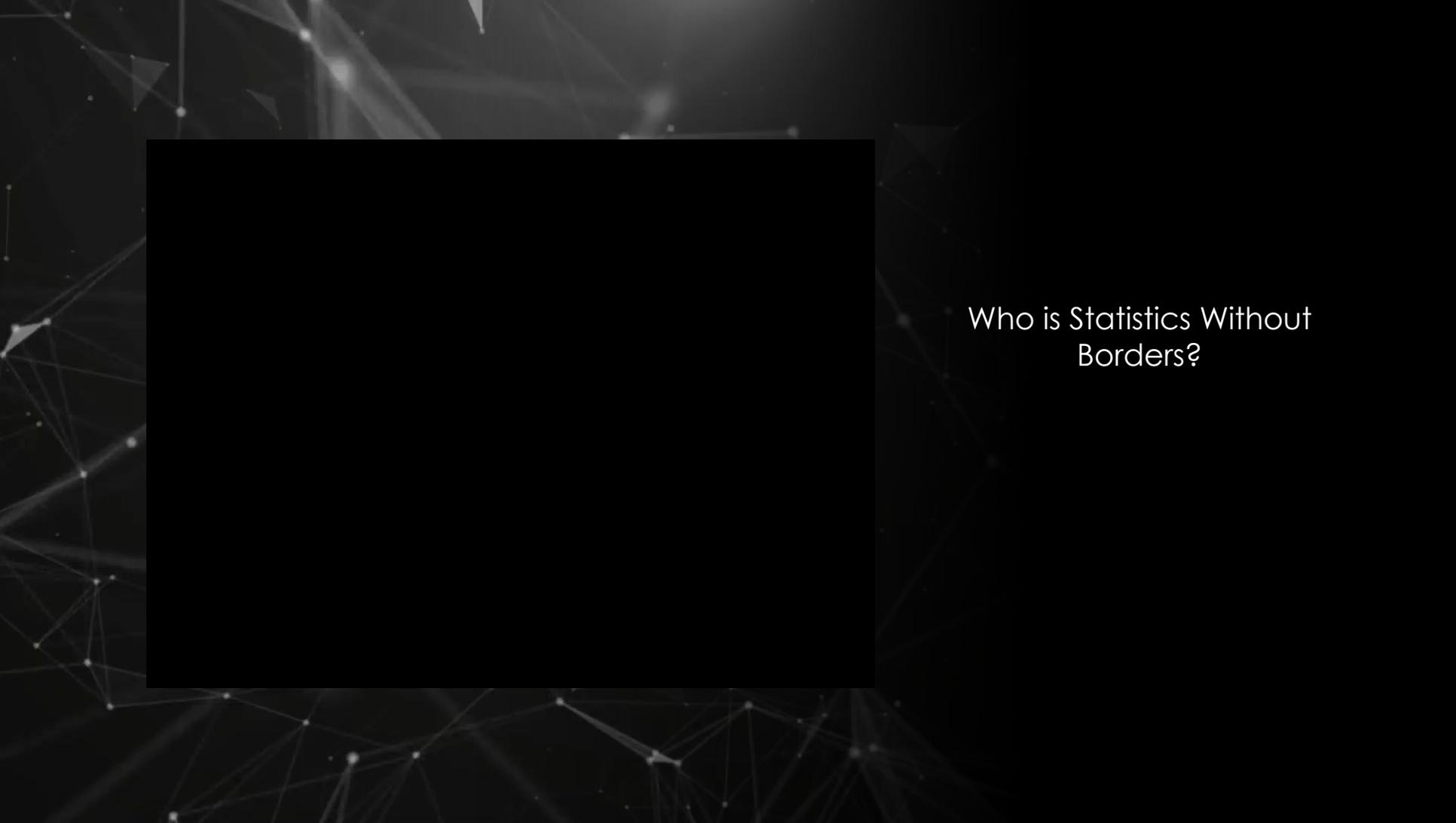




Data Science  
For Social Good  
Statistics Without Borders



Who is Statistics Without  
Borders?

