# Push-Button Publishing in RStudio Connect

*Jeff Allen 11/2017*

*RStudio*

@trestleJeff
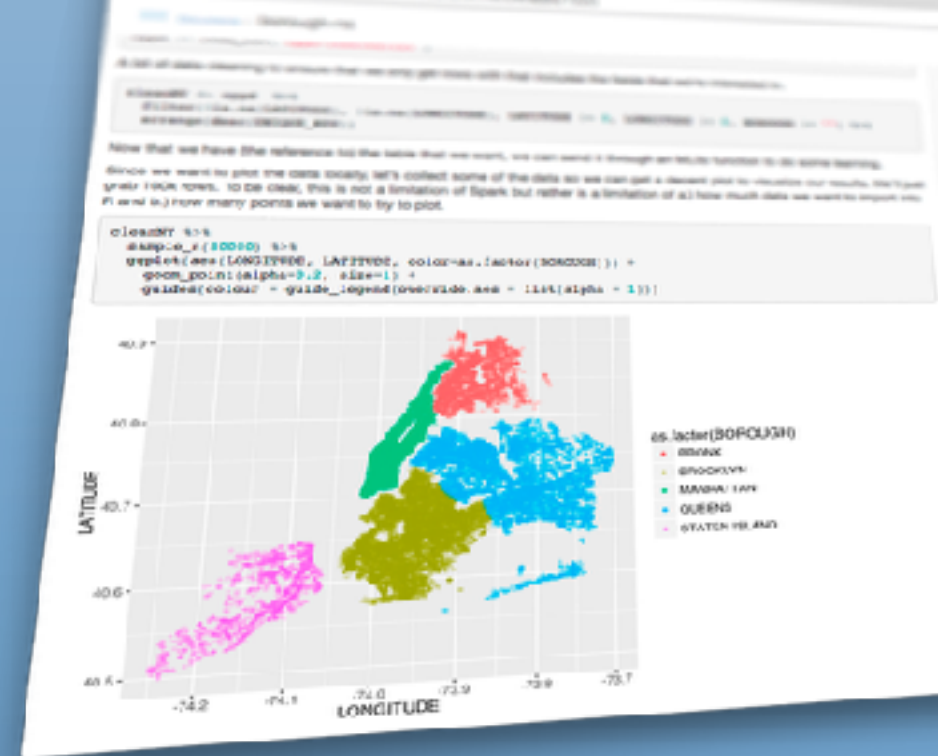
1

Publish Document...

Manage Accounts...

# RStudio Connect

- Push-button publishing from the RStudio IDE
- Manages all the content types you produce in R:

    Shiny, APIs, R Markdown, plots, etc.
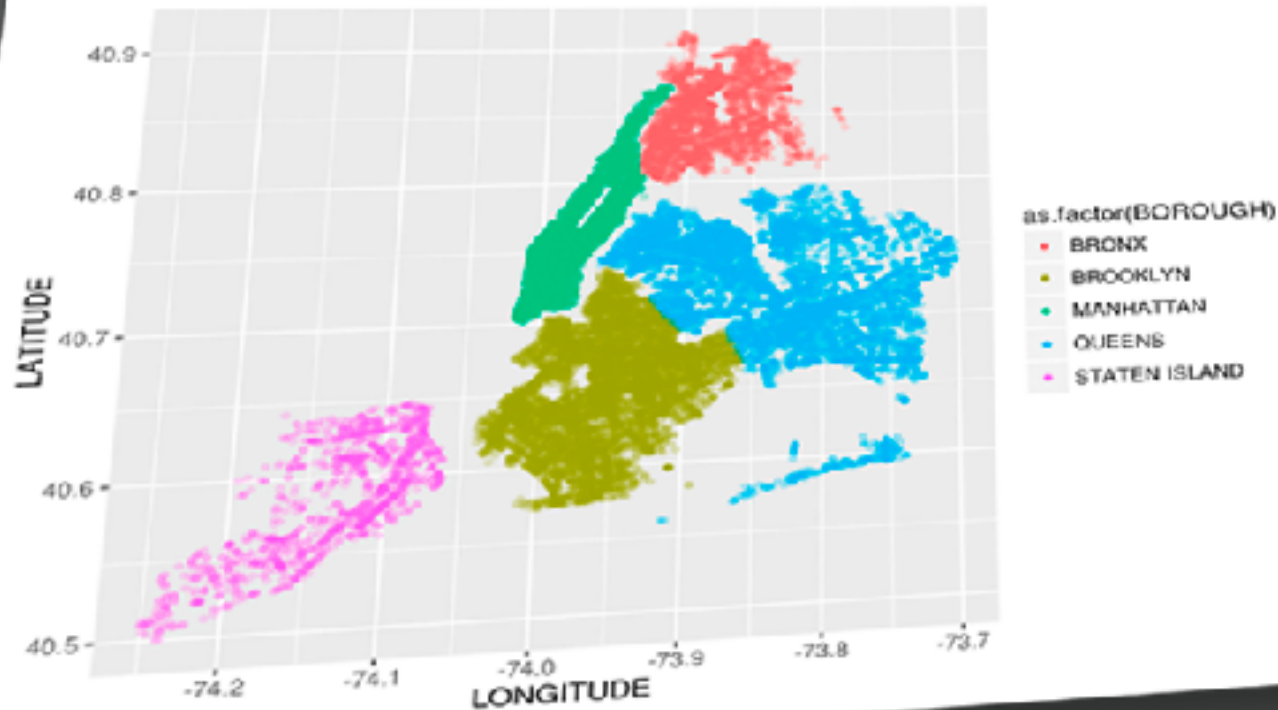
- On-premises, commercial
- Share data science artifacts

