# Building A Data Science Team From Scratch

*... cómo iniciar un equipo de data science cuando todos usan Excel*

*[how to start a data science team when everyone is still using Excel]*

# About Me_

twitter @catherinezh

#rstats

#rladies

#codecademy

- **Based in New York City**
- **Data Scientist Manager @ Codecademy**
  - Company's first data science hire
- **10+ years of statistical programming**
  - 6 yrs of professional work
  - 4 yrs of technical hiring & interviewing

codecademy

# Lessons Learned from
# Both Sides of the Hiring Table

codecademy

# Agenda

**By the end of this talk you will be able to:**

- **Hire the right type** of data professional
- Decide how to **structure your data team**
- Conduct a **data science technical screen**
- **Identify the right tools** for assessments
- When and how to **scale up your team**

codecademy

# Data Science Hiring_

codecademy

**Everybody wants to hire data scientists,
but no one can agree on a job title and description.**

MENU

DATA

# Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

FROM THE OCTOBER 2012 ISSUE

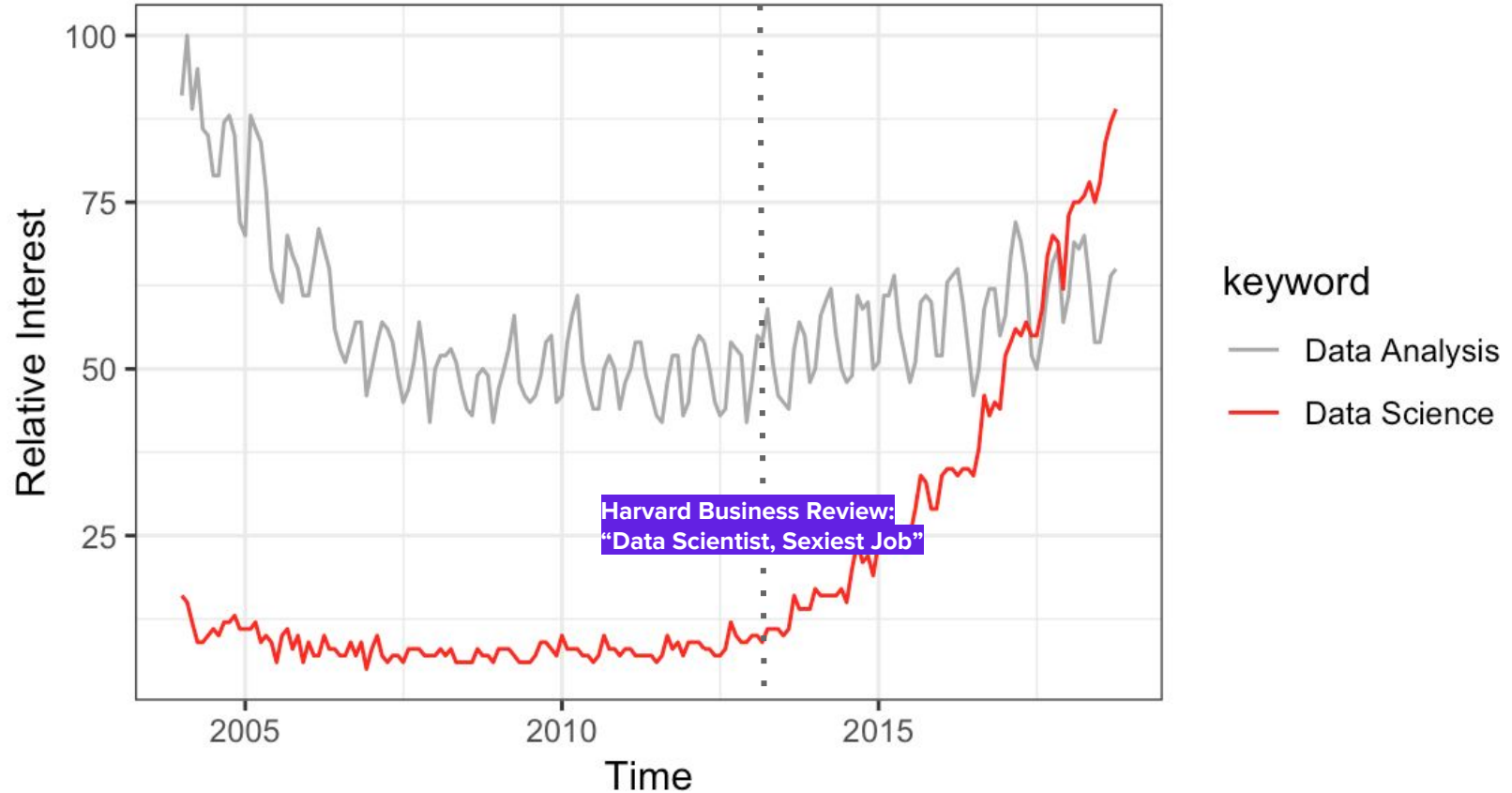SUMMARY | SAVE | SHARE | COMMENT (4) | TEXT SIZE | PRINT | $8.95 BUY COPIES

**W**hen Jonathan Goldman arrived for work in June 2006 at LinkedIn, the business networking site, the place still felt like a start-up. The company had just under 8 million accounts, and the number was growing quickly as existing members invited their friends and colleagues to join.

# Google Trends Data

## Search volume between Jan'04-Oct'18

# Data janitor

A **data janitor** is a person who works to take big data and condense it into useful amounts of information. Also known as a "data wrangler,"

**Nicholas Chamandy**  Follow

Scientific Director at Lyft

Apr 18 · 5 min read

# What's in a name?

The semantics of Science at Lyft

At Lyft, we're rebranding our Data Analyst function as *Data Scientist*, and our Data Scientist function as *Research Scientist*. In this short post, we describe the reasoning behind the change, which we believe will set Lyft up to make better decisions and build better products as we scale.

