

R Packages for Team Collaboration

Stephanie Kirmer, Data Scientist

@data_stephanie
www.stephaniekirmer.com
<https://github.com/skirmer/odsc2018/>

About Me and Uptake

My role includes: modeling, internal tool development, exploratory analysis, customer interaction

Uptake:

- Analytics company serving a variety of areas of industry (rail, mining, oil & gas, trucking, etc)
- We provide data science and analytics insights to help reduce downtime, improve safety, and increase efficiency for our customers.
- 60+ person data science team, with expertise across the field
- Data science is at the core of what our company does!

With a team this big, we have strong techniques for collaboration and sharing knowledge.

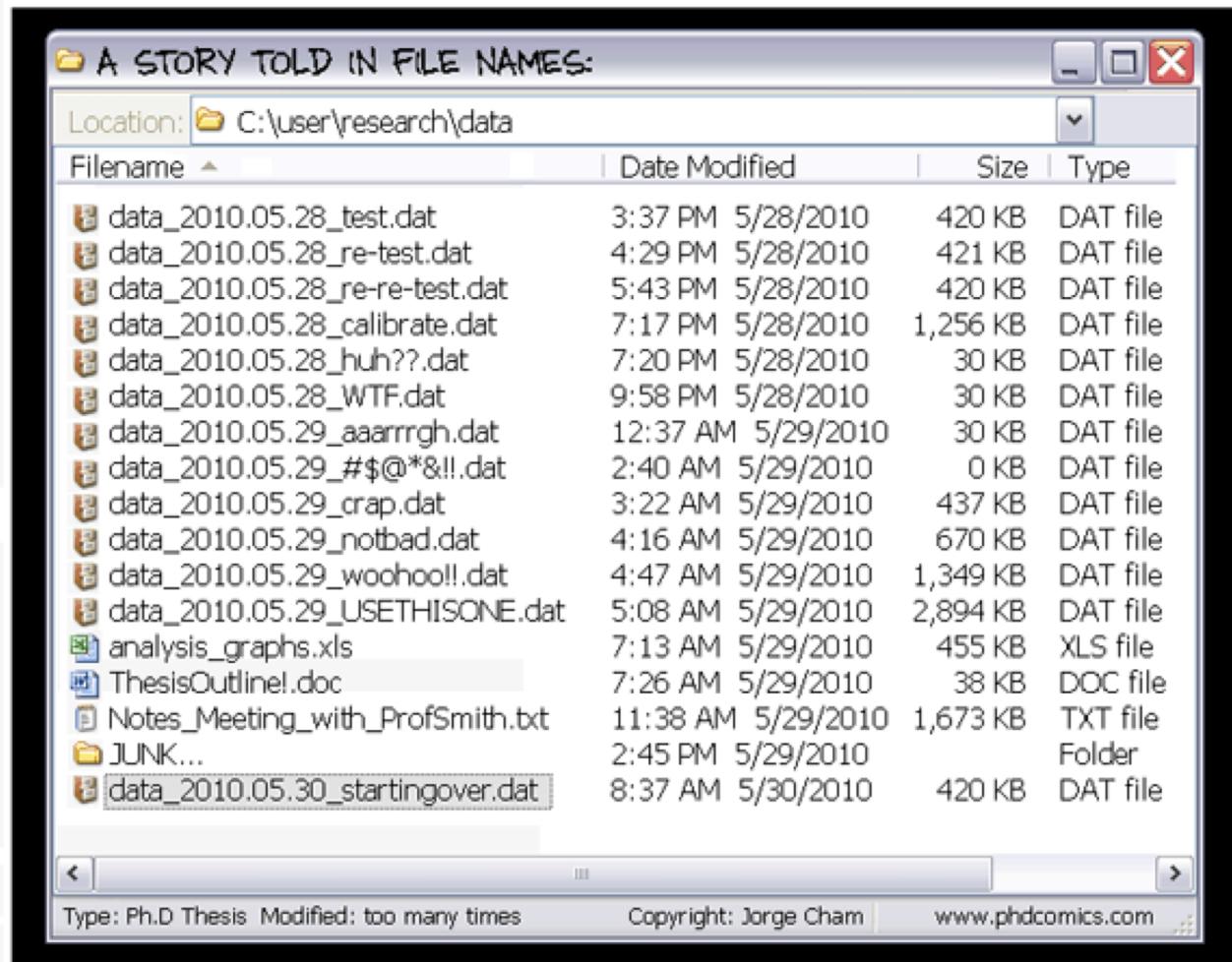


Problem:

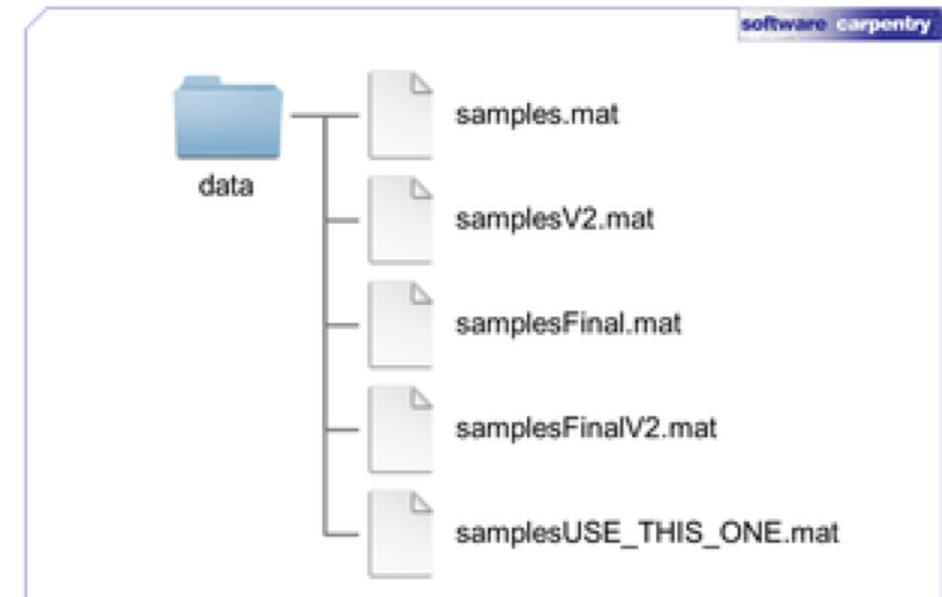
Data science teams without structured, intentional collaboration leak knowledge and waste resources.