R Studio

# Three stages of R adoption in an organization

1. Bespoke analysis

2. Interactive tools
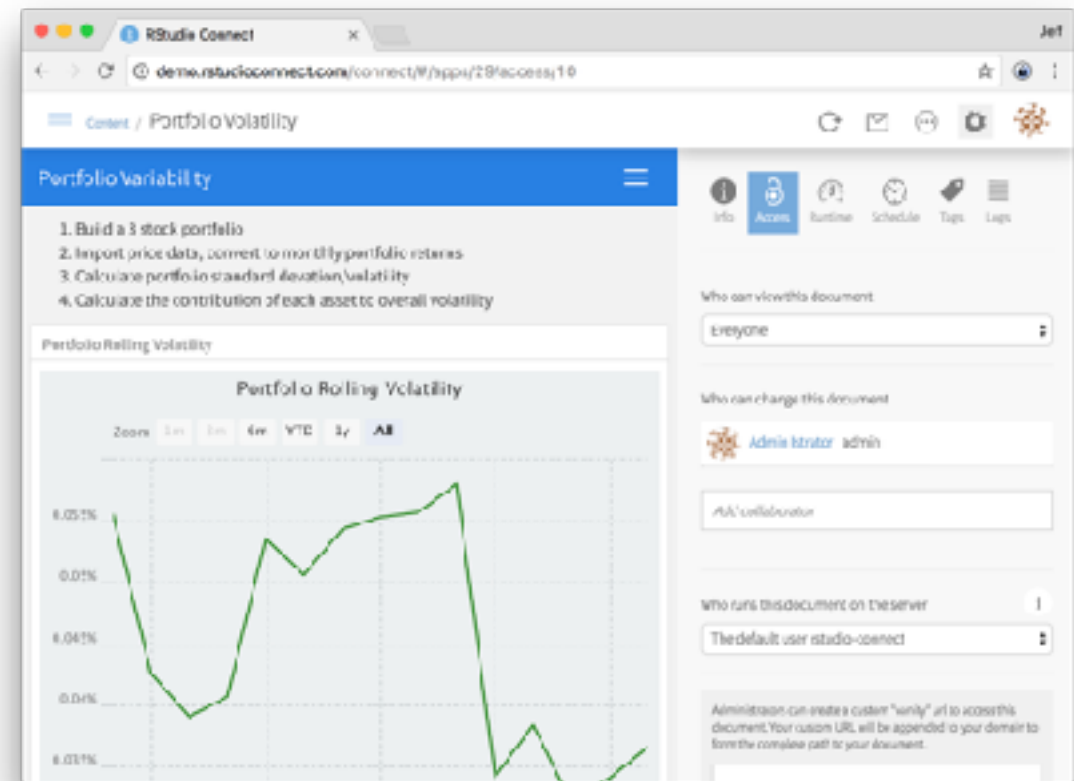
3. Fully integrated

# Stage 1: Bespoke Analysis

# Stage 1: Overview

- Data scientist is a black box

- Delivers a static result (plot, table, document)

