

Cloud and Classroom Technology for Data Science



COLLEGE OF
BEHAVIORAL &
SOCIAL SCIENCES

TLTC Experiential Learning Grant

Project Title:

Cloud-based Active Learning for BSOS Statistics and Data Science Courses

- Program-level grant
- Aimed at improving data science education within BSOS

Grant Personnel

Dr. Brian Kim, BSOS/JPSM, Co-Director of Social Data Science major

Dr. Sarah Croco, GVPT, Director of Honors Global Communities LLP

Dr. Candace Turitto, GVPT, Director of Applied Political Analytics

Dr. Jesse Klein, INFO

Dr. Madeline Brown, ANTH

Dr. Taylor Oshan, GEOG

Dr. Zubin Jelveh, CCJS

Graduate Assistants:

Ujjayini Das

Joe Hoskisson

Motivation for the Project

- Data Science is growing in all fields, not just Computer Science or Statistics.
- Needs to be incorporated into social science curricula
- Technology can facilitate incorporating data science into existing courses

Goals for the Project

- Develop a **cloud-computing environment** so that students can get practice programming without needing to install anything
- Develop **modular tutorials** that can be used in many different classes to teach basic concepts
- Build up a **data repository** of social science datasets which can be used in classes to put the concepts in context

JupyterHub and HyFlex for Data Science



COLLEGE OF
BEHAVIORAL &
SOCIAL SCIENCES

Social Data Science major (SDSC)

- Officially launched Fall 2022
- Joint major between INFO and BSOS
- Students learn about data science as well as a social science track
- Focus on application and human context

See <https://sdsc.umd.edu> for more information!

BSOS 233

- Core class for SDSC major
- Designed to be an introductory course on using Python for Data Science

Challenges:

- Need to make programming accessible
- Students from a variety of backgrounds
- Scalability for growth

Addressing the Challenges

- **JupyterHub Cloud-Computing Environment** to make the programming as accessible as possible.
- **HyFlex (Hybrid-Flexible)** format to make learning in the classroom as accessible as possible.

JupyterHub

- Hosted on BSOS OACS servers
- Cloud-based platform
- Everything is done through the browser -- no need to install anything
- Students work completely within the environment, with all work saved in the cloud
- Submit assignments in ELMS by downloading an HTML file

JupyterHub

The screenshot displays the JupyterLab web interface. On the left is a file browser showing a directory structure with files like 'confusion_', 'lab12.ipynb', and 'movement...'. The 'lab12.ipynb' file is selected. The top bar includes a 'Launcher' tab and a 'lab12.ipynb' tab. The main area shows the content of 'lab12.ipynb', which is a Jupyter notebook. The notebook has a title 'Lab 12: Classification' and a text cell explaining the lab's purpose: 'Please complete this lab by providing answers in cells after the question. Use **Code** cells to write and run any code you need to answer the question and **Markdown** cells to write out answers in words. After you are finished with the assignment, remember to download it as an **HTML file** and submit it in **ELMS**.' Below this is a code cell with the following Python code:

```
[1]: import numpy as np
from datascience import *

from sklearn.neighbors import KNeighborsClassifier

from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score, precision_score, recall_score
```

Below the code cell is a text cell stating: 'This lab will cover the main steps for applying machine learning models.'


- Create train and test sets
- Fit the models using train set
- Predict using the test set
- Evaluate models using metrics such as precision and recall
- Make your conclusions

The bottom status bar shows 'Simple' mode, 'Python 3 (ipykernel) | Idle', and 'Mode: Command Ln 1, Col 1 lab12.ipynb'.



Jupyter Notebooks

- Combine executable code and narrative text together



The screenshot displays a Jupyter Notebook interface. At the top, a tab is labeled 'Lab 12: Classification'. Below the tab, a text block provides instructions: 'Please complete this lab by providing answers in cells after the question. Use **Code** cells to write and run any code you need to answer the question and **Markdown** cells to write out answers in words. After you are finished with the assignment, remember to download it as an **HTML file** and submit it in **ELMS**.' Below this text is a code cell containing the following Python code:

```
11: import numpy as np
from datascience import *

from sklearn.neighbors import KNeighborsClassifier

from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score, precision_score, recall_score
```

Below the code cell, a text block states: 'This lab will cover the main steps for applying machine learning models.'