Don't do this.



<http://phdcomics.com/comics/archive.php?comicid=1323>



<http://v4.software-carpentry.org/data/mgmt.html>

**Not in your personal life.
Definitely not in business.**

Solution:

**Build your team's culture,
infrastructure, and tools so that
retaining and sharing knowledge
is easy and valued.**

**R packages are a great tool to
help you get there.**



Do this instead.

The screenshot shows the GitHub organization page for `UptakeOpenSource`. The top navigation bar includes links for `Repositories 3`, `People 17`, `Teams 0`, and `Projects 0`. Below the navigation is a search bar with placeholder text `Search repositories...` and filters for `Type: All` and `Language: All`, along with a green `New` button. The main content area displays two repository cards:

pkgnet
R package for analyzing other R packages via graph representations of their dependencies

R dependency-analysis dependency-graph graph-theory
R ★ 16 8 1 issue needs help Updated 21 hours ago

uptasticsearch
An Elasticsearch client tailored to data science workflows.

elasticsearch data-science r etl nosql document-database
R ★ 22 16 4 issues need help Updated a day ago

On the right side of the page, there are two boxes: **Top languages** (R and Python) and **People** (17), which lists 17 user profiles.

<https://github.com/UptakeOpenSource>

Version control

Collaborative development

Packaged tools

 **jameslamb** Merge pull request #55 from UptakeOpenSource/feature/release_v0.2.0 ...

 R	Made it a merge instead of rep
 docs	Moved R package up to repo root
 inst/testdata	Moved R package up to repo root
 man	renamed doc_null to doc_shared
 tests	Made it a merge instead of rep
 .Rbuildignore	speed up unpack_nested_data #42
 .gitignore	Deleted "*.html" from gitignore since we'll need to track the docs
 .travis.yml	Moved R package up to repo root
 CONDUCT.md	Updated CONDUCT.md
 DESCRIPTION	Added release v0.2.0
 LICENSE	Moved R package up to repo root
 NAMESPACE	Remove handling of mixed atomic/data.table columns
 NEWS.md	Added release v0.2.0
 README.md	Moved R package up to repo root
 _pkgdown.yml	Moved R package up to repo root
 cran-comments.md	Added release v0.2.0

<https://github.com/UptakeOpenSource/uptasticsearch>

Code can all be installed with a single command.

Versions are automatically managed and are always accessible.

Folders have standard organizing scheme.



Progression of Team Collaboration Infrastructure

None	Little	Moderate	High
No version control systems.	Minimal, unsystematic use of version control.	Git is in place, but used inconsistently.	Git is an integral part of the workflow for the entire team.
Code sharing via email, when it is shared at all.	Code sharing via email or links/filepaths on a shared file space.	Code sharing is done with git when it happens.	Code is shared & reviewed through git, and continuous integration (CI) system is used.
Code stored locally on hard drives, not accessible to colleagues.	Code stored on shared drive space or a cloud service.	Git is the ideal place to share code, but shared drives are still used.	Packages are shared through an internal CRAN-type system.



Progression of Team Collaboration Infrastructure

None	Little	Moderate	High
No version control systems.	Minimal, unsystematic use of version control.	Git is in place, but used inconsistently.	Git is an integral part of the workflow for the entire team
Code sharing via email, when it is shared at all.	Code sharing via email or links/filepaths on a shared file space.	Code sharing with git happens	<p>This is where using team R packages becomes a real possibility</p> <p>it, item</p>
Code stored locally on hard drives, not accessible to colleagues.	Code stored on shared drive space or a cloud service.	Git is the ideal place to share code, but shared drives are still used.	Packages are shared through an internal CRAN-type system.



Why?

Improve quality and reproducibility by using peer reviewed, validated tools in package form.

- Rewriting code = opportunity for mistakes.
- A shared set of tools creates consistency across projects.
 - Definitions of measures and metrics
 - Data cleaning and management steps

Tutorial

Curating Research Assets: A Tutorial on the Git Version Control System



Matti Vuorre¹ and James P. Curley^{1,2}

¹Department of Psychology, Columbia University, and ²Department of Psychology, University of Texas at Austin