# Types of Data Scientists_

codecademy

# DATA SCIENTIST TYPE A

## VS.

# DATA SCIENTIST TYPE B

- *Michael Hochster, Data Science Director at Stitch Fix (Quora 2014)*

codecademy

# DATA SCIENTIST TYPE <mark>ANALYZE</mark>

## DATA SCIENTIST TYPE B

# VS.

- *Michael Hochster, Data Science Director at Stitch Fix (Quora 2014)*

codecademy

# DATA SCIENTIST TYPE ANALYZE

Primarily concerned with making sense of data or working with it in a fairly static way. Very similar to a statistician (and may be one) but knows all the practical details of working with data that aren't taught in the statistics curriculum: data cleaning, methods for dealing with very large data sets, visualization, deep knowledge of a particular domain, writing well about data, and so on.

## VS.

# DATA SCIENTIST TYPE BUILD

- *Michael Hochster, Data Science Director at Stitch Fix (Quora 2014)*

codecademy

# DATA SCIENTIST TYPE <u>A</u>NALYZE

Primarily concerned with making sense of data or working with it in a fairly static way. Very similar to a statistician (and may be one) but knows all the practical details of working with data that aren't taught in the statistics curriculum: data cleaning, methods for dealing with very large data sets, visualization, deep knowledge of a particular domain, writing well about data, and so on.

**VS.**

# DATA SCIENTIST TYPE <u>B</u>UILD

- *Michael Hochster, Data Science Director at Stitch Fix (Quora 2014)*

code|cademy

# DATA SCIENTIST TYPE <u>A</u>NALYZE

Primarily concerned with making sense of data or working with it in a fairly static way. Very similar to a statistician (and may be one) but knows all the practical details of working with data that aren't taught in the statistics curriculum: data cleaning, methods for dealing with very large data sets, visualization, deep knowledge of a particular domain, writing well about data, and so on.
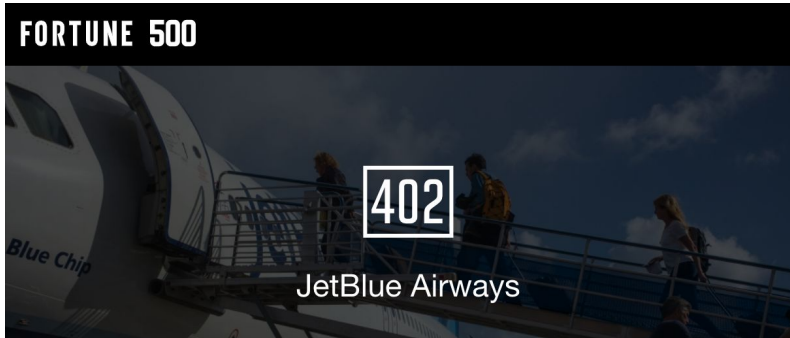
## VS.