- Create train and test sets
- Fit the models using train set
- Predict using the test set
- Evaluate models using metrics such as precision and recall
- Make your conclusions

Technical Details

- Accessible through <https://bsos233.umd.edu>
- Account needs to be created and approved manually, working on implementing campus authentication
- Assignments, labs, other documents distributed via link on ELMS (pulls from GitHub repository behind the scenes)

HyFlex

Motivation: Students may learn better through different modes of learning

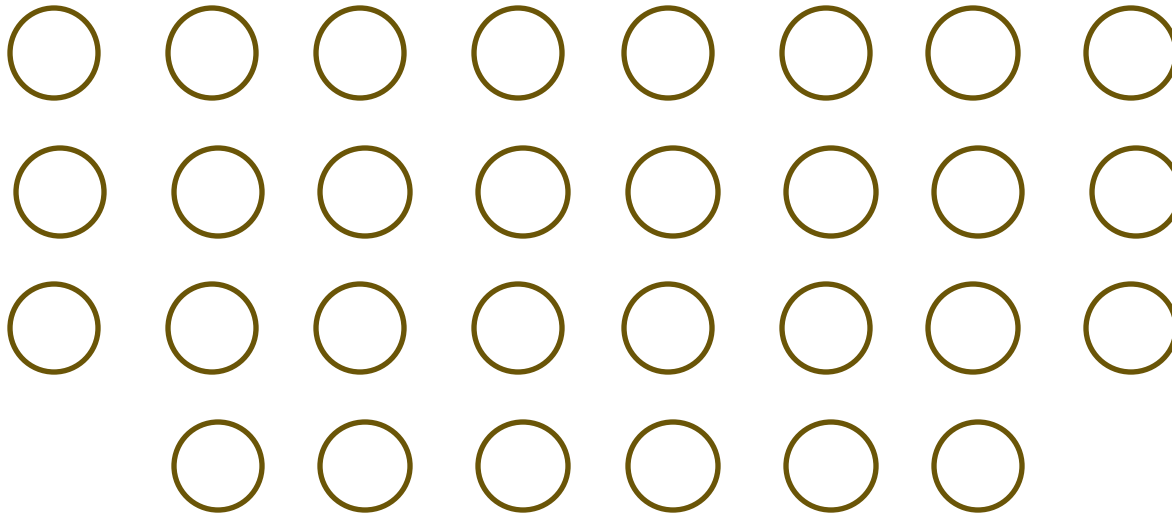
- In-person synchronous
- Online synchronous
- Online asynchronous

HyFlex (Hybrid-Flexible) uses all of these methods of learning and lets the student choose whichever mode they like at any point throughout the semester.

Room Layout

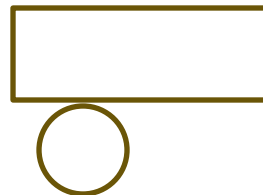
Screen

Screen



Screen

Screen



HyFlex



HyFlex Lessons Learned

- Many students participated **even if they did not turn cameras on** in Zoom.
- Students **took advantage of the flexibility** to choose when to come in person and when to attend online.
- After around Spring Break, **the vast majority of students who attended class did so online.**

Modular Interactive Tutorials for R



COLLEGE OF
BEHAVIORAL &
SOCIAL SCIENCES

Motivation

- Data science is taking over the world; why should we stay behind?
- Programming is hard!! There should be an easier way to deal with it; isn't it a good time to find out how?
- Sounds like a lot of work; where to begin?

Let's try to figure out!

What is the plan?

- Build a set of modular data science tutorials in R, which use real-world data.
- Host them on BSOS cloud servers (with the help of OACS) to easily use as assignments with back-end completion tracking.
- Use them as refresher material in more advanced courses to make sure students meet the requisite level of programming.

What is the benefit here?