- Revise question or ask new questions and start over

- Individual secretly installing R on your laptop to a group managing their own RStudio Server
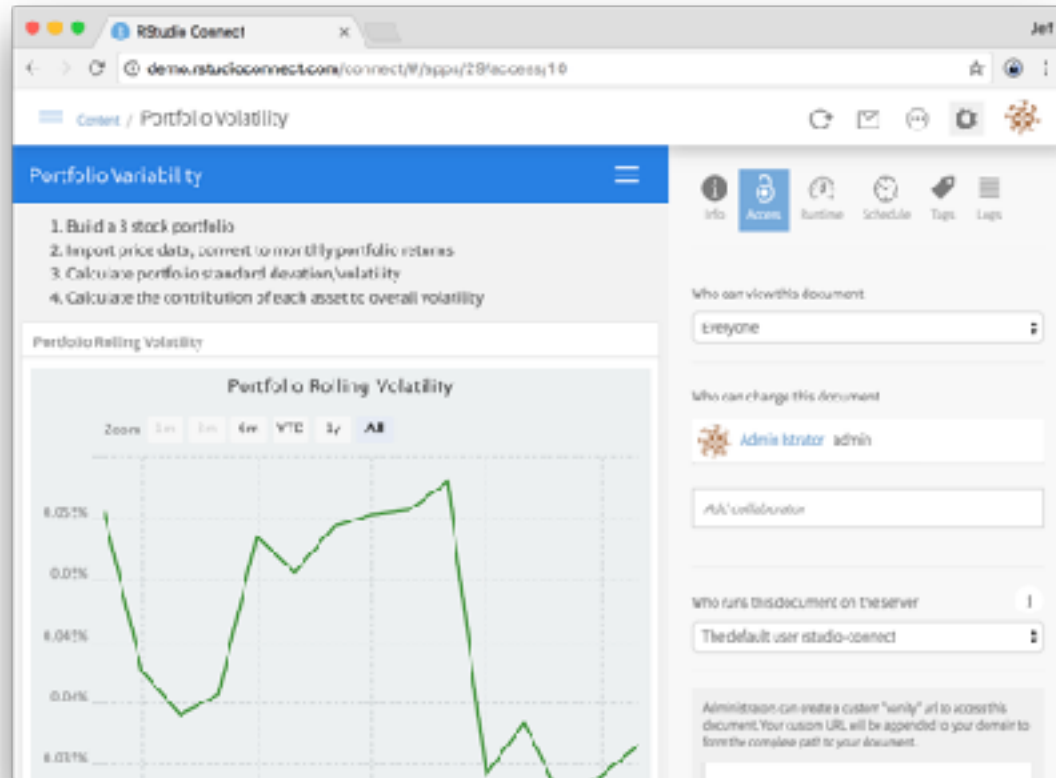
# Stage 1: Use Cases

- Sales reporting
- Investment analysis, reproducible decision making
- Game usage metrics reporting