<http://journals.sagepub.com/eprint/CxmDEDrWCmqw6sZTGNEh/full>

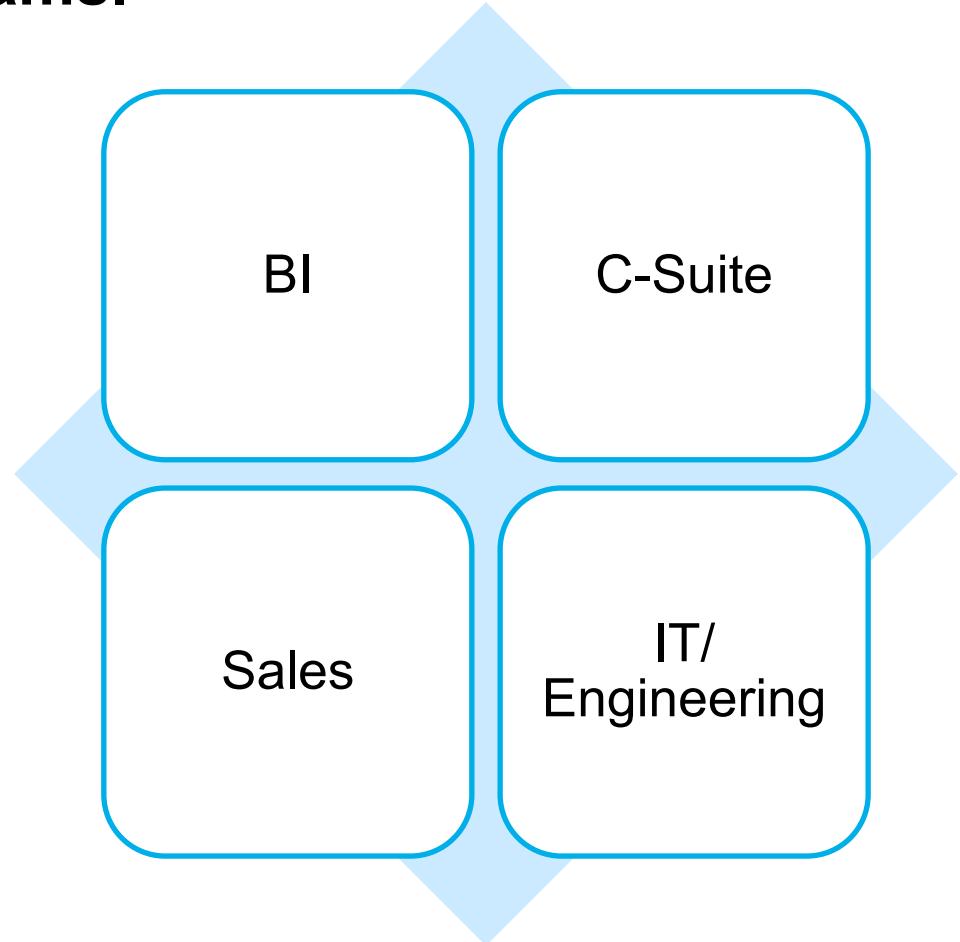


Advances in Methods and Practices in Psychological Science
1–18
© The Author(s) 2018
Reprints and permissions:
sagepub.com/journalsPermissions.nav
DOI: 10.1177/2515245918754826
www.psychologicalscience.org/AMPPS



Build package elements for use by other teams.

- Package functionalities can include:
 - Data ETL
 - Basic analysis
 - Visualizations
 - .. and more!
- Teams can also build useful tools for their specific tasks



Empower your colleagues to do routine work, and free up data science to do the complex and innovative projects.

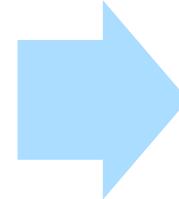


Don't wait until personnel changes to think about your information.

With robust information infrastructure:

New Hire

- Goes from beginner to contributor – fast!
- Has a welcoming onboarding experience.