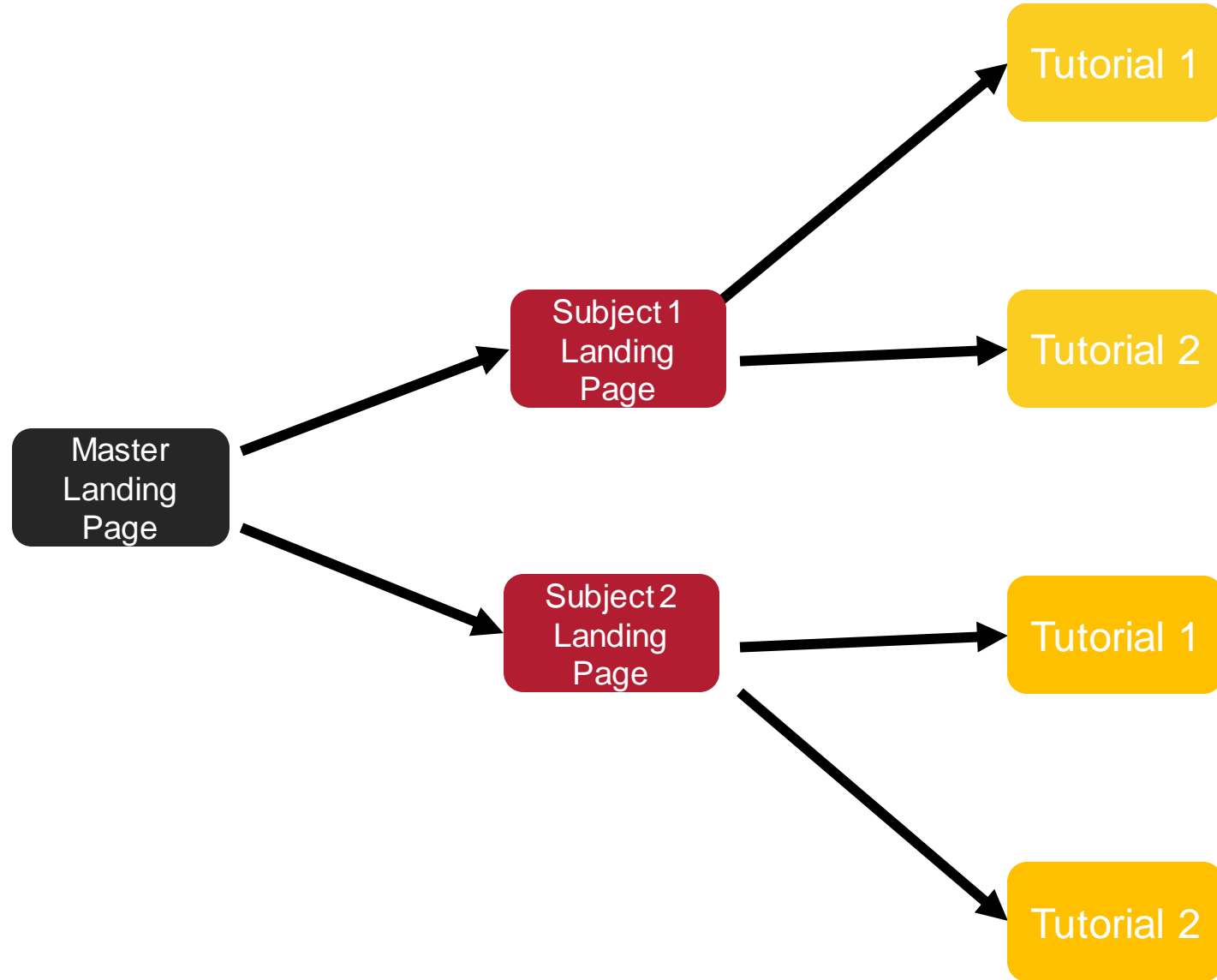
- No need to install software in personal computers, hence less burden for students to get started!
- More streamlined, improved accessibility
- Step-by-step guide on chosen topics based on different levels of courses
- Interactive tutorials, hands on experience in coding
- Exercises with sequential difficulty levels to try out
- Repeated attempts allowed to encourage true learning

Tools up our sleeves!