# DATA SCIENTIST TYPE <u>B</u>UILD

Share some statistical background with Type A, but they are also very strong coders and may be trained software engineers. Mainly interested in using data "in production." They build models which interact with users, often serving recommendations (products, people you may know, ads, movies, search results).
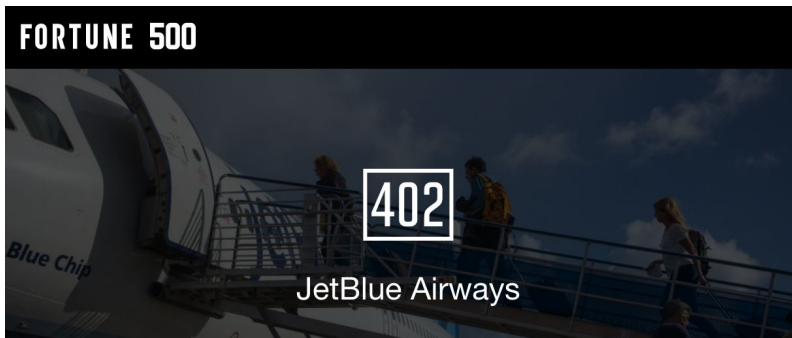
- *Michael Hochster, Data Science Director at Stitch Fix (*<u>Quora 2014</u>*)*

code|cademy

The stage your organization is in will guide the type of data scientist you should hire.

codecademy

**402**

JetBlue Airways

| | |
|---|---|
| Industry | Airlines |
| HQ Location | Long Island City, N.Y. |
| Website | www.jetblue.com |
| Years on Fortune 500 List | 6 |
| Employees | 17,424 |

**FORTUNE 500**

**402**

JetBlue Airways

| | |
|---|---|
| Industry | Airlines |
| HQ Location | Long Island City, N.Y. |
| Website | www.jetblue.com |
| Years on Fortune 500 List | 6 |
| Employees | 17,424 |

**code|cademy**

**Join the Millions**
**Learning To Code**

New York City

Education · Curated Web

51-200 employees

www.codecademy.com/   f   in

WHAT STAGE ARE YOU IN?

Assess what type of Data Science work your organization needs.

code|cademy

# DATA SCIENCE FOCUS AREAS

Early Stage:

- **Data collection phase**: database design, ETL
- Build foundation: biz intelligence, reporting

**DATA SCIENCE FOCUS AREAS**

Early Stage:

- Data collection phase: database design, ETL
- Build foundation: biz intelligence, reporting

Growth Stage:

- Define & measure key metrics, **produce insights**
- **Run experiments** on new and existing features

codecademy

## DATA SCIENCE FOCUS AREAS

Early Stage:

- Data collection phase: database design, ETL
- Build foundation: biz intelligence, reporting

Growth Stage:

- Define & measure key metrics, produce insights
- Run experiments on new and existing features

At Scale:

- Develop and deploy models into production.
- Automate analysis at scale

codecademy

## WHAT IS THE STATE OF YOUR DATA WAREHOUSE?

**Important note:**

The age of the company is not inherently tied to the state the data warehouse is in.