# Stage 1: Issues

- Disorganized output

- Managing access

- Irreproducible
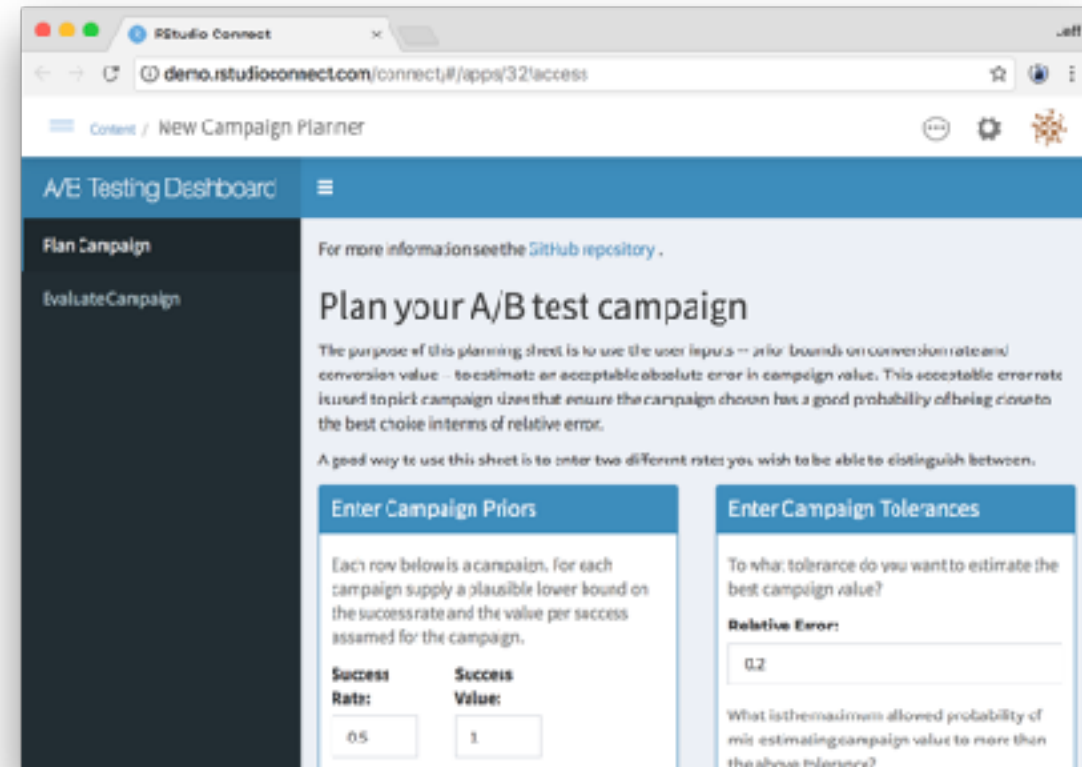
- No ability to automate

# Stage 1 Demo

# Stage 2: Interactive Tools

13

# Stage 2: Overview

- Some formal acceptance of R
  - Internal user-groups, training
  - Optimizing the on-boarding experience
- Non-R-users consume tools written in R
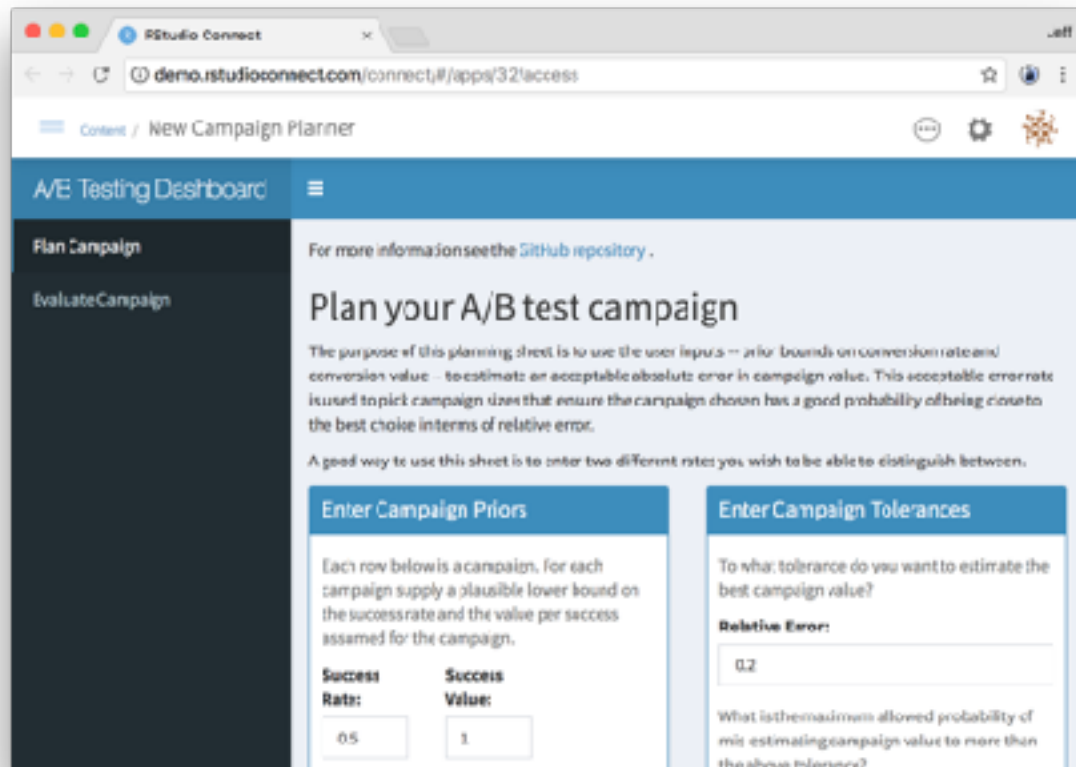  - Self-service
  - Shiny, Parameterized R Markdown