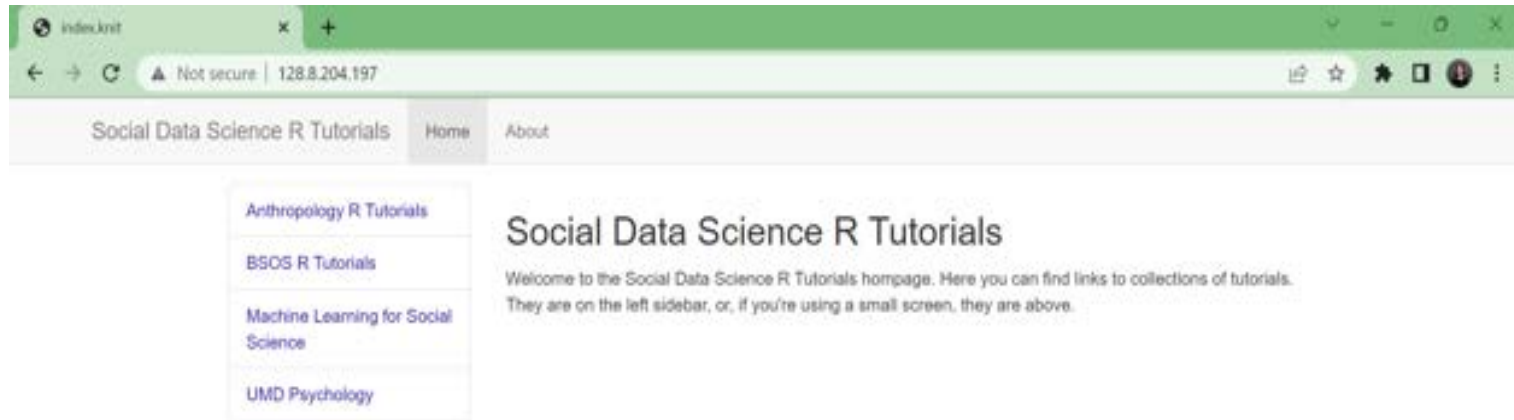
- **Software used**
 - R Markdown
- **Main Packages**
 - learnr
 - gradethis
- **Hosting Server**
 - R Shiny