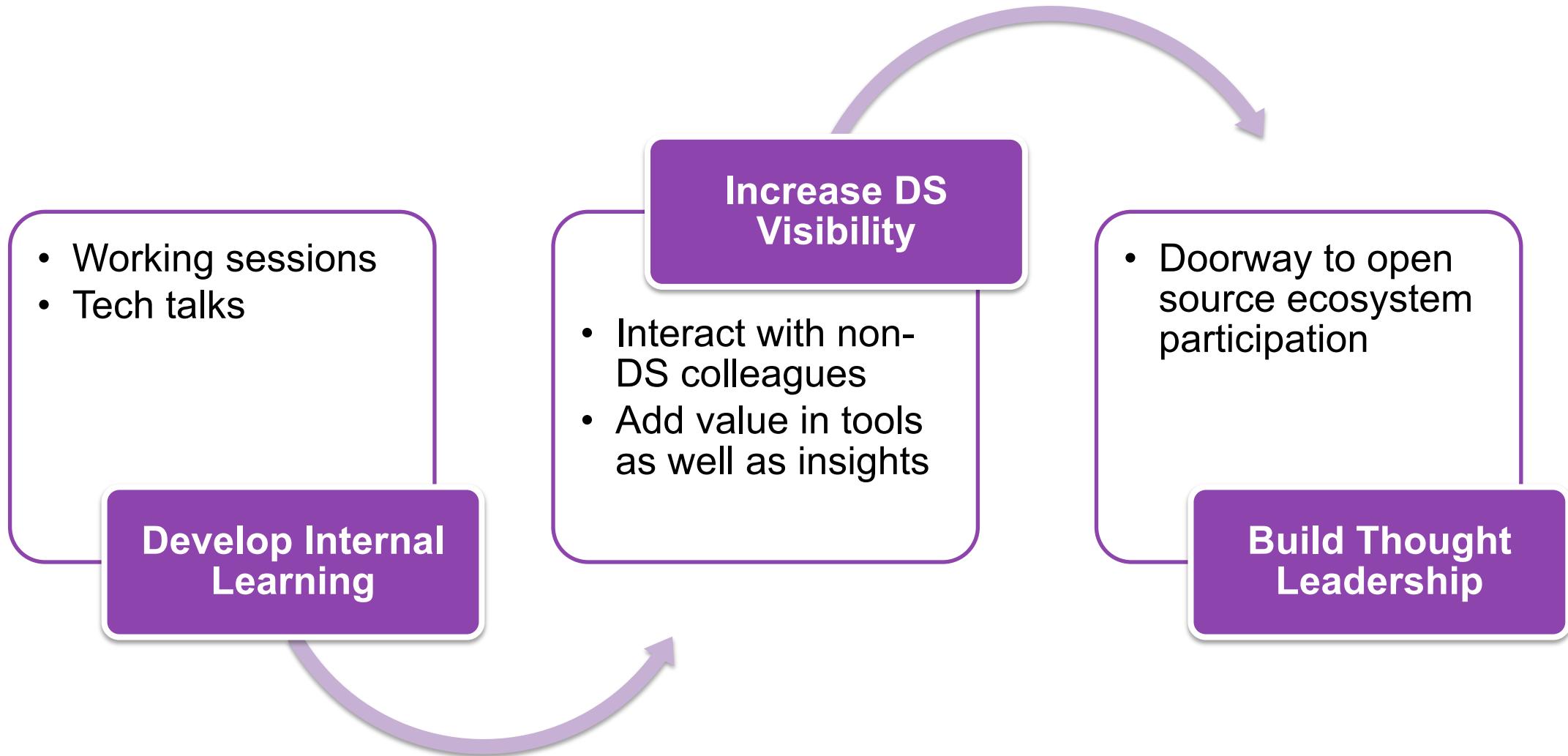
Staff Departure

- Causes no panic or scramble.
- Results in no lost information or knowledge.

Community tools increase the stability of your team for the long term.



Expand a culture of education across your organization.



How?

A group of individual data scientists != a collaborative team of data scientists.

Develop a paradigm of thinking that is productive and works towards your objectives.

- Common, understood goals
- Trust in colleagues
- Willingness to share information and help each other



Leadership has to commit to changes.

People need to be incentivized to maintain and contribute to new tools through:

- Time
- Compensation
- Recognition



Adam McElhinney

Yesterday at 9:48 AM

in #uptake

Team, Data Science and the Uptake Open Source Committee is pleased to announce that Data Science had their second R package accepted to CRAN! This is a great opportunity for us to showcase our Data Science skills to prospective customers and candidates. Special shoutout to

[@brian.burns](#) [@james.lamb](#)

[@patrick.boueri](#) [@jay.qi](#) who worked nights and weekends to get this initiative done. (edited)

1 reply



Ganesh Bell 19 hours ago

brilliant.. I want a TShirt "My data scientists work on more impactful problems than yours!"

Have a vision for what sort of system you want your team to have, and build infrastructure that can scale and evolve.



Start by thinking and planning.

Explore your undocumented resources

- Learn your institutional history/knowledge.
- Investigate existing best practices.

Catalog user needs

- Who do you want using your tools? What will they need?

Think for the long haul

- Your needs today will change tomorrow.



Build out the tech infrastructure you need thoughtfully.

Don't take on trendy tools you don't really need!

Internal git infrastructure

- Set up a system for versioning and code sharing, and communicate it to your team.

Building functions and packages

- After planning, construct the actual code.
- Many useful tutorials exist to help you along (see reference slide)



Write functions and documentation intelligently.

Make modular functions, follow best practices

- One operation/task = one function.
- Think about readability and future editability.
- Parameterize any aspects the user may want to specify.

(Remember to plan for growth and change!)

Good jobs for one function:

- Calculate a value
- Generate one plot
- Produce one dataframe

BAD job for one function:

- Calculate seventeen values, produce six dataframes, and produce four plots based on those dataframes



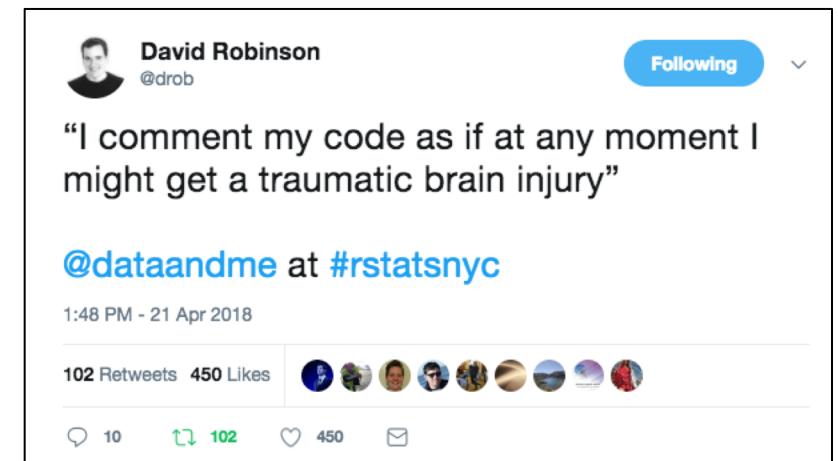
Make robust documentation a norm and a habit.

At the earliest stages of team infrastructure, documentation is haphazard.