**Older companies often possess technical debt from ties to legacy software and old infrastructure.**

codecademy

**Start with DS generalists before specialists. Strong teams come from diverse backgrounds.**

code|cademy

# DATA SCIENTIST TRACKS

## Data Scientist – Analytics

Defines and monitors metrics, creates data narratives, builds tools

## Data Scientist – Algorithms

Builds and interprets algorithms that power data products

## Data Scientist – Inference

Establishes causal relationships with statistics

- *Elena Grewal, Head of Data Science at Airbnb (2018)*

PART TWO

# Technical Interviews_

codecademy

**Structure effective data science interviews
to weed out the data science "experts".**

# CODECADEMY HIRING PROCESS

1. Resume Review
2. Phone Interview
3. Pair Programming
4. Take-Home Project
5. On-Site Interview

Screening for:
Technical Skill &
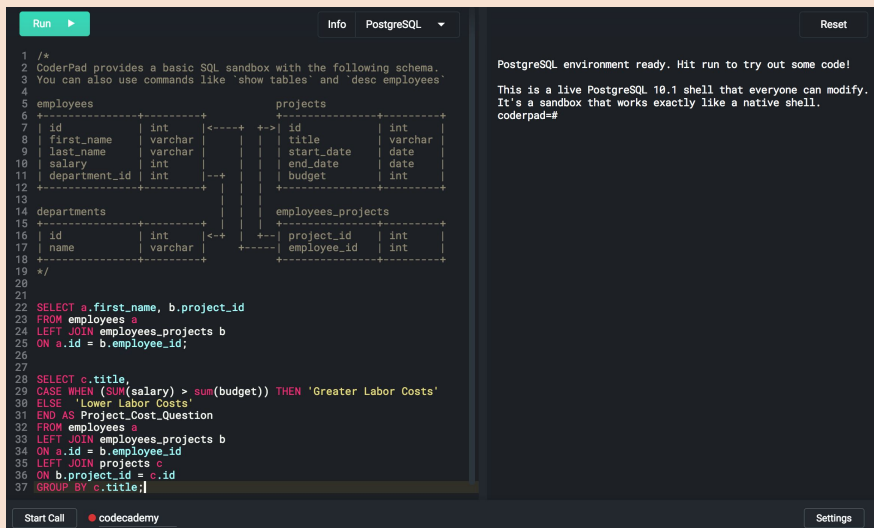Business Sense_

codecademy

# Phone interview

Question 1:

**What's your favorite R or Python package, and why?**

Question 2:

**Tell me about an analytical project you worked on recently.**

codecademy

# CODECADEMY HIRING PROCESS



## PAIR PROGRAMMING

- Last 15 min of 45 min phone interview
- Test ability to write advanced SQL queries using CoderPad.
- Bonus points if SQL queries are optimized

codecademy

# TAKE-HOME PROJECT

- Provide toy dataset of company data

- Candidates have a week to send back a write-up

- Findings should include business conclusions

- Can include supporting visualizations and models

- Assess their code, creativity, statistical soundness

code|cademy

# ON-SITE INTERVIEW

- **Take-Home Project Presentation**
  - **Present findings to cross-functional panel**
- **In-Person Whiteboarding Challenge**
  - **Design a database schema**
- **1-1 interviews with the rest of the team**

code cademy

Don't limit your hiring pool by focusing on previous job titles. Screen applicants based on their skills.

code|cademy

## FOCUS ON SKILLS, NOT JOB TITLES

Look beyond just the data scientist title. Many industries have been slow to adopt the data scientist title. Data scientists can also be called analysts, engineers, quants, specialists -- even data janitors, apparently!  It's a relatively new field and titles are not standardized.

