Data
Science
Skills

# Team Tidy

Alec McCabe

Chinedu Onyeka

Cliff Lee

Preston Peck

Santiago Torres

# Agenda

| Section | Team Member | Time |
|---|---|---|
| Approach<br>• Tools<br>• Assumptions | Santiago | 1 minute |
| Data Collection<br>• Web Scraping<br>• Persistent Storage | Alec & Cliff | 3 minutes |
| Data Transformation<br>• EMSI | Preston | 1 minute |
| Data Analysis | Chinedu & Preston | 2 minutes |
| Conclusion | Santiago | 1 minute |

# Approach

# Approach

For DATA 607 Project 3, all teams must use data to answer the question, "Which are the most valued data science skills?" Consider your work as an exploration; there is not necessarily a "right answer."

Gather job postings and survey what the job market finds valuable

# Collaboration Tools

Communication
- Slack
- Zoom
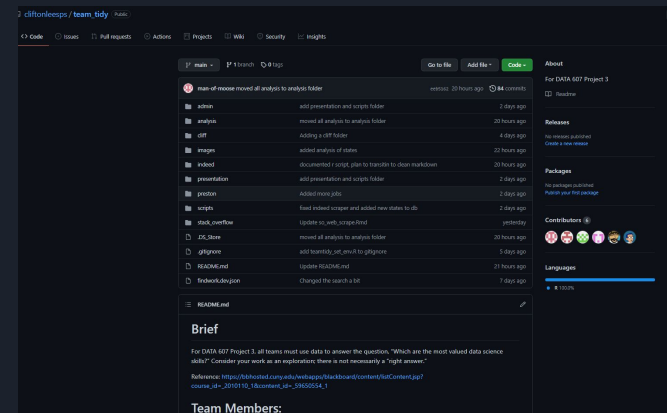
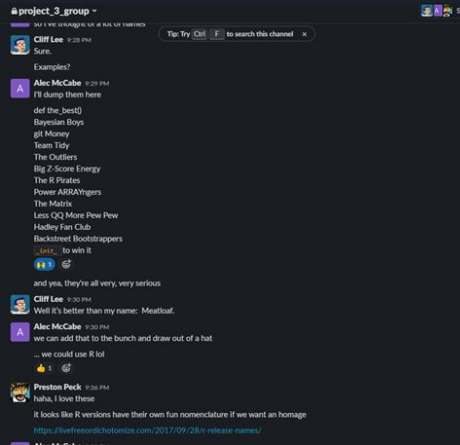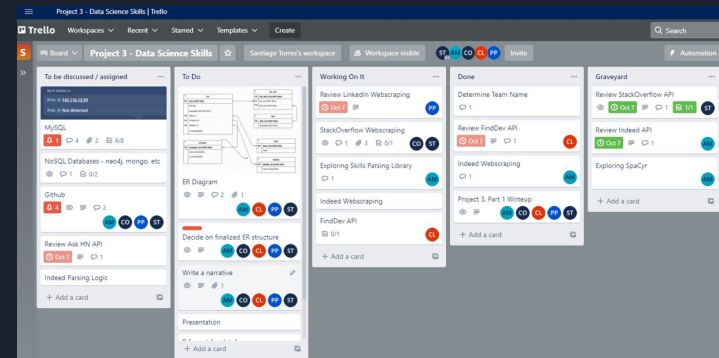Project Documentation
- Google Drive

Project Management
- Trello

Code Sharing
- GitHub
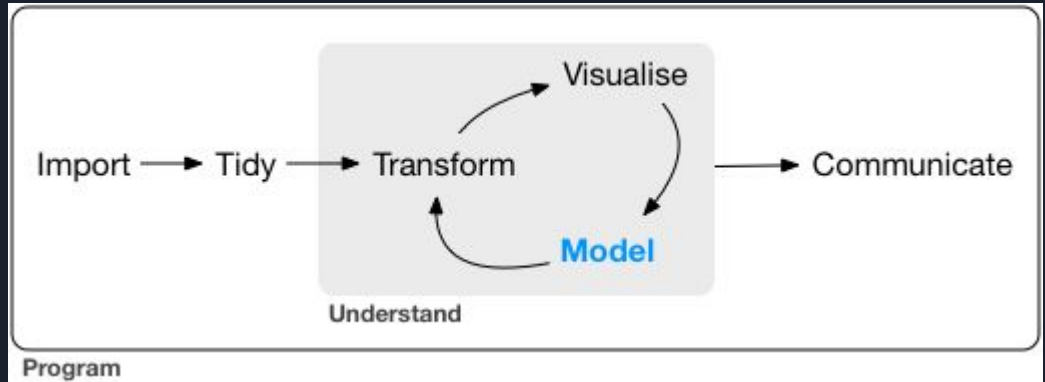
https://github.com/cliftonleesps/team_tidy