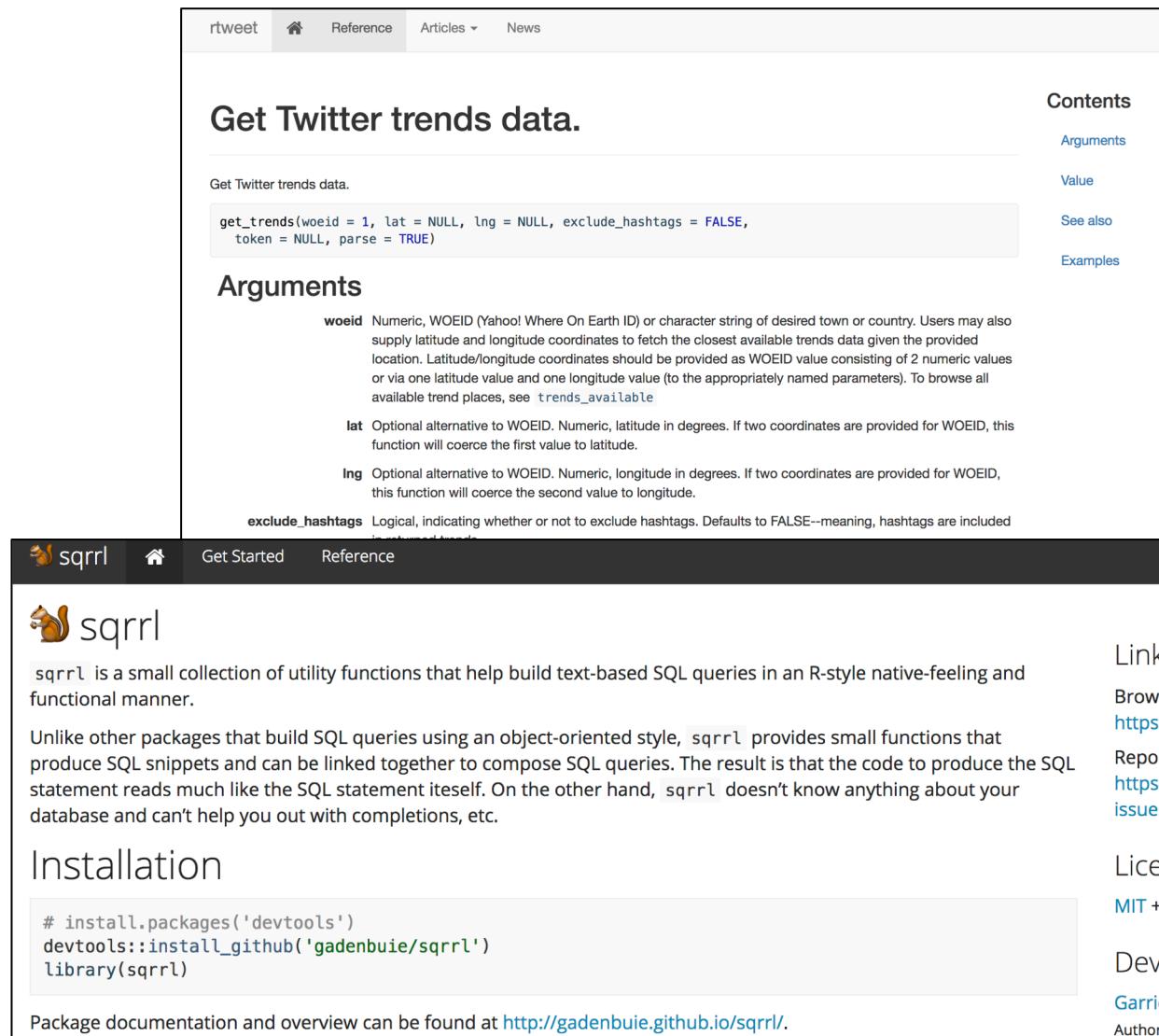
Write for the layperson or for ‘future you’.

R packages have a huge assortment of documentation possibilities!

- Function help docs
- Vignettes
- Package level descriptive documents
- Websites via pkgdown (<http://pkgdown.r-lib.org/>)



Great documentation



The screenshot shows the documentation for the `rtweet` package. At the top, there's a navigation bar with links for `rtweet`, `Home`, `Reference`, `Articles`, and `News`. Below the navigation, a section titled "Get Twitter trends data." contains a brief description and a code example:

```
get_trends(woeid = 1, lat = NULL, lng = NULL, exclude_hashtags = FALSE,  
          token = NULL, parse = TRUE)
```

On the right side, there's a sidebar with links for `Contents`, `Arguments`, `Value`, `See also`, and `Examples`. Below the sidebar, there's a "Links" section with links to source code and bug reports.