<https://pkgs.rstudio.com/learnr/articles/exercises.html>

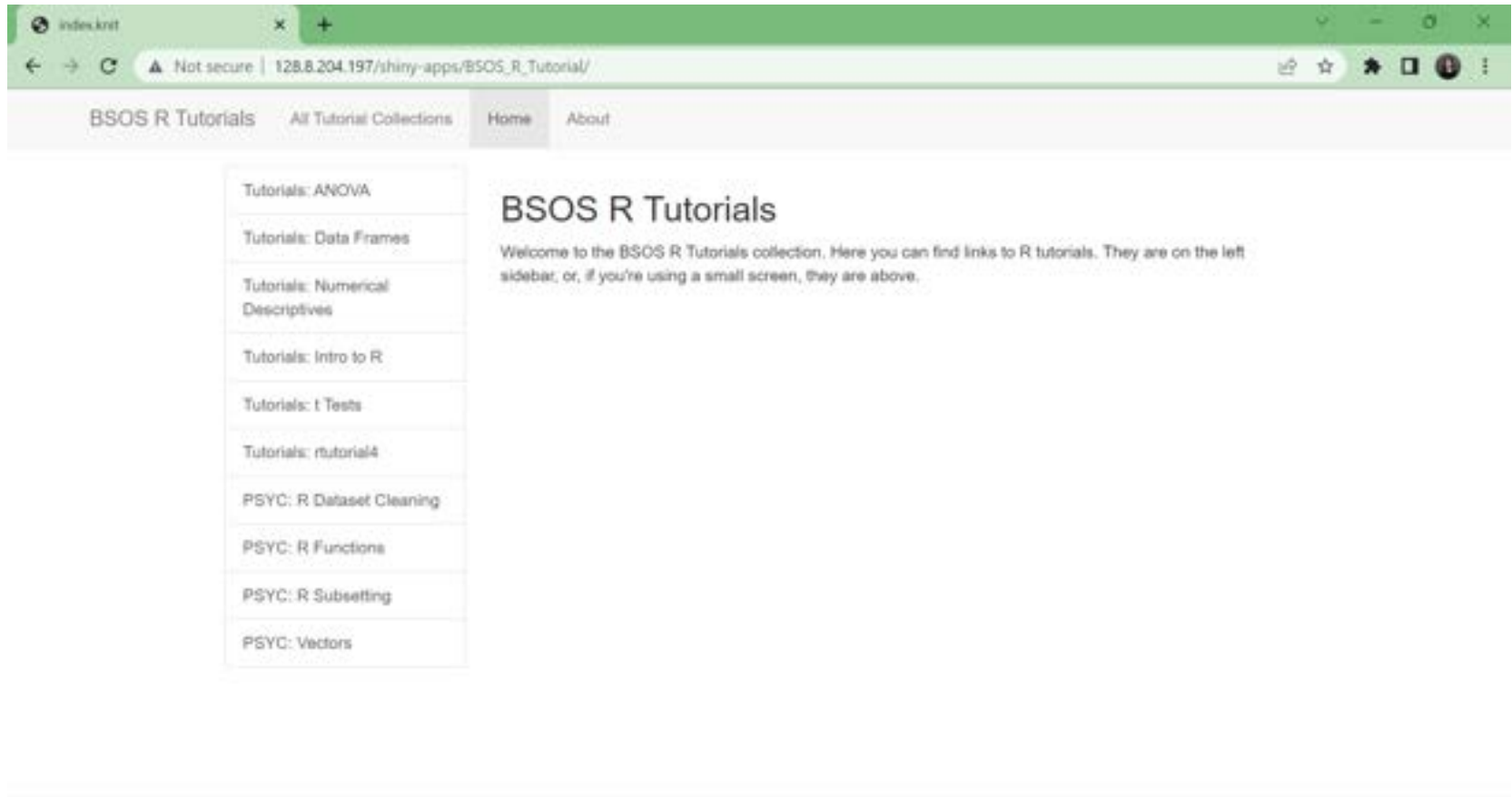
User Interface



Look and Feel (1)



Look and Feel (2)



Look and Feel (3)

Subsetting in R

Section One : Extracting Rows and Columns

Section Two : Extracting & Building Vectors / Dataframes

Section Three : Subsetting

Start Over

Section One : Extracting Rows and Columns

In this module, we will learn some basic data set operations.

- We will start this lesson by using the class dataset we used in Lesson 3.
- Let's start by using the csv class data file.

Opening/loading Data set

```
dat.org <- read.csv("../data_psych/STAT200_Data_Update_4.8.csv")
dat.work <- dat.org
```

#Note: Make sure you are using readxl instead of read.csv
#You will need to redo this part to open the dataset in your folder

Revision of Related Functions

In this section we will talk about a few exploratory functions that can help you explore the data set more. Carefully look at the following code chunk to understand what each function does.

```
#Name of the columns
colnames(dat.work) #See all the variable names in dataset
```

```
## [1] "Participant" "Scores" "Sleep" "Motivation"
## [5] "Note_Type_Name" "Stat.Anxiety" "Missing_Data"
```

Limitations

- Students do not get to see the RStudio environment
- Hides file path, can be difficult for novice students to understand the mechanism
- Cumbersome to program in some functionalities like a stopping rule for attempts or providing feedback properly if someone gets something wrong

Data Curation Fellowship



COLLEGE OF
BEHAVIORAL &
SOCIAL SCIENCES

Data Curation Fellows Selection Process

- Had nearly 40 people apply!
- Applications were primarily from BSOS and Journalism, but we also had some CS students.
- We were happy to see such interest from students at different points in their college careers.

Training the Fellows

- 3 training sessions
 - 2 in person, one online
- Topics included:
 - Attribution
 - Data documentation
 - Codebooks
 - Methodology
 - Data cleaning
 - Missing data
 - Transforming variables
 - Practice with everyone cleaning the same dataset

Special Qualities of Our Program

- Done with an eye towards reproducible science
- Emphasis on DEI data and talking about what this means.
- Getting student curators to think like teachers as they choose dataset and clean them.
 - Nature of the variables
 - What types of analyses could these be used for?