# Assumptions

We'd see skills that fall along equally on each part of the model:
- Import / Tidy / Transform skills
- Visualization skills
- Modeling skills
- Communication skills

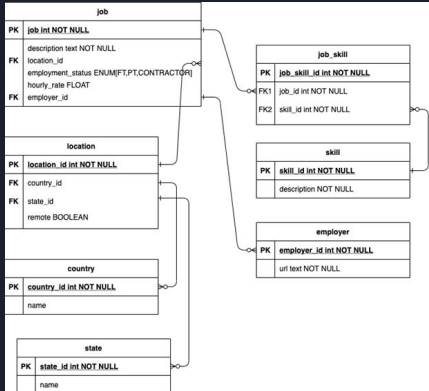Hadley's Data Science Model

# Data Collection

# Persistent Storage



## ER Diagram



| Job Title | State | Company | Original Source | Description | Type |
|-----------|-------|---------|-----------------|-------------|------|
| IT Data Scientist | California | JPL/NASA | indeed | Basic Math | Soft Skill |
| | | | | Communications | Soft Skill |
| | | | | Creative Thinking | Soft Skill |
| | | | | Creativity | Soft Skill |
| | | | | Curiosity | Soft Skill |
| | | | | Innovation | Soft Skill |
| | | | | Management | Soft Skill |
| | | | | Mentorship | Soft Skill |
| | | | | Presentations | Soft Skill |
| | | | | Problem Solving | Soft Skill |
| | | | | Professionalism | Soft Skill |
| | | | | Resourcefulness | Soft Skill |
| | | | | Self Starter | Soft Skill |
| | | | | Self-Awareness | Soft Skill |
| | | | | Agile Methodology | Hard Skill |
| | | | | Analytics | Hard Skill |
| | | | | Application Specific Integrated Circuits | Hard Skill |
| | | | | Auditing | Hard Skill |
| | | | | Business Operations | Hard Skill |
| | | | | Business Process | Hard Skill |
| | | | | Computer Science | Hard Skill |
| | | | | Data Analysis | Hard Skill |
| | | | | Data Integration | Hard Skill |
| | | | | Data Mining | Hard Skill |
| | | | | Data Modeling | Hard Skill |
| | | | | Data Quality | Hard Skill |
| | | | | Data Science | Hard Skill |
| | | | | Data Visualization | Hard Skill |
| | | | | Data Wrangling | Hard Skill |
| | | | | Machine Learning | Hard Skill |

*Sample Data
(Soft Skills in Red)*

# Webscraping



Distribution of Jobs by source

**Web-scraped 220 jobs from:**
- Indeed
- LinkedIn
- StackOverfow

**Technologies used:**
- Xml2
- Rselenium

**Similar structure between websites:**
- Jobcards collected from page
- Iterate through pages

# Data Transformation

# Data Transformation
## Extracting meaning

- With all our **job description text**, how do we extract the key **skill data** effectively and efficiently?
- We could parse manually or hope for predictable patterns to parse by, but this is tedious and inconsistent
- We explored options for Natural Language Processing (NLP) services with Named Entity Recognition (NER) models...

**prodigy**

## Text Classification

Whether you're doing intent detection, information extraction, semantic role labeling or sentiment analysis, Prodigy provides easy, flexible and powerful annotation options. Active learning keeps you efficient even if your classes are heavily imbalanced.
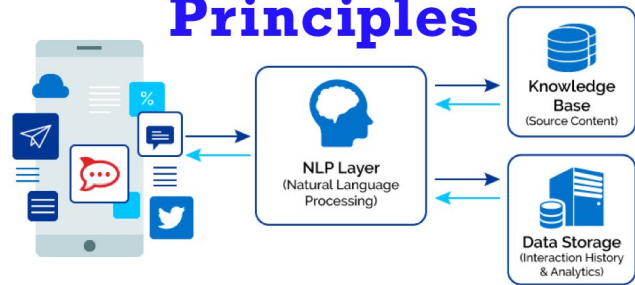
spaCy