# FOCUS ON SKILLS, NOT JOB TITLES

**Job Seekers:**

- **When reviewing job postings, don't focus on the job title.**

- **Focus on the responsibilities listed in the job description.**

- **Search keywords for specific skills: SQL, R, Python, clustering, AB testing, regression, NLP.**
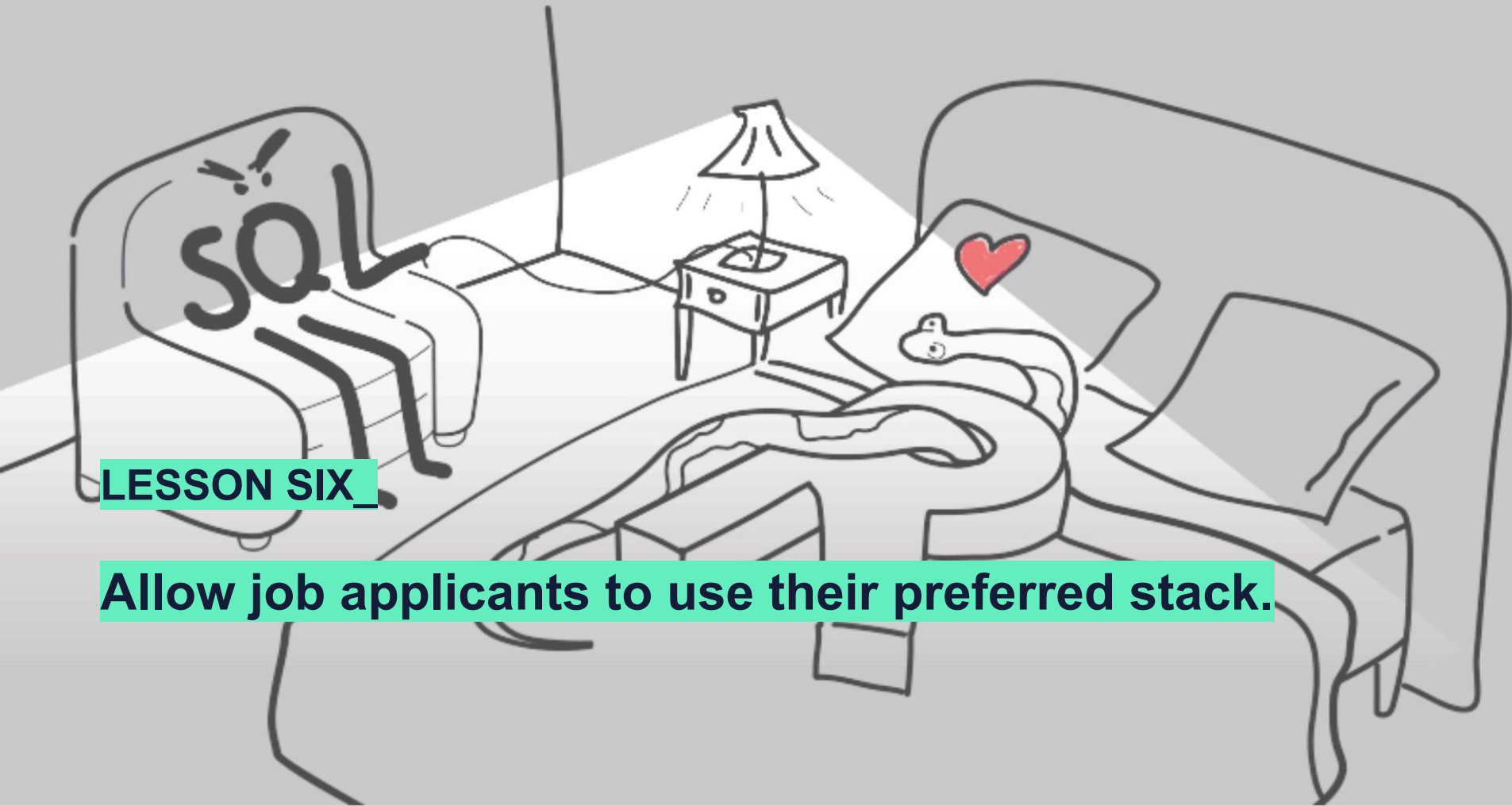
codecademy

# FOCUS ON SKILLS, NOT JOB TITLES

**Job Seekers:**

- **When reviewing job postings, don't focus on the job title.**

- **Focus on the responsibilities listed in the job description.**

- **Search keywords for specific skills: SQL, R, Python, clustering, AB testing, regression, NLP.**

**Interviewers:**

- **DS title hasn't been around for long, and titles vary by industry.**

- **Over-focusing on previous job titles will limit your applicant pool.**

- **Write truthful job descriptions: list specific projects and tools.**

codecademy

**LESSON SIX_**

**Allow job applicants to use their preferred stack.**

## WHAT LANGUAGE SHOULD I USE?

**Job Seekers:**

- **Use the tools you're most comfortable with to showcase your skills.**
- **Now is not the time to experiment with programs or methods you're not familiar with.**

code|cademy

# WHAT LANGUAGE SHOULD I USE?

**Job Seekers:**

- Use the tools you're most comfortable with to showcase your skills.
- Now is not the time to experiment with programs or methods you're not familiar with.

**Interviewers:**

- Be flexible. Let candidates use their preferred tools for their assessment.
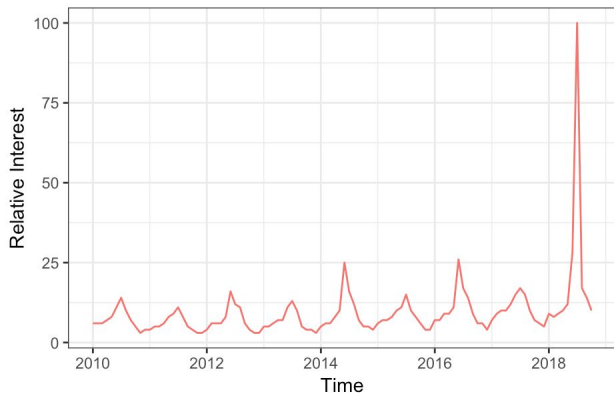- Their preferred stack might even include non-programming tools, and sometimes that's okay!