Some Glimpses of Work of the Data Curation Fellows (1): Data on food consumption

cleanedFoodData

File Edit View Insert Format Data Tools Extensions Help

90% 10

	A	B	C	D	E	F	G	H	I	J	K	L	M
	GPA	Gender	breakfast	calories_chicken	calories_day	calories_score	coffee	comfort_food	comfort_food_reasons	comfort_food_cook	comfort_food_cuisine		
1	3.894	1	1	810	3	420	2	chocolate chips Stress, boredom	1	3	1	1	
2	3.3	1	1	720	4	420	2	frozen yogurt, all stress, sadness	1	1	1	1	3
3	3.5	1	1	720	2	420	2	ice cream, stress Stress, boredom	1	1	1	1	2
4	3.8	2	1	810	3	420	2	Chocolate, ice cream Stress, boredom	1	3	1	1	1
5	3.3	1	1	720	3	420	1	ice cream, stress I eat comfort food	1	3	1	1	1
6	3.3	1	1	400	3	315	2	Milk and cheese Stress, anger at	1	3	1	1	1
7	3.5	1	1	810	3	880	2	Pasta, gardenia Boredom	2	1	2	1	1
8	3.894	1	1	720	4	420	2	chocolate pasta sadness, stress	3	3	3	1	1
9	3.4	2	1	400	3	420	2	Cookies, popcorn Sadness, boredom	3	3	3	1	1
10	3.8	1	1	810	3	420	2	ice cream, cake stress, boredom	1	2	1	1	1
11	3.3	2	1	810	3	420	2	Pizza, fruit, apple Friends, emotion	2	3	2	1	1
12	3.8	2	1	400	3	880	2	chips, cookies, I usually only eat	2	3	2	1	1
13	3.4	1	1	720	3	880	1	Chocolate, ice cream Sadness, stress	3	3	3	1	1
14	3.3	2	1	810	3	880	2	Fast food, pizza happiness, satisfaction	7	3	7	1	1
15	3.7	2	1	810	3	420	1	burgers, chips, sadness, depression	3	3	3	1	2
16	3.7	2	2	810	3	420	2	Chili, soup, just Stress and boredom	1	4	1	1	1
17	3.8	1	2	720	3	420	2	chocolate, ice cream boredom	2	3	2	1	1
18	Take out	2 Medical office		4	1	1	3	3	3	1	1	3	1188
19	3.7	2	1	810	3	420	1	Chips, ice cream Boredom, happy	2	3	2	1	1
20	3	2	1	810	4	880	2	Chicken fingers, Boredom	2	4	2	1	1

cleanedFoodData

Codebook

File Edit View Insert Format Data Tools Extensions Help

90% 10

A	B	C
Variables	Description	Type of Variable
GPA		numerical, actual GPA
Gender	1 = Female 2 = Male	
breakfast	1 = cereal option 2 = donut option	the participants are shown the following pictures and asked which one of these pictures they associate with the word "breakfast"
calories_chicken		guessing calories in chicken padina
calories_day	1 - I dont know how many calories I should consume 2 - it is not at all important 3 - it is moderately important 4 - it is very important	Importance of consuming calories per day
comfort_food_reasons	1 - stress 2 - boredom 3 - depression/sadness 4 - hunger 5 - laziness 6 - cold weather 7 - happiness 8- watching tv 9 - none	What are some of the reasons that make you eat comfort food?
	1 - Every day 2 - A number of times a week	

