

Scaling Data Science Education

Sean Kross
May 10, 2016

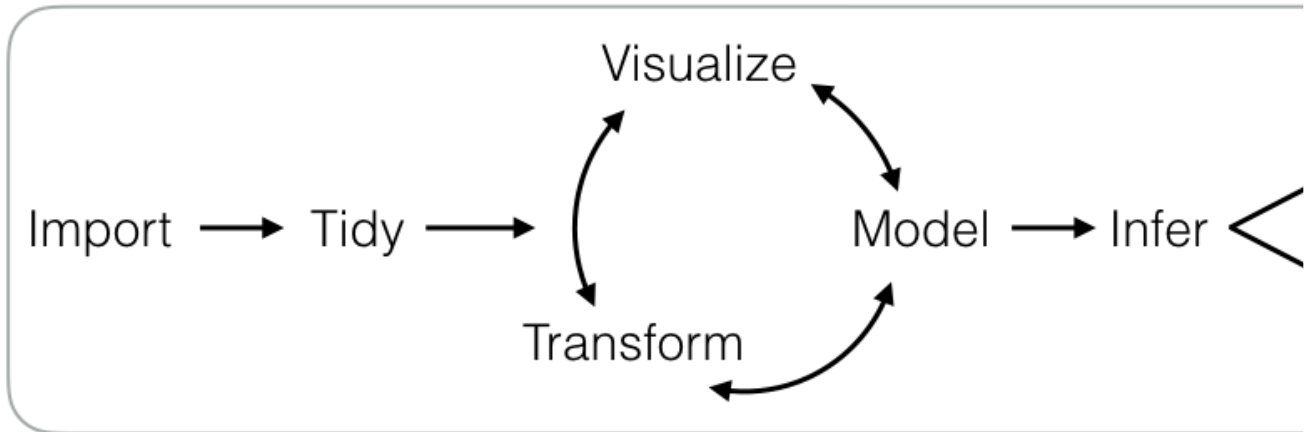
Who am I?

- Genomics at New York University
- Computer Science at The University of Maryland
- Statistics at Johns Hopkins University
- Currently I work in the Department of Biostatistics at JHU
- The path I've taken is problematic.

What is this talk about?

- What is data science?
- What is data science education?
- How can we scale data science education?
(And by "we" I mean you and I because I need your help!)

What is Data Science?: Process



<https://github.com/rstudio/RStartHere>

What is Data Science?: Skills

The tip of the iceberg:

(Probably and American colloquialism (sorry))

- Computer Science
 - Machine Learning
 - Distributed Computing
 - Parallel Computing
 - Algorithms
- Statistics
 - Inference
 - Regression
- + the scientific methods from a specific field (genomics, sociology, molecular gastronomy, etc)

Data Science Education: Motivation

a data scientist should be able to
run a regression, write a sql query, scrape a web
site, design an experiment, factor matrices, use a
data frame, pretend to understand deep learning,
steal from the d3 gallery, argue r versus python,
think in mapreduce, update a prior, build a
dashboard, clean up messy data, test a hypothesis,
talk to a businessperson, script a shell, code on a
whiteboard, hack a p-value, machine-learn a model.
specialization is for engineers.

JOEL GRUS



Jenny Bryan
@JennyBryan

Follow

A data scientist should be able to ...
My favourite slide from [@joelgrus](#)'s fun talk [#ddtx16](#)
8:33 PM - 17 Jan 2016

520

694

What is the underlying concept here?

What is Data Science Education?

- How can we most effectively teach students to manipulate and interrogate data using a computer?
- How can we scale this kind of education to millions of people?
- How can we improve upon issues that are currently present in computer science and mathematics education?

Attempts at Scaling Data Science Education



What is swirl?



{ swirl

What is swirl?

- swirl is an R package that turns the R console into an interactive learning environment for data science.
- swirl takes advantage of the R console's call-and-response behavior.
- swirl provides an authentic environment, there is no separation between where a student is learning data science and where they will go on to practice data science.

What problems does swirl solve?

- swirl is free and open source, and it always will be.
- swirl is available in multiple languages (we're always looking for volunteer translators!)
- Anyone can make and distribute their own swirl course for free.

What is a swirl course?

- A swirl course is a collection of lessons.
- Each lesson contains about 20 minutes of instruction about a particular topic in data science.
- Lessons are written and structured using yaml and R.

Demo

swirl

```
> library(swirl)
> swirl()
```

```
| Welcome to swirl! Please sign in. If you've been here before, u
| as you did then. If you are new, call yourself something unique
```

```
What shall I call you? sean
```

```
| Thanks, sean. Let's cover a couple of quick housekeeping items |
| our first lesson. First of all, you should know that when you s
| means you should press Enter when you are done reading and read
```

```
... <-- That's your cue to press Enter to continue
```

```
| Also, when you see 'ANSWER:', the R prompt (>), or when you are
| from a list, that means it's your turn to enter a response, the
| continue.
```

swirl

Selection: 1

|

| In this lesson, we will explore some basic building blocks of the R programming language.

...

|===

| If at any point you'd like more information on a particular topic related to R, you can type `help.start()` at the prompt, which will open a menu of resources (either within RStudio or a web browser, depending on your setup). Alternatively, a simple web search often yields the information you're looking for.

...

|=====

| In its simplest form, R can be used as an interactive calculator. Type `5 + 7` and press the `enter` key.

>

Coming to swirl this summer:

- Improved lessons
- More courses
- Better channels for distributing swirl courses
- Feedback for course instructors to assess student's performance

We want you!

The screenshot shows the GitHub profile for the 'swirl development team'. At the top, there's a navigation bar with 'This organization' and a search box, followed by links for 'Pull requests', 'Issues', and 'Gist'. The main header features the organization's logo (a blue curly brace with an 'S') and the name 'swirl development team', with the tagline 'Home of the swirl R package.' and contact information: 'http://swirlstats.com' and 'info@swirlstats.com'. Below this is a menu with 'Repositories', 'People 8', 'Teams 3', and 'Settings'. A search bar for repositories is present, along with a '+ New repository' button. The repository list includes:

- swirl_courses**: A collection of interactive courses for the swirl R package. Updated a day ago. 1,721 stars, 4,406 forks.
- swirlify**: A toolbox for writing swirl courses. Updated 11 days ago. 33 stars, 37 forks.
- swirl**: Learn R, in R. Updated 26 days ago. 482 stars, 342 forks.

On the right side, there is a 'People' section with a profile picture and an 'Invite someone' button.

Acknowledgements

- Co-Creators:
 - Nick Carchedi, Bill Bauer, Gina Grdina
- Generous support from:
 - Brian Caffo, Roger Peng, Jeff Leek
- Financial support from:
 - The United States [National Institutes of Health](#)

Questions?

- Find me:
 - On Twitter: @seankross
 - On GitHub: seankross
 - Online: <http://seankross.com>
 - Email me: sean@seankross.com
- Find swirl:
 - On Twitter: @swirlstats
 - On GitHub: swirldev
 - Online: <http://swirlstats.com>
 - Email swirl: info@swirlstats.com