Sqrrl Documentation

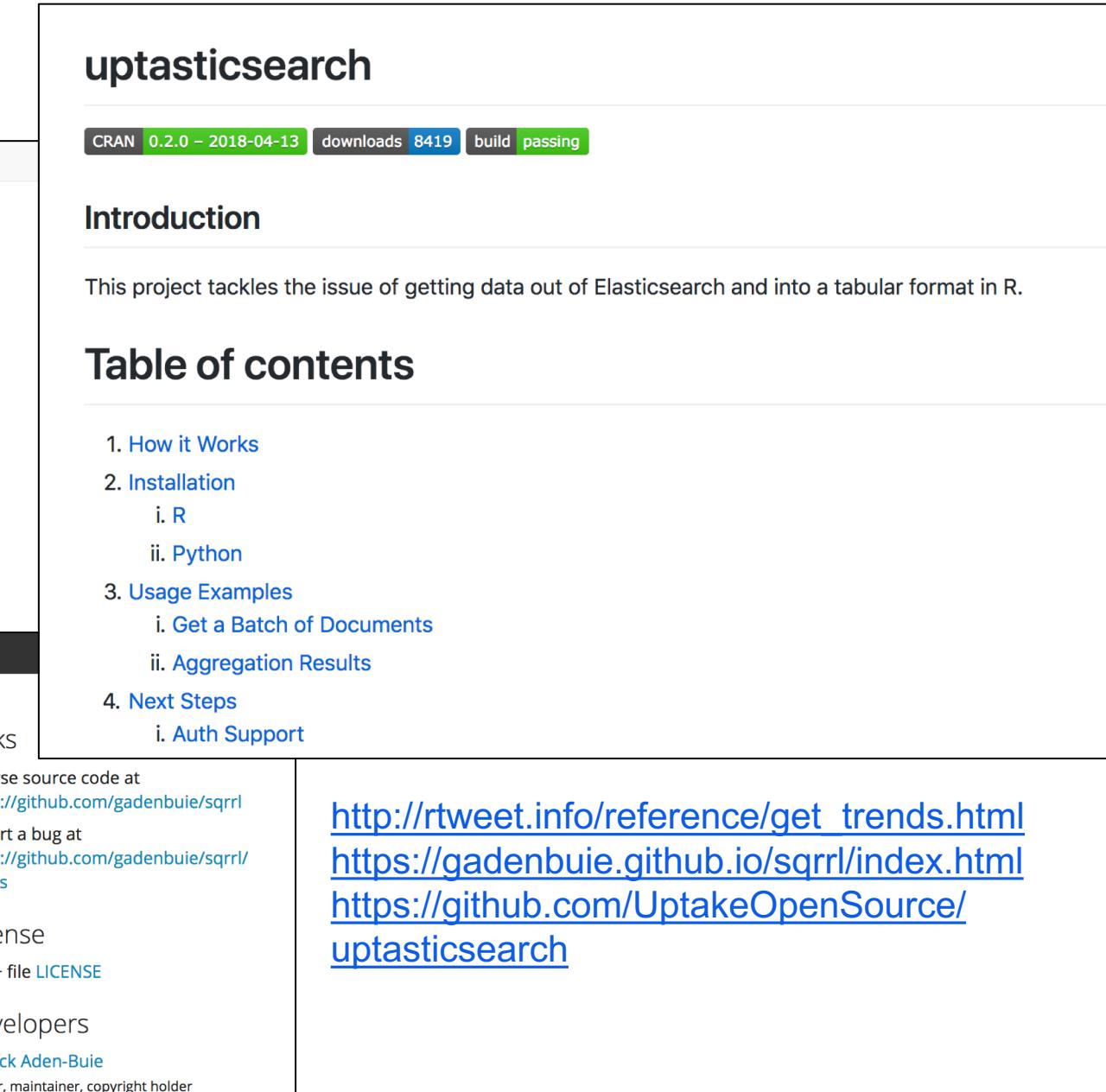
The Sqrrl documentation page has a header with a squirrel icon, "sqrrl", and links for "Get Started" and "Reference". The main content area starts with a heading "sqrrl" and a brief description of what it does. It then discusses how it differs from other SQL query builders:

Unlike other packages that build SQL queries using an object-oriented style, `sqrrl` provides small functions that produce SQL snippets and can be linked together to compose SQL queries. The result is that the code to produce the SQL statement reads much like the SQL statement itself. On the other hand, `sqrrl` doesn't know anything about your database and can't help you out with completions, etc.

Installation

```
# install.packages('devtools')  
devtools::install_github('gadenbuie/sqrrl')  
library(sqrrl)
```

Package documentation and overview can be found at <http://gadenbuie.github.io/sqrrl/>.



The screenshot shows the documentation for the `uptasticsearch` package. At the top, there's a header with the package name and CRAN information: "CRAN 0.2.0 – 2018-04-13", "downloads 8419", and "build passing".

Introduction

This project tackles the issue of getting data out of Elasticsearch and into a tabular format in R.

Table of contents

- [How it Works](#)
- [Installation](#)
 - [R](#)
 - [Python](#)
- [Usage Examples](#)
 - [Get a Batch of Documents](#)
 - [Aggregation Results](#)
- [Next Steps](#)
 - [Auth Support](#)

Links

Browse source code at <https://github.com/gadenbuie/sqrrl>
Report a bug at <https://github.com/gadenbuie/sqrrl/issues>

License
[MIT + file LICENSE](#)