# Statistics Without Borders Helps UNICEF

1 JULY 2016 809 VIEWS NO COMMENT

*Stephanie Eckman, Monica Dashen, Aliou Diouf Mballo, and Robert Johnston*

## Experts Leverage Statistical Methods to Investigate Human Trafficking

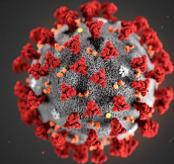
Use of GPS-Enabled Mobile Devices to Conduct Health Surveys: Child Mortality in Sierra Leone

• [Columns, Here's to Your Health](#)

## Haiti after the earthquake Statistics Without Borders

When a major disaster strikes, urgent needs may be food, water, shelter, medicines – and data. Unless you know the numbers of people involved and how their lives have been affected, giving efficient help is impossible. Statistics Without Borders tries to provide the data. The team that worked on a project in Haiti describe one effort.

Our Work



# COVID-19

# California COVID-19 By The Numbers

July 27, 2020

Numbers as of July 26, 2020

## CALIFORNIA COVID-19 SPREAD

# 460,550 (+6,891)

### TOTAL CASES

#### Ages of Confirmed Cases

- 0-17: 41,148
- 18-49: 278,295
- 50-64: 88,800
- 65+: 51,772
- Unknown/Missing: 535

#### Gender of Confirmed Cases

- Female: 230,423
- Male: 227,545
- Unknown/Missing: 2,582

# 8,445 (+29)

### Fatalities

### Hospitalizations\*

Confirmed COVID-19  
**6,935/2,012**  
Hospitalized/in ICU

Suspected COVID-19  
**1,484/209**  
Hospitalized/in ICU

For county-level  
hospital data:  
[bit.ly/hospitalsca](http://bit.ly/hospitalsca)

USA  
**4,280,135**  
TOTAL CASES

CDC | Updated: Jul 28 2020 2:46PM

USA  
**147,672**  
TOTAL DEATHS

CDC | Updated: Jul 28 2020 2:46PM

USA  
**1,306**  
Cases per 100,000  
People

CDC | Updated: Jul 28 2020 2:46PM

### Total Cases by State/Territory

State/Territory	Total Cases	Confirmed	Probable
California	460,550	N/A	N/A
Florida	427,698	N/A	N/A
Texas	385,923	N/A	N/A
New York City*	225,541	220,907	4,634

Your actions **save lives.**

For county-level data:  
[data.chhs.ca.gov](http://data.chhs.ca.gov)  
[covid19.ca.gov](http://covid19.ca.gov)



#### Territories

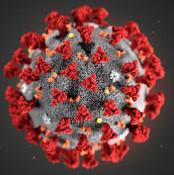
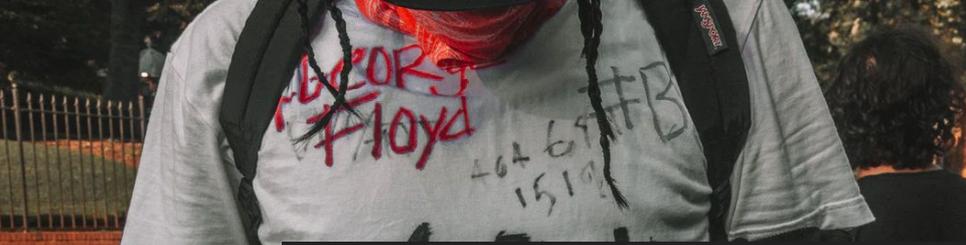
AS FSM GU MP PR PW RMI VI