# Stage 2: Use Cases

- BI dashboarding
- Monitoring & analysis of customer support
- Domain-specific tooling
  - Split testing

https://github.com/WinVector/CampaignPlanner_v3

# Stage 2: Issues

- Hosting
  - Internal vs external
  - Authentication
  - Scalability & performance
- Publishing applications
- Testing

# Stage 2 Demo

# Stage 3: Fully Integrated

# Stage 3: Overview

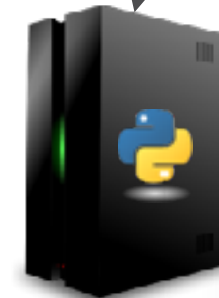- R is fully embraced
  - Permissive access to data
  - Share analysis broadly
  - Non-data-scientists get involved

- Incorporate R into other systems

Warehouse

# Stage 3: Issues

- Enabling other systems to communicate with R

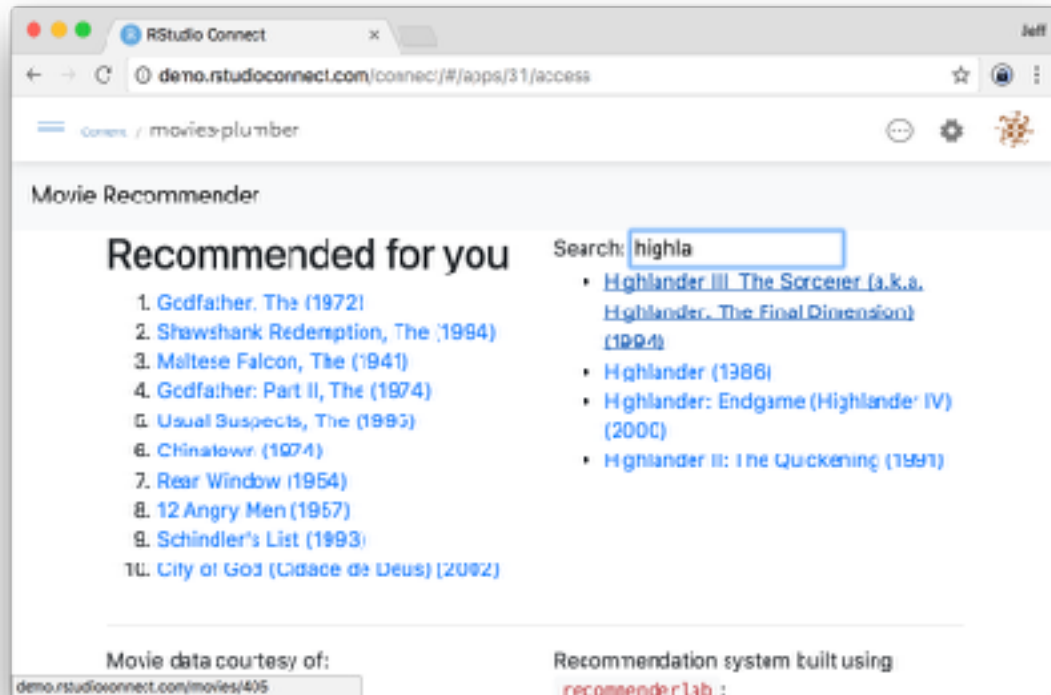- Code sharing, package management

- Data governance

# Plumber

- Open-source R package
- Make your R code available over a web API
- Consumed by systems in any other programming language
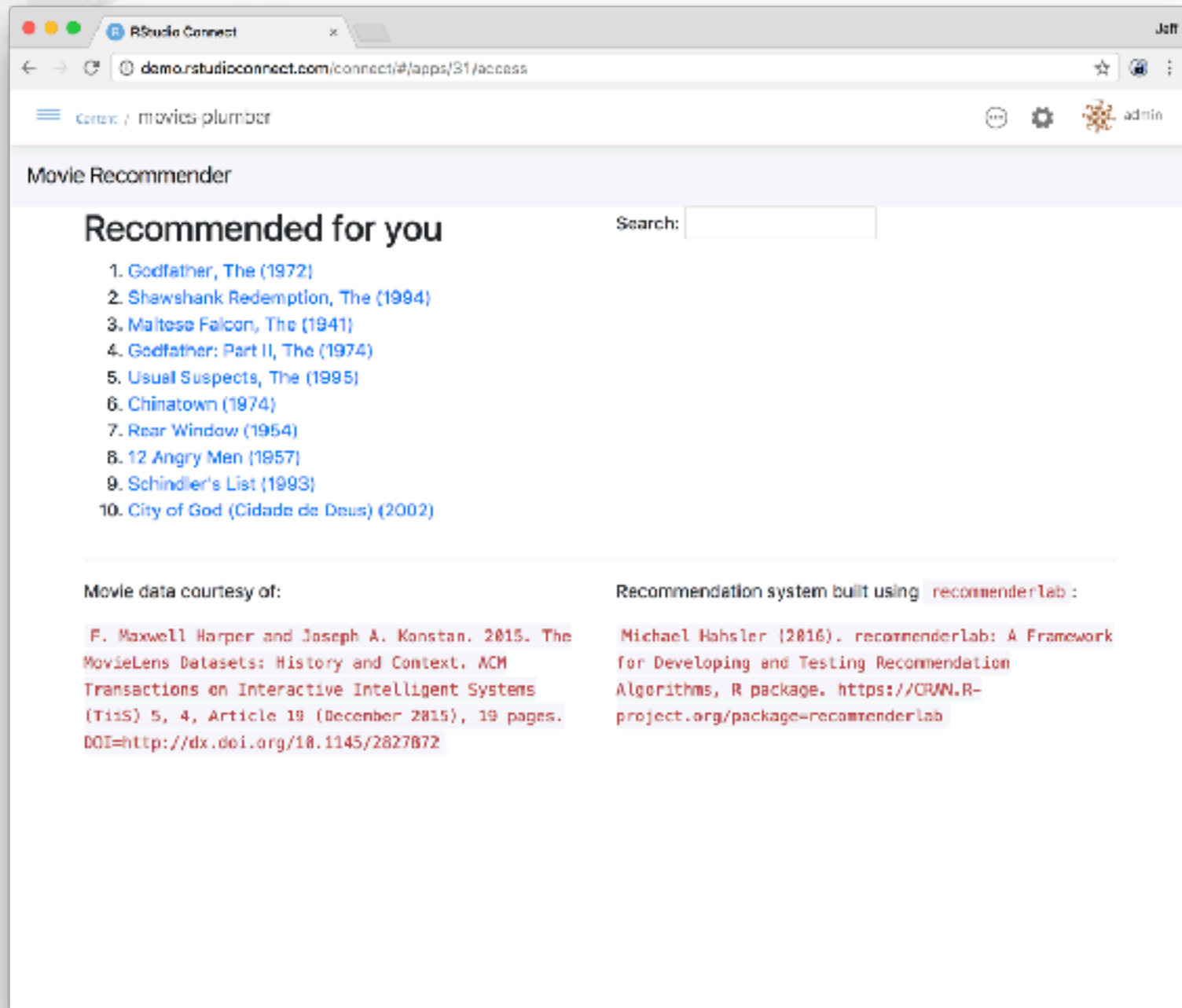- https://www.rplumber.io/

