# Introduction to GitHub:

## A Crash Course for Social Science Researchers

Brian Heseung Kim

University of Virginia

May 3, 2021

# Lesson Outline

- **Setting the Stage**
  - Why GitHub? Why Now?
- **Frameworks and Principles**
  - Git, GitHub, and GitHub Desktop, Oh My
- **Getting Started**
  - Level 1: GitHub as Google Drive but Worse
  - Level 2: GitHub as a Dissemination Tool
  - Level 3: GitHub as a Version Control Platform
  - Level 4: GitHub as a Facilitator of Complex Collaboration
- **Next Steps**

# Setting the Stage

UNIVERSITY *of* VIRGINIA | EdPolicyWorks

# Why GitHub? Why Now?

# Why GitHub? Why Now?

- In today's education research space:

# Why GitHub? Why Now?

- In today's education research space:
  - Collaboration across methods, languages, and disciplines is becoming more and more common

UNIVERSITY *of* VIRGINIA | EdPolicyWorks

# Why GitHub? Why Now?

- In today's education research space:
  - Collaboration across methods, languages, and disciplines is becoming more and more common
  - Analyses are becoming more complex and multi-stage

UNIVERSITY *of* VIRGINIA | EdPolicyWorks

# Why GitHub? Why Now?

- In today's education research space:
  - Collaboration across methods, languages, and disciplines is becoming more and more common
  - Analyses are becoming more complex and multi-stage
  - Entire analyses can change on the back of one line of code

UNIVERSITY *of* VIRGINIA | EdPolicyWorks

# Why GitHub? Why Now?

- In today's education research space:
  - Collaboration across methods, languages, and disciplines is becoming more and more common
  - Analyses are becoming more complex and multi-stage
  - Entire analyses can change on the back of one line of code
  - Research is increasingly conducted by multiple analysts building on the same codebase

UNIVERSITY of VIRGINIA | EdPolicyWorks

# Why GitHub? Why Now?

- In today's education research space:
  - Collaboration across methods, languages, and disciplines is becoming more and more common
  - Analyses are becoming more complex and multi-stage
  - Entire analyses can change on the back of one line of code
  - Research is increasingly conducted by multiple analysts building on the same codebase
  - Researchers can benefit from the existing code of other researchers, even in unrelated projects

UNIVERSITY of VIRGINIA | EdPolicyWorks

# Why GitHub? Why Now?

- In today's education research space:
  - Collaboration across methods, languages, and disciplines is becoming more and more common
  - Analyses are becoming more complex and multi-stage
  - Entire analyses can change on the back of one line of code
  - Research is increasingly conducted by multiple analysts building on the same codebase
  - Researchers can benefit from the existing code of other researchers, even in unrelated projects

=

Frictions and dangers in analytic work abound

UNIVERSITY of VIRGINIA | EdPolicyWorks

# Frameworks and Principles

UNIVERSITY of VIRGINIA | EdPolicyWorks

# Git, GitHub, and GitHub Desktop

# Git, GitHub, and GitHub Desktop

UNIVERSITY of VIRGINIA | EdPolicyWorks

# Git, GitHub, and GitHub Desktop

- For now, you can think of Git, GitHub, and GitHub Desktop roughly as:
  - GitHub = Box
  - GitHub Desktop = Box Sync
  - Git = Crazy back-end magic (JavaScript, Python, SQL)

# Git, GitHub, and GitHub Desktop

- For now, you can think of Git, GitHub, and GitHub Desktop roughly as:
  - GitHub = Box
  - GitHub Desktop = Box Sync
  - Git = Crazy back-end magic (JavaScript, Python, SQL)
- And just as there are many cloud storage services…

UNIVERSITY of VIRGINIA    EdPolicyWorks

# Git, GitHub, and GitHub Desktop

- For now, you can think of Git, GitHub, and GitHub Desktop roughly as:
  - GitHub = Box
  - GitHub Desktop = Box Sync
  - Git = Crazy back-end magic (JavaScript, Python, SQL)
- And just as there are many cloud storage services…

| Version Control Systems | File Repository Services | Desktop Clients |
| --- | --- | --- |
| Git | GitHub | GitHub Desktop |
| Subversion | BitBucket | SourceTree |
| … | … | … |

University of Virginia | EdPolicyWorks

# Getting Started

UNIVERSITY *of* VIRGINIA | EdPolicyWorks

# Level 1: GitHub as Google Drive but Worse

- GitHub can work as a file storage solution for code, documentation, and smaller datasets
- Can do this manually via GitHub
- Can do this semi-automatically via GitHub Desktop

UNIVERSITY *of* VIRGINIA | **EdPolicyWorks**

# Level 1: GitHub as Google Drive but Worse

- GitHub can work as a file storage solution for code, documentation, and smaller datasets
- Can do this manually via GitHub
- Can do this semi-automatically via GitHub Desktop

- **Quick Demo**: COVID Dashboard

UNIVERSITY *of* VIRGINIA | **EdPolicyWorks**

# Level 2: GitHub as a Dissemination Tool

- Level 1 plus:
  - Can share public "repositories" (repos) easily and freely
  - Can keep it updated at all times
  - Can write simple "read-me's" for public use
  - Looks fancier than it is