Netherlands

Pakistan

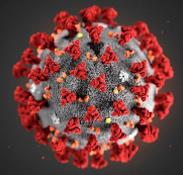
Sweden



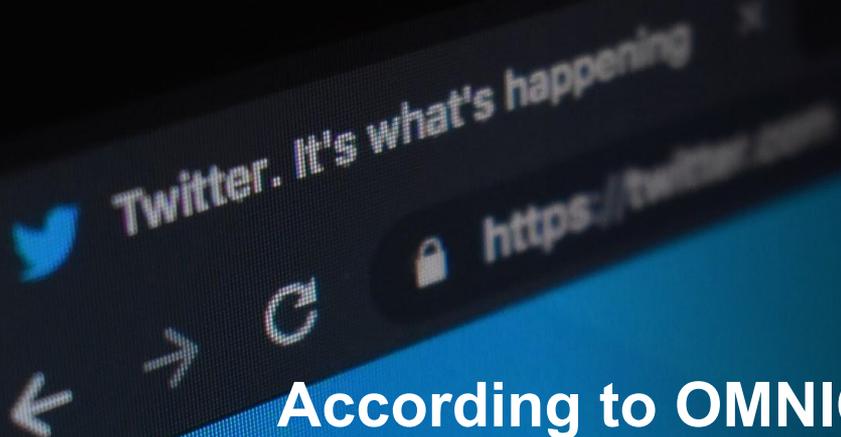


COVID-19





SWB + MCCERT

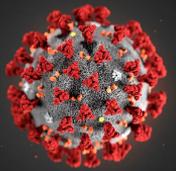


**According to OMNICORE (2020):**

- > 330 million active Twitter users
- > 500 million tweets are posted per day
- 71% of users say they use Twitter to get their news (Pew Research Center, 2019)

**Client Organization: Montgomery County, Community Emergency Response Team**

- MCCERT asked SWB to try this methodology in an independent geographic area around Palo Alto, California
- **Steve Peterson**
  - Virtual Emergency Response Team
  - Steve developed a framework (Peterson et al., 2019) to utilize Twitter data to inform emergency response in the National Capital Region using George Mason University's streaming analytics system, Citizen Helper.



SWB + MCCERT



**3-Month Project**

**8 SWB Volunteers**

**Julia Reid, PCM**

**Keri Wheatley**

**Heli Vora**

**Satyajeet Pradhan**

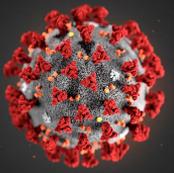
**Rachel Doehr**

**Qingyuan Wang**

**Harshit Sharma**

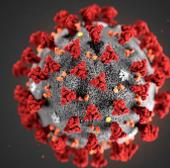
**Lena Lickteig, DQA**

**Collaboration by video, email, phone, text, and chat**



SWB + MCCERT



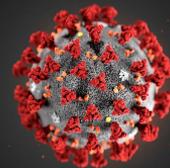


SWB + MCCERT



## Objective:

- Gather tweets and sort by relevance to COVID-19
  - Specific locations of interest
  - Specific terms of interest
    - Keywords associated with:
      - Prevention, Symptoms, etc.

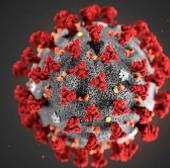


SWB + MCCERT



## The Process

- Establish a flow of targeted tweets
  - Web scraping and data engineering

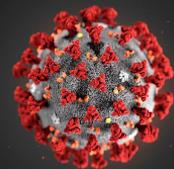


SWB + MCCERT



## The Process

- Establish methods for predicting the relevance of Tweets for emergency response
  - Data wrangling, conditional statements
  - Natural Language Processing, Pre-processing
  - Modeling Approaches: Supervised, Unsupervised



SWB + MCCERT



← Tweet



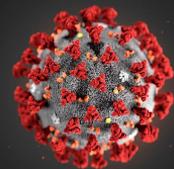
This stresses me out. Why? Because my clinic has a shortage of supplies as well. We also had to lock up supplies because people are stealing them. PPE like gloves & masks are vital! [#seattlecovid19](#)



· Mar 13

A hospital in Seattle area has sent out a note to staff, shared with me, suspending elective surgery and warning that "our local COVID-19 trajectory is likely to be similar to that of Northern Italy." The hospital is down to a four-day supply of gloves.

