf>

Web: R

- Is it Worth Learning R: <https://www.quora.com/Is-it-worth-learning-R-programming/answer/Aaron-Kumar-35>
- Ditch the Spreadsheet: <https://lucidmanager.org/spreadsheets-for-data-science>
- Data Exploration Package: <https://boxuancui.github.io/DataExplorer/>
- Using GitHub with R: <https://happygitwithr.com/>
- Feature Engineering with the Tidyverse: <https://bookdown.org/max/FES/>
- Programming with the Tidyverse: <https://speakerdeck.com/lionelhenry/programming-in-the-tidyverse>
- R tutorials: <https://teachingr.com/>
- EDA w/ R: <https://www.r-bloggers.com/part-2-simple-eda-in-r-with-inspectdf/>
- From Excel to R: <https://rfortherestofus.com/2019/06/a-guide-to-r-for-excel-users/>

Web: Viz

- Data Viz with R: <https://socviz.co/>
- Geocomputation with R: <https://geocompr.robinlovelace.net/>
- Gallery of ggplot2: <https://www.r-graph-gallery.com/ggplot2-package.html>
- How to pick the right plot: <https://www.data-to-viz.com/>
- Shiny Contest Winners: <https://blog.rstudio.com/2019/04/05/first-shiny-contest-winners>
- Shiny+Plotly: https://plotly-r.com/plotly_book.pdf
- Ggplot and Gganimate: <https://djenavarro.github.io/satrdajoburg/slides/>



Videos: Data Science

- General:
 - Just Build Something w/ Data: <https://youtu.be/MOdlp1d0PNA>
- Viz:
 - Tour of Data Viz: <https://youtu.be/9F1pZ4X8xPE>
 - Chess Players over Time: <https://youtu.be/z2DHPw79w0Y>

Videos: R

- Gentle Intro to Stats w/ R/Tidyverse: <https://resources.rstudio.com/webinars/a-gentle-introduction-to-tidy-statistics-in-r>
- R then and now: <https://youtu.be/c075aRNmkUk>
- Wrangling w/ R: <https://youtu.be/DwWH1mTerOc>
- R vs Excel: <https://youtu.be/S8wO4ppE9L8>
- Deploying R APIs w/ Docker: <https://youtu.be/cO3QvioMESo>
- Six Reasons R is Ideal for DS: <https://youtu.be/D4LPLS9yCmc>