| PERSON ₁ | ORG ₂ | PRODUCT ₃ | DATE ₄ |

In a **March 2014** DATE interview , **Apple** ORG designer **Jonathan Ive** PERSON used the **iPhone** PRODUCT as an example of **Apple** ORG 's ethos of creating high - quality , life - changing products .
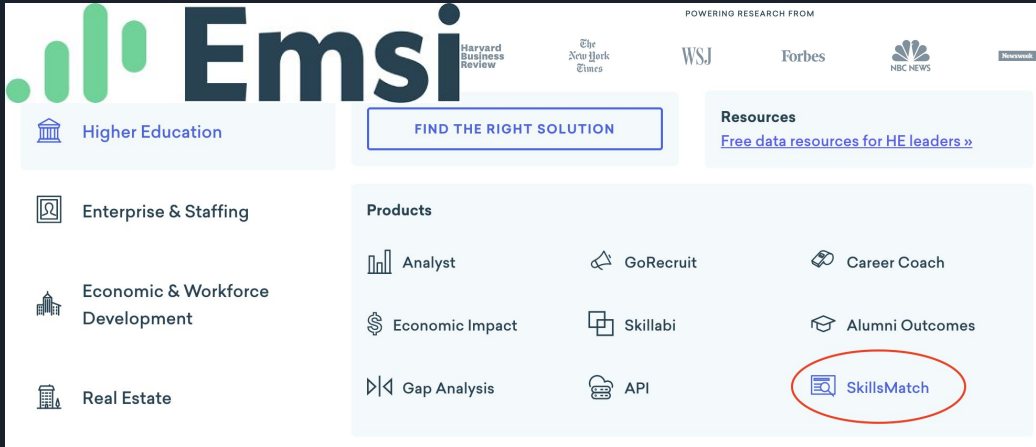
# Natural Language Principles

Knowledge Base
(Source Content)

NLP Layer
(Natural Language Processing)

Data Storage
(Interaction History & Analytics)

# Data Transformation
## NLP using Labor Market Analytics & Economic Data API



POWERING RESEARCH FROM

Harvard Business Review | The New York Times | WSJ | Forbes | NBC NEWS | Newsweek

**Higher Education**
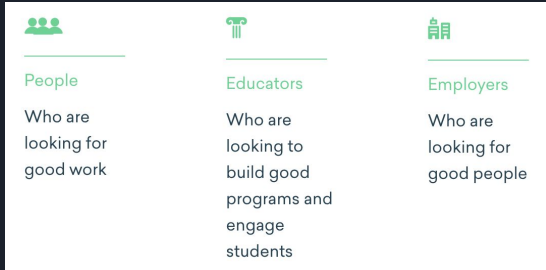
FIND THE RIGHT SOLUTION

**Resources**
Free data resources for HE leaders »

Enterprise & Staffing

**Products**

Economic & Workforce Development

- Analyst
- GoRecruit
- Career Coach
- Economic Impact
- Skillabi
- Alumni Outcomes

Real Estate

- Gap Analysis
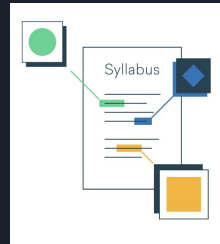- API
- SkillsMatch

**People**
Who are looking for good work

**Educators**
Who are looking to build good programs and engage students

**Employers**
Who are looking for good people

Syllabus

## OpenSkills API
*"30,000+ skills that we've collected from hundreds of millions of job postings, resumes, and online profiles"*

### What is included with free API access?

- Access to every skill in Emsi's Open Skills Library
- Skill Names, unique machine-readable IDs, and types (technical skill, human skill, certification/license)
- Access to every title in Emsi's Open Titles Library
- Title Names, unique machine-readable IDs, and more
- Autocomplete search
- Limited access to skill extraction

# Data Transformation
## Emsi Skills API

**/versions/{version}/extract**

`POST` Extract skills from document

Client ID: 7k_____yuoeqhda9
Secret: vi___RH
Scope: emsi_open

**+**

### URL Parameters

| Name | Description |
|------|-------------|
| version string | The skills classification version. Example: `latest` |

