

Data 607 Project 3: What Are the Most Valuable Data Science Skills?

Group members: Sean Connin, Ethan Haley, Daniel Moscoe, Zachary Safir.

Proposed Approach:

- Our project workflow includes the following:
 - Establish shared Github Repository for code, files, and related materials (complete)
 - Establish Slack channel and Google Drive to share information (complete) - we are meeting regularly via Zoom to debrief and identify next steps as well as assignments
 - Develop R code to scrape job listings and descriptions from Indeed (complete)
 - Develop code to identify and extract “skills” identifiers from job descriptions (in progress)
 - Clean data and convert to Tidy format (in progress)
 - Develop and populate relational database (in progress)
 - Export local database to RDB if time permits
 - Prepare analyses, visualizations, and final report

Collaboration tools

- MySQL database -- we are currently experimenting with various ways to share access during the development phase.
- Google Drive, for non-technical documents like this proposal;
- Slack, for asynchronous team communication;
- Zoom;
- GitHub, for shared R code and final product. We have already established a repository and individual branches for each collaborator.

Data sources

Job listings hosted by Indeed - we have a working script that returns search results as well as full job descriptions. While we have run into encoding errors associated with some search results, the data we have procured to date is sufficient to proceed.

Database design

Please see the following page for a tentative entity-relationship diagram. We expect to make modifications to this design as we progress. Time permitting, we will host the database on a remote server such as AWS.