Sheet1

Some Glimpses of Work of the Data Curation Fellows (2): Data from Twitter

codebook ☆ 📁 📄

File Edit View Insert Format Data Tools Extensions Help

90% \$ % .00 123 Default... - 10 + B I

A1 Variables

	A	B	C	D
1	Variables	Type	Type (additional)	Description
2	id	numeric	continuous	id of user account
3	user_name	text/string		account name
4	user_location	text/string		scraped location data (self-reported)
5	user_description	text/string		account bio information
6	user_created	date		date of account creation
7	user_followers	numeric	continuous	number of account followers
8	user_friends	numeric	continuous	number of account friends
9	user_favorites	numeric	continuous	number of favorites
10	user_verified	numeric	binary	binary value of whether account is verified (0 = FALSE; 1 = TRUE)
11	date	date		date of tweet
12	text/string	text/string		text/string of tweet
13	hashtag.1	text/string		hashtag text/string
14	hashtag.2	text/string		hashtag text/string
15	hashtag.3	text/string		hashtag text/string
16	hashtag.4	text/string		hashtag text/string
17	hashtag.5	text/string		hashtag text/string
18	hashtag.6	text/string		hashtag text/string
19	hashtag.7	text/string		hashtag text/string
20	hashtag.8	text/string		hashtag text/string
21	source	text/string		device used to post tweet

+ ≡ Sheet1

nonvoters_cleaned ☆ 📁 📄

File Edit View Insert Format Data Tools Extensions Help

100% \$ % .00 123 Default... - 10 + B I A

A1 Respid

	A	B	C	D	E	F	G	H	I
1	Respid	weight	Q1	Q2_1	Q2_2	Q2_3	Q2_4	Q2_5	Q2
2	470001	0.7516	1	1	1	2	4	1	
3	470002	1.0267	1	1	2	2	3	1	
4	470003	1.0844	1	1	1	2	2	1	
5	470007	0.6817	1	1	1	1	3	1	
6	480008	0.991	1	1	1	-1	1	1	
7	480009	1.0591	1	3	2	3	4	1	
8	480010	1.1512	1	1	1	2	3	1	
9	470008	1.0174	1	1	1	2	2	1	
10	470010	0.8184	1	1	1	1	3	1	
11	470011	1.1853	1	1	1	2	1	1	
12	470013	0.9517	1	1	2	2	2	2	
13	470016	1.374	1	1	1	1	4	1	
14	480020	0.8438	1	3	4	2	2	4	
15	470017	1.2745	1	1	1	3	3	1	
16	470018	0.9742	1	1	1	1	1	1	
17	470019	0.9893	1	1	1	1	1	1	
18	470022	0.8473	1	1	1	1	1	1	

+ ≡ nonvoters_cleaned



Some Glimpses of Work of the Data Curation Fellows (3): Salaries for people in STEM fields

STEM Salaries Codebook

File Edit View Insert Format Data Tools Extensions Help

100% Default

AI Name

	A	B	C	D	E	F	G	H	I	J	K
1	Name	Min Value	Max Value	Description	Variable Type	Notes					
7	company	1,042 unique responses		What company they work at	free_response	"Might need cleaning, there are a lot of companies that have sub-sections"					
8	level	2,592 unique responses		What level the observation is at	free_response	"Doesn't make a lot of sense, some responses are letters, some are numbers, some are combinations of letters and numbers"					
9	title	15 unique responses		Role title	categorical	Business Analyst, Data Scientist, Hardware Engineer, Human Resources, Management Consultant, Marketing					
10	totalyearlycompensation	10,000	4,980,000	Total yearly compensation	internal						
11	location	1,048 unique responses		Job location	categorical	"Structured 'city state' for locations in the US. For outside the US, it's structured 'city, [local], country' (in some cases, the local is missing)"					
12	yearsexperience	0	69	Years of experience	internal	"There are some weirdly specific decimals which is interesting but not necessarily wrong (someone said 15.5 years)"					
13	yearsatcompany	0	69	Years of experience at said company	internal	"Again there are some really specific decimals"					
14	tag	2,818 unique responses		tag	free_response	"A lot of things are repeated (for ex. with an extra 'X') and some are irrelevant (hashtags and question marks)"					
15	basepay	0	1,859,870	Base salary	internal	"I want to know if these salaries are all in the same currency"					
16	stockgrantsvalue	0	2,800,000	Stock grant value	internal						
17	bonus	0	1,000,000	Bonus	internal						
18	gender	4 unique responses		What gender they identify as	free_response	"Male, Female, NA (could assign numbers), one person wrote 'Title: Senior Software Engineer' in the gender field"					
19	otherdetails	12,769 unique responses		Other details	free_response	"Really random, can't really be cleaned or used for statistical analysis"					
20	cityid	0	47,506	City ID	internal	"Is ID number internal or categorical? I'm not sure if it's possible for a city to have an ID of 0 so maybe people use 1 for no city"					
21	rowNumber	0	881	row number	internal	"1-881"					
22	race	8 unique responses		What race they identify as	categorical	Asian, Black, Hispanic, White, Two Or More, NA					
23	education	8 unique responses		What education they have fulfilled	categorical	Highschool, Some College, Bachelor's Degree, Master's Degree, PhD, NA					
24	timestamp			When the data was recorded							