**+**

httr-package: 'httr' makes http easy. →

### Request Body

| Property | Description |
|----------|-------------|
| text string | Document to be used in the skills extraction process |
| confidenceThreshold number | Filter out skills with a confidence value lower than this threshold |
| | Hide details |
| | This is an optional attribute. |
| | Example: `0.6` |
| | Minimum: `0` |
| | Maximum: `1` |
| | Default: `0.5` |

We decided on 0.4

### Response Examples
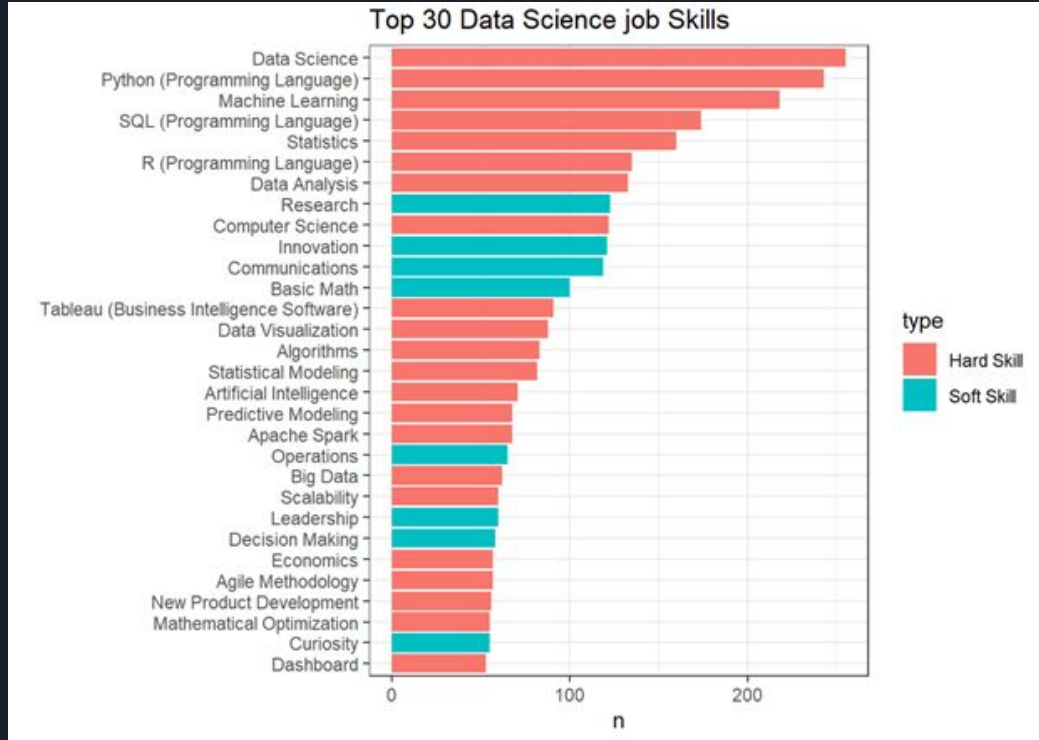
| Property | Description |
|----------|-------------|
| attributions array (objects) | Data attribution information |
| *attributions[]*.name string | Attribution name |
| *attributions[]*.text string | Licensing information |
| data array (objects) | List of extracted skill information |
| *data[]*.confidence number | A number between 0 and 1 representing the confidence of the skill classification |
| *data[]*.skill object | Extracted skill information object |
| *data[]*.skill.type object | Skill type information object |
| *data[]*.skill.type.id string | Skill type ID |
| *data[]*.skill.type.name string | Skill type name |
| *data[]*.skill.id string | Skill ID |
| *data[]*.skill.name string | Skill name |
| *data[]*.skill.tags array (objects) | List of tag information of the skill |
| *data[]*.skill.tags[].key string | Skill tag key |
| *data[]*.skill.tags[].value string | Skill tag value |
| *data[]*.skill.infoUrl string | URL for a publicly accessible web page that includes information about the skill |