Developers
[Gerrick Aden-Buie](#)
Author, maintainer, copyright holder

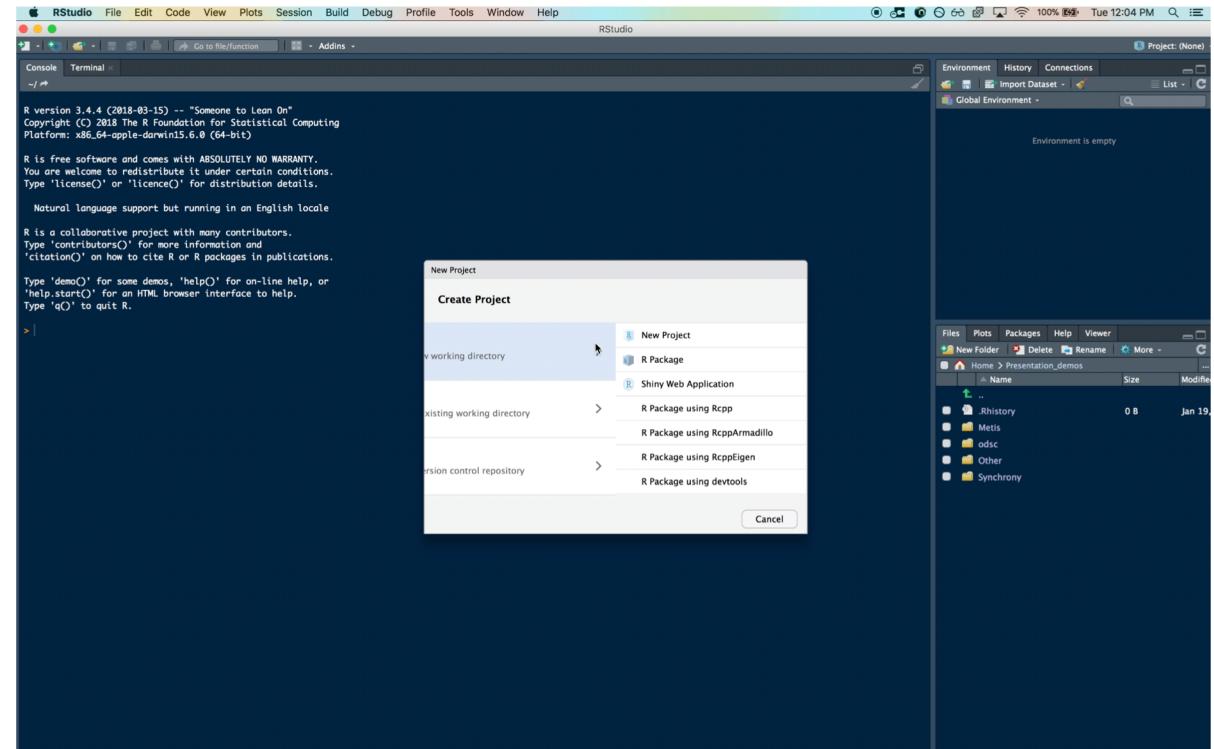
http://rtweet.info/reference/get_trends.html
<https://gadenbuie.github.io/sqrrl/index.html>
<https://github.com/UptakeOpenSource/uptasticsearch>



Building a Simple R Package - Tutorial

Essential Steps

1. Write or collect some functions to start with
2. Environment setup (libraries)
3. Initialize package skeleton
4. Write documentation (each function and the overall package)
5. Build
6. Share!



See my video tutorial and scripts at
<https://github.com/skirmer/odsc2018> to follow along step by step.



R Package Infrastructure

Files:

- DESCRIPTION: package overview
- NAMESPACE: dependencies/exporting management

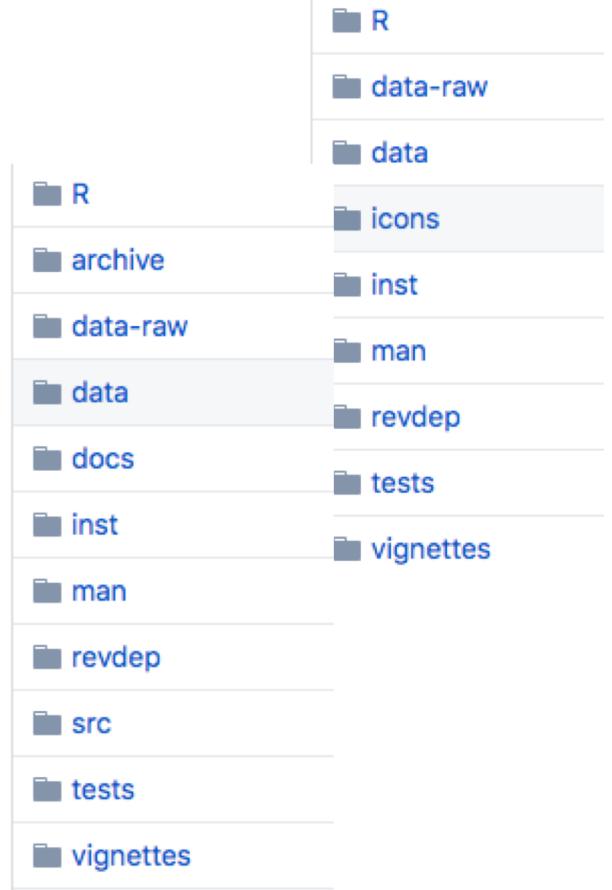
Folders:

- R: Function scripts
- man: Function help docs
- vignettes: Rmd files explaining the functions' operation in depth
- tests: unit testing code to ensure correct function behavior

AND THEN: you can include pretty much anything else you like.

inst, data, images, docs, etc.

Even the most complex, powerful R packages out there all use this same construction.



When you have tools ready to be used, teach your colleagues.

Let data scientists who authored functions teach their peers.

Proactively teach BI, analytics, executive, and sales staff

- What is this, and why should they care?
- Empower your colleagues to be your collaborators.

Over time, add more things that people need.