# Stage 3: Use Cases

- Leverage R from existing enterprise system: Java, Python, Ruby, etc.
  - Compare to porting your R code into another language
- Recommendation/forecasting engine
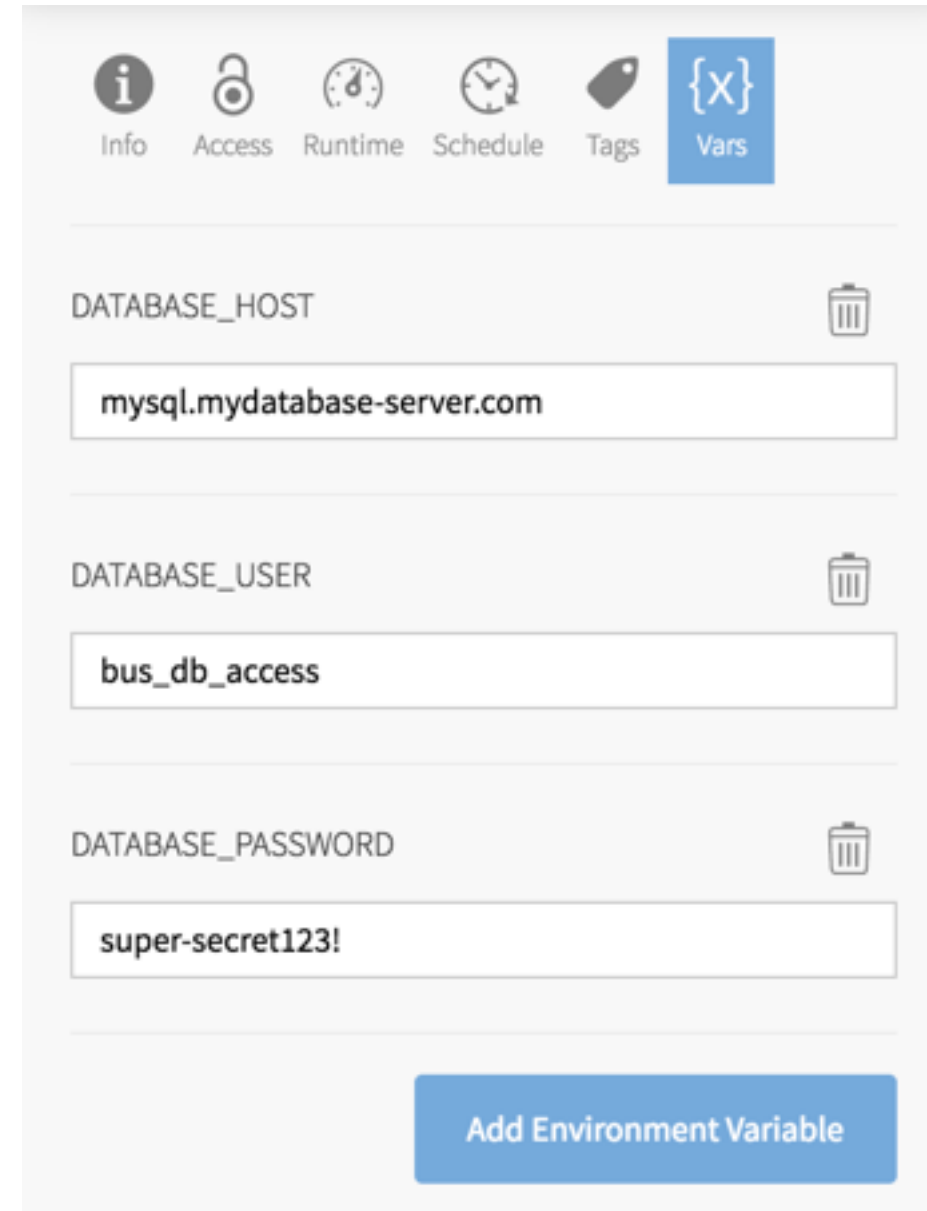- Python ETL jobs in a scheduled RMD