# Data Transformation
## Summary

Input: *Job Description*

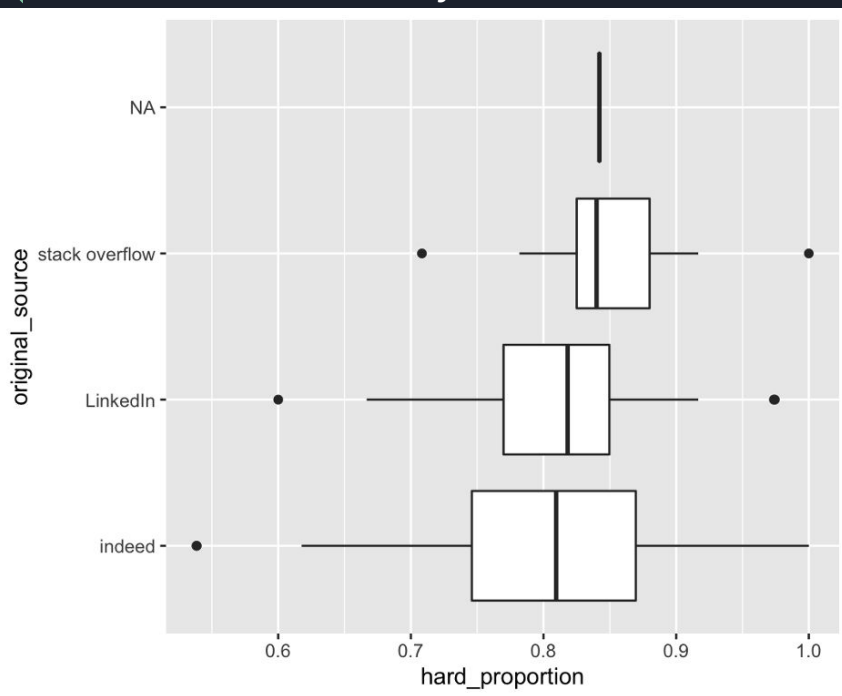Output: *Skills*

# Data Analysis

# Top 30 Data Science Job Skills



Top 30 Data Science job Skills

- >900 Unique Skills
- Hard vs Soft Skills
    - ⅓ of top 12 are soft skills compared to 20% in overall dataset
- Programming skills highly ranked in top 30 skills

# Statistical Modeling

**"Does Stack Overflow's distribution of hard_proportion in jobs have a statistically significant difference than the jobs of other sources"**



1. Visualize proportional differences
2. Stack Overflow displays a larger proportion of hard skills with a shorter spread
3. Is there a statistical difference?
4. Two Sample T-Test -> p-value = 0.1778
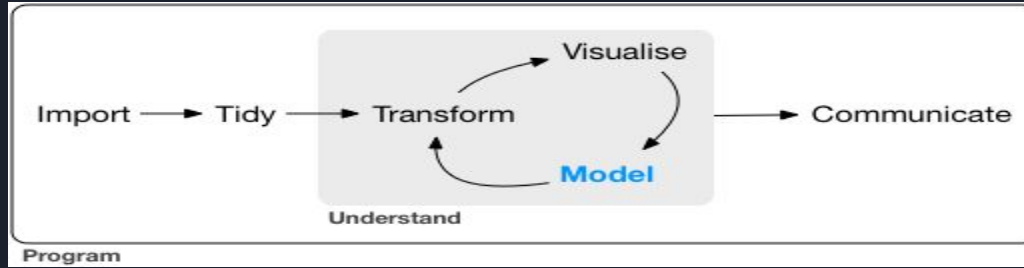5. Fail to reject the null hypothesis

- Two Sample t-test
- data: so_props and other_props
- t = 1.3536, df = 161, p-value = 0.1778
- 95 percent confidence interval:
  a. -0.0168604 0.0903353
- sample estimates:
  a. mean of x mean of y
  b. 0.8475232 0.8107858

# Conclusion

# Most Valuable Data Science Skills
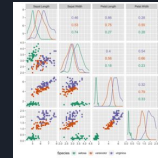
Hadley's Data Science Model



Programming

Statistics & Machine Learning

Communication

Import
Tidy
Transform
Visualise

Model

Communicate

# Next Steps

1. Expand to other jobs / job sources
    a. POC for Data Science - what do other technology jobs look like?
    b. Increase sample sizes
2. Analyzing job titles
    a. We believe we may see different hard to soft skill proportion based on job title (Senior vs Staff Data Scientist)
3. Remove duplicates from sources
    a. Currently no checks in place to enforce uniqueness
4. Build our own model for skills categorization
    a. We depend heavily on EMSI to produce results