# Data science isn't always the 'sexiest job'.
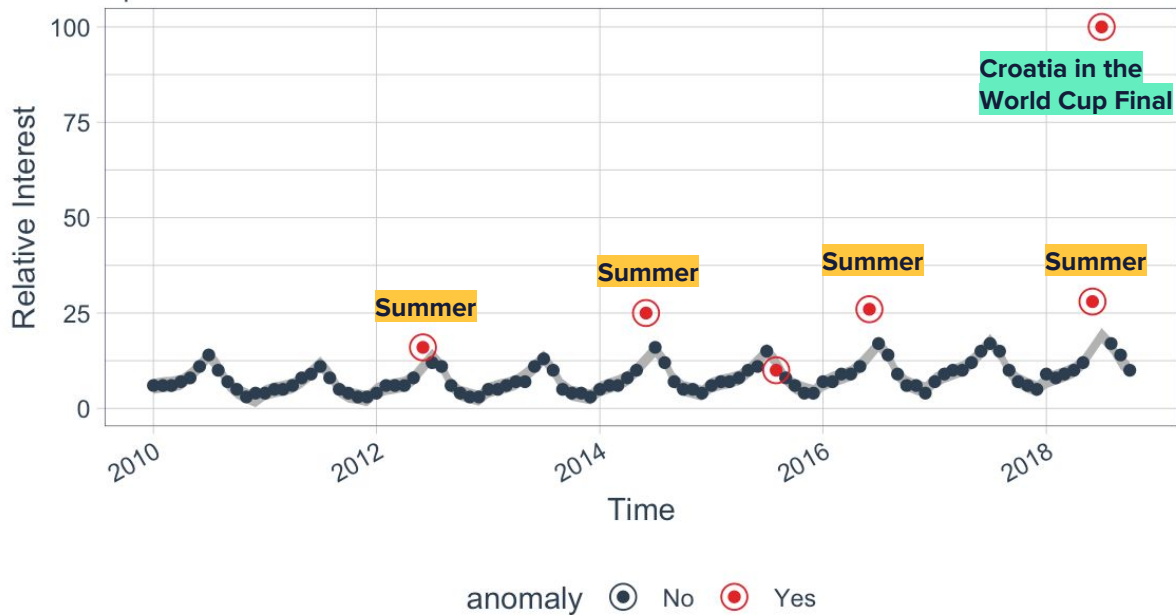
codecademy

# EXPECTATION



Google Trends Data
Spain search volume for 'Croatia' between Jan'10-Oct'18



Google Trends Data - Anomaly Detection
Spain search volume for 'Croatia' between Jan'10-Oct'18

Croatia in the World Cup Final

Summer

anomaly    No    Yes

# REALITY

**catherine** 🐱 4:43 PM

my typical workflow:                                          🔥 🔥 🔥

1) start working on an analysis i'm excited about

2) fire drill, everything else is derailed

3) somehow still working on the fire drill and other related issues

4) think longingly about the analysis i was planning to work on

For job seekers and hiring managers, prioritize finding a good mutual fit.

codecademy

# JOB SEEKERS

- Know what type of role you want:
  - Generalist vs Specialist
  - Early vs Late Stage Company
- What skills do you bring to the table?

code|cademy

# JOB SEEKERS

- Know what type of role you want:
  - Generalist vs Specialist
  - Early vs Late Stage Company
- What skills do you bring to the table?

# HIRING MANAGERS

- Assess candidates based on company needs for the year
- For skills that are "nice to have", screen for potential
- Ensure that you offer job seekers mutual growth and satisfaction
- Keep your data scientists happy and interested!

codecademy

# Additional Resources

- [Codecademy DS Reddit AMA](#)
- [Github: Slides For This Talk](#)
- [Airbnb: Data Science Tracks](#)
- [Doing Data Science At Twitter](#)
- [Lyft: Renaming Data Scientists](#)
- [Data Analyst vs Scientist](#)

codecademy

# ¡HASTA LA PRÓXIMA, MADRID!

twitter @catherinezh

github @cattystats

#rstats

#rladies

codecademy