# Stage 3 Demo

# Updates

# What's New?

- Load balancing & high availability
- Floating licenses
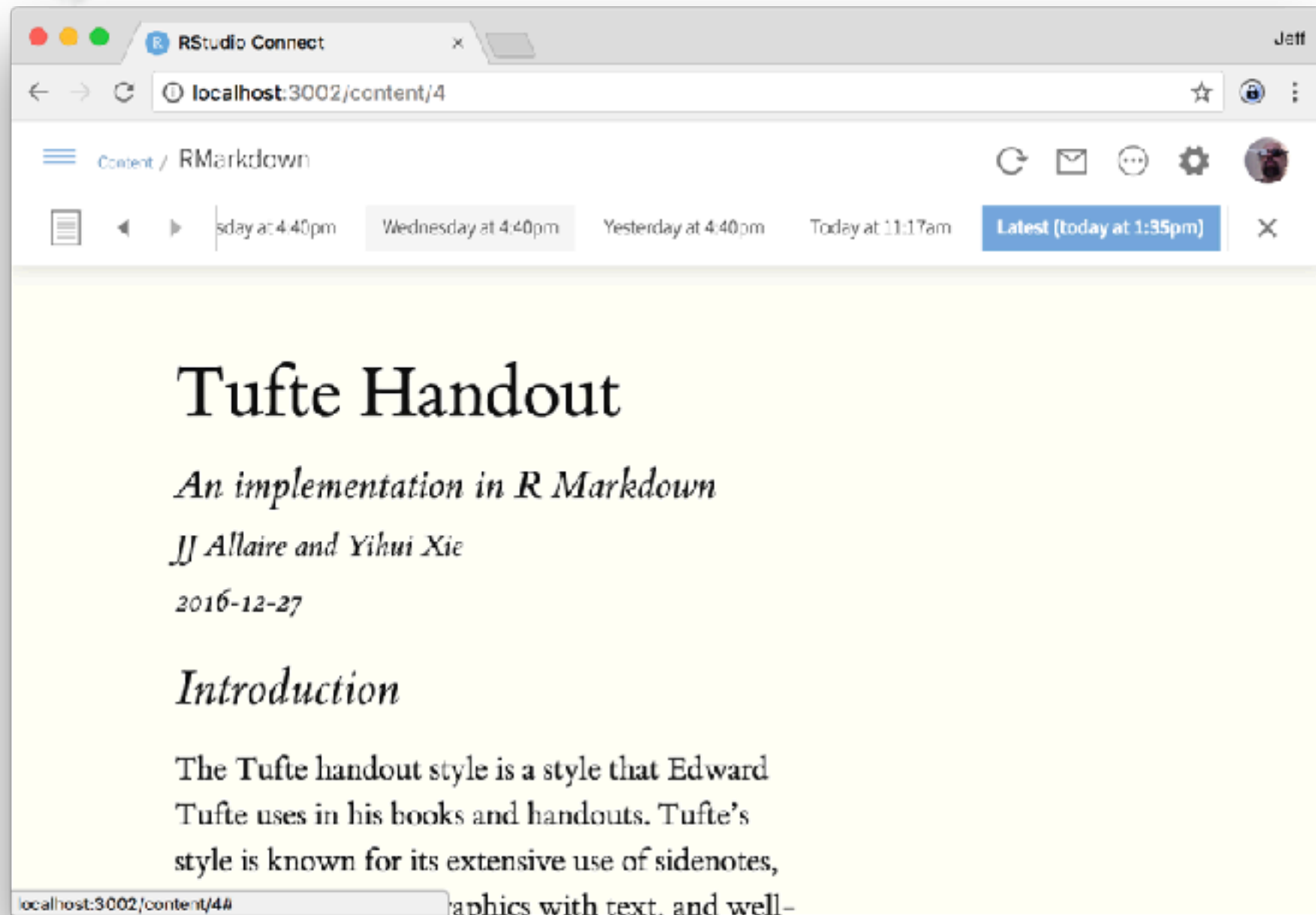- Kerberos
- Content collaboration & continuous integration

# Upcoming…

- View old versions of a report
- API to control RStudio Connect
- In-UI Management of Environment Variables
- SUSE support
- App-specific event logs

# Questions?

- Download & 45-day free trial: http://rstd.io/rsc

- Admin Guide: http://rstd.io/rsc-admin

- Slides: http://rstd.io/earl2017

@trestleJeff