Sheet1 4/25 Update

stem (4/25 update)

File Edit View Insert Format Data Tools Extensions Help

100% Default

AI row_number

	A	B	C	D	E	F	G	H	I	J	K
1	row_number	order_rank	timestamp	company	level	title	totalyearlycompensation	location	city	nation	yearsexperience
2	70138	52474	5/1/21 8:12	Amazon	L4	Software Engineer	114000	Aachen, NRW	Aachen	Germany	8
3	22027	17072	2/8/20 8:44	Amazon	SOE II	Software Engineer	150000	Aachen, NRW	Aachen	Germany	3
4	11236	10485	8/16/18 8:42	Amazon	L3	Software Engineer	112000	Aachen, NRW	Aachen	Germany	5
5	88850	60817	8/12/14 4:57	Uber	L8	Software Engineer	418000	Aarhus, AR	Aarhus	Denmark	10
6	90007	37173	12/4/20 9:08	Uber	L5a	Software Engineer	370000	Aarhus, AR	Aarhus	Denmark	4
7	16297	12848	10/15/18 15:27	Uber	Software Engineer	Software Engineer	200000	Aarhus, AR	Aarhus	Denmark	4
8	14890	11742	8/18/18 8:52	Uber	L4	Software Engineer	200000	Aarhus, AR	Aarhus	Denmark	3
9	14720	11623	9/14/19 13:41	Uber	3	Software Engineer	120000	Aarhus, AR	Aarhus	Denmark	3
10	89008	11629	5/8/21 13:27	Parsons	L1	Software Engineer	78000	Aberdeen Proving Ground, MD	Aberdeen Proving Ground	United States	2
11	82913	48089	3/23/21 8:07	GrabHub	Senior Engineer	Software Engineer	170000	Arlington, MD	Arlington	United States	8
12	72480	54274	8/4/21 23:55	Andela	ES	Software Engineer	88000	Accra, AA	Accra	Ghana	8
13	32824	24586	7/8/20 3:01	PeCo	Director	Management Director	180000	Adelaide, SA	Adelaide	Australia	13
14	25220	18296	3/20/20 2:28	QNC Technology	3	Software Engineer	48000	Adelaide, SA	Adelaide	Australia	4
15	18420	12546	10/16/18 2:57	Beijing	L2	Software Engineer	88000	Adelaide, SA	Adelaide	Australia	4
16	38884	25249	8/10/20 15:31	HPE	MD3	Hardware Engineer	88000	Aguadilla, PR	Aguadilla	Puerto Rico	11
17	57122	42440	2/10/21 1:18	Abbott	15	Software Engineer	300000	Alameda, CA	Alameda	United States	21
18	50705	37859	12/9/20 8:12	Abbott	Grade 16	Hardware Engineer	198000	Alameda, CA	Alameda	United States	13
19	38188	27096	8/12/20 10:22	Abbott	Grade 16	Hardware Engineer	200000	Alameda, CA	Alameda	United States	13
20	69756	52188	5/10/21 17:34	Red Ventures	Senior	Software Engineer	175000	Albany, NY	Albany	United States	10
21	66719	48699	4/11/21 15:14	IBM	T1	Software Engineer	171000	Albany, NY	Albany	United States	11

stem424v2

THANK YOU!

Dr. Brian Kim
kimbrian@umd.edu

Ujjayini Das
ujstat@umd.edu

Dr. Sarah Croco
scroco@umd.edu

Joe Hoskisson
jhoskiss@umd.edu