data
diagnostics
docs
graphs
lib
logs
munge
profiling
reports

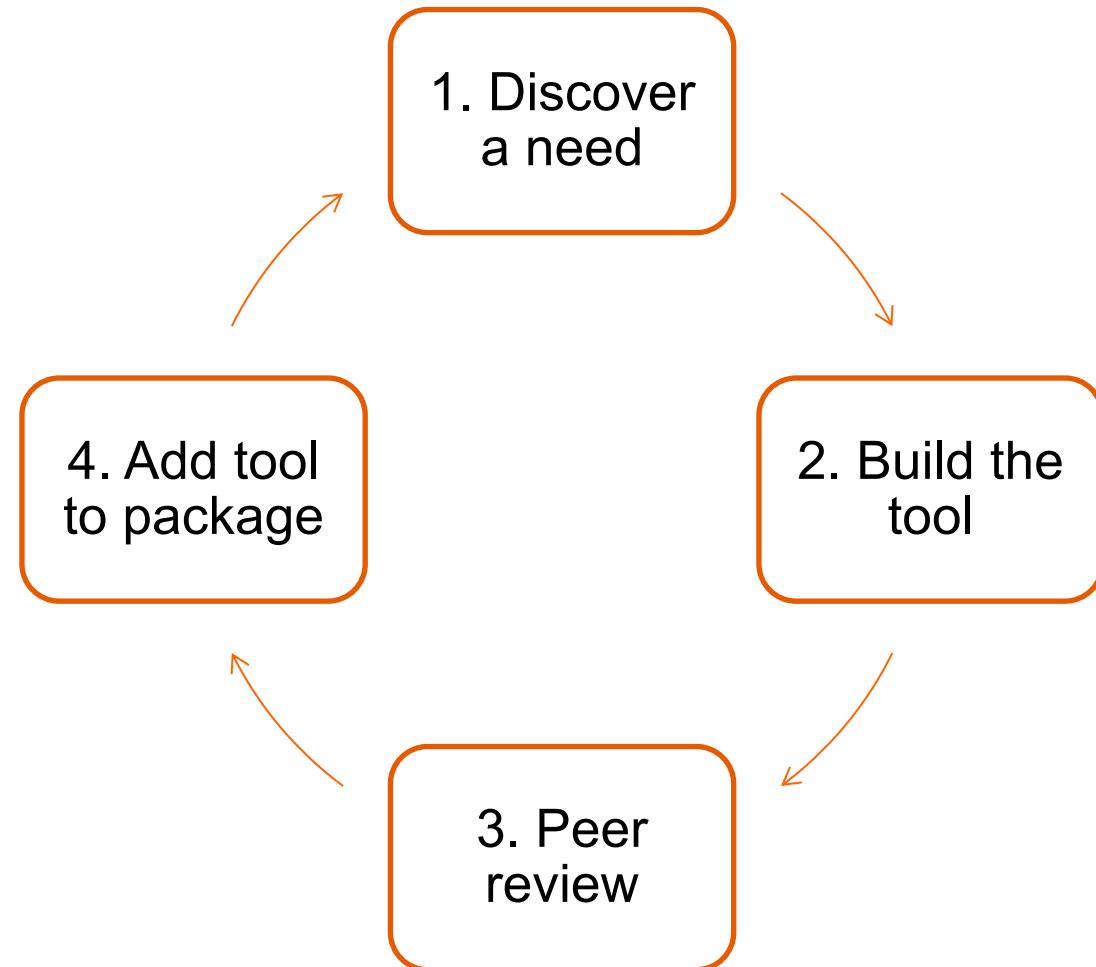
data
diagrams
extra
images
rmarkdown
screenshots

Add new features when the need arises, in a systematic way.

Ideas for new features can come from any user – welcome feedback

Don't skimp on the peer review element!

Do occasional housecleaning to ensure things are still working, efficient, and relevant



Make collaboration a core priority across the team.

Support improving the package, not just creating the package.

Stagnating documentation/tools < than no tools at all.

Move from single user ownership to team ownership.

- How you talk about the tool matters. Get people invested in it!

The screenshot shows a GitHub issue page for the repository 'f2m2 / apidocs'. The 'Issues' tab is selected, showing 2 open issues. The specific issue is titled 'Outdated dependency. #35'. It was opened by 'Addgod' on Mar 6, with 2 comments. The first comment is from 'Addgod' on Mar 6, saying 'Hi.' and asking if there's a reason for using an old version of reflection-docblock. The second comment is from 'f2m2rd' on Mar 6, edited, stating 'Thanks for the feedback. The version of reflection-docbloxk was current at the time of development.' and encouraging pull requests.



Other Technical Notes

Unit testing

- Programmatic checks that automatically run when a package is built
- Prevents unintentional breaking changes
- When to write unit tests, and how many, is a topic of debate.

Internal CRAN

- Set up an internal-only repository where your packages live
- They can be installed with a traditional `install.packages()` call with a few small setup steps
- <https://cran.r-project.org/doc/manuals/R-admin.html#Setting-up-a-package-repository>

Continuous Integration

- Contributions are built automatically in test environment to ensure constant compatibility



Wrapping Up

Having smooth functioning, collaborative data science is possible!

You can do work that is reproducible, high quality, and less effort.

Your entire organization can be more data literate and self supporting, freeing up data scientists for innovation.

You can build a culture of learning and education for your whole organization, starting from data science.

Suggested Roadmap

1. Develop collaborative culture;
2. Get leadership support;
3. Plan and think ahead;
4. Build tools and write code thoughtfully;
5. Don't skip documentation;
6. Educate your organization;
7. Maintain the tools!



Reference Material

How-To:

- Write your own functions:
<https://nicercode.github.io/guides/functions/>
- Write great documentation: <http://r-pkgs.had.co.nz/man.html>
- Construct an R package:
 - https://github.com/skirmer/r_packages
 - <https://github.com/jtleek/rpackages>
- Build websites from package documentation:
<http://pkgdown.r-lib.org/>
- Get familiar with git:
 - <https://product.hubspot.com/blog/git-and-github-tutorial-for-beginners>
 - <https://guides.github.com/>
 - <https://try.github.io/levels/1/challenges/1>
- Write unit tests:
 - <https://katherinemwood.github.io/post/testthat/>
 - <http://r-pkgs.had.co.nz/tests.html>

Case Studies:

Capital One:

https://github.com/emilyriederer/references/blob/master/tidycf_Riederer_rstudio.pdf

Airbnb: <https://peerj.com/preprints/3182/>

Tools:

- <https://bitbucket.org/>
- <https://about.gitlab.com/>
- <https://jenkins.io/>
- <https://travis-ci.org/>

@data_stephanie

www.stephaniekirmer.com

<https://github.com/skirmer/odsc2018/>