UNIVERSITY *of* VIRGINIA | EdPolicyWorks

# Level 2: GitHub as a Dissemination Tool

- Level 1 plus:
    - Can share public "repositories" (repos) easily and freely
    - Can keep it updated at all times
    - Can write simple "read-me's" for public use
    - Looks fancier than it is

- **Quick Demo**: Education Deserts w/ Dan

# Level 3: GitHub as a Version Control Platform

- Level 2 plus:
    - Can be edited simultaneously by multiple users w/o conflict
    - Easily identify who changed what, and when
    - Easily revert erroneous or catastrophic changes to code
    - Force yourself to write good notes on each update

UNIVERSITY of VIRGINIA | EdPolicyWorks

# Level 3: GitHub as a Version Control Platform

- Level 2 plus:
  - Can be edited simultaneously by multiple users w/o conflict
  - Easily identify who changed what, and when
  - Easily revert erroneous or catastrophic changes to code
  - Force yourself to write good notes on each update

- **Quick Demo**: N2FL Text Analysis

# Level 4: GitHub as a Facilitator of Complex Collaboration

- Level 3 plus:
  - Can formally request changes, document bugs, and keep track of future changes needed
  - Can incorporate suggestions *from* total strangers, or suggest your own changes *to* total strangers ("open source!")
  - Can allow for the same codebase to exist in many forms at once, be many things to different people
  - Automate codebase updates


- Quick Demo: [Parttree Package](#)

UNIVERSITY *of* VIRGINIA | EdPolicyWorks

# Next Steps

University of Virginia | EdPolicyWorks

# Resources to Explore

- [GitHub's Hello World tutorial](#)
- [Similar tutorial from HubSpot](#)
- [Similar tutorial, but using more of the back-end command line](#)
- Just play around with it!

UNIVERSITY *of* VIRGINIA | **EdPolicyWorks**

# Any Questions?

- **Setting the Stage**
  - Why GitHub? Why Now?
- **Frameworks and Principles**
  - Git, GitHub, and GitHub Desktop, Oh My
- **Getting Started**
  - Level 1: GitHub as Google Drive but Worse
  - Level 2: GitHub as a Dissemination Tool
  - Level 3: GitHub as a Version Control Platform
  - Level 4: GitHub as a Facilitator of Complex Collaboration
- **Next Steps**

UNIVERSITY *of* VIRGINIA | EdPolicyWorks

# Thank you!

## Brian Heseung Kim

@brhkim

brian.kim
@virginia.edu

brhkim.com

@brhkim

File   Edit   View   Repository   Actions   Tools   Help

employment_pathways_vccs   n2fl_nlp   cc_job_recommendations

Commit   Pull   Push   Fetch   Branch   Merge   Stash   Discard   Tag   Git-flow   Remote   Terminal   Explorer   Settings

WORKSPACE
- File Status
- History
- Search

BRANCHES
- master

TAGS

REMOTES

STASHES

Pending files, sorted by file status

**Staged files**    Unstage All   Unstage Selected

- 00_n2fl_nlp_core.do
- 04_process_analytic_data.do
- 06b_forest_accuracy_analysis.do

**Unstaged files**    Stage All   Stage Selected

- logs/n2fl_nlp_log_15Apr2021.log
- logs/n2fl_nlp_log_16Apr2021.log

Search

00_n2fl_nlp_core.do

Hunk 1 : Lines 57-63    Unstage hunk

```
57   57          local switch_topic_modeling = 0
58   58          local switch_analysis = 0
59   59          local switch_forest_accuracy_analysis = 0
-         local switch_sentiment_masking_comparisons = 0
+    60   local switch_sentiment_masking = 0
61   61          local switch_interaction_descriptives = 0
62   62          local switch_sentiment_graphs = 0
63   63          local switch_topic_modeling_graphs = 0
```

Hunk 2 : Lines 81-87    Unstage hunk

```
81   81          rscript using "${scripts}/02_advisor_flagging.R"
82   82      }
83   83
-         //SENTIMENT ANALYSIS SCRIPTS ARE TO BE RUN HERE (03 scripts)
+    84   //SENTIMENT ANALYSIS SCRIPTS ARE TO BE RUN HERE (03 scripts)
85   85
86   86      if `switch_process_analytic_data' == 1 {
87   87          do "${scripts}/04_process_analytic_data.do"
```

Hunk 3 : Lines 103-109    Unstage hunk

```
103  103          do "${scripts}/06b_forest_accuracy_analysis.do"
104  104      }
105  105
-         if `switch_sentiment_masking_comparisons' == 1 {
+    106  if `switch_sentiment_masking' == 1 {
107  107          do "${scripts}/06c_sentiment_masking_comparisons.do"
108  108      }
109  109
```

04_process_analytic_data.do

Hunk 1 : Lines 96-103    Unstage hunk

brhkim <brhkim@gmail.com>

Commit options...

Create automated process to create individual qual-coding validation datasets

☐ Push changes immediately to origin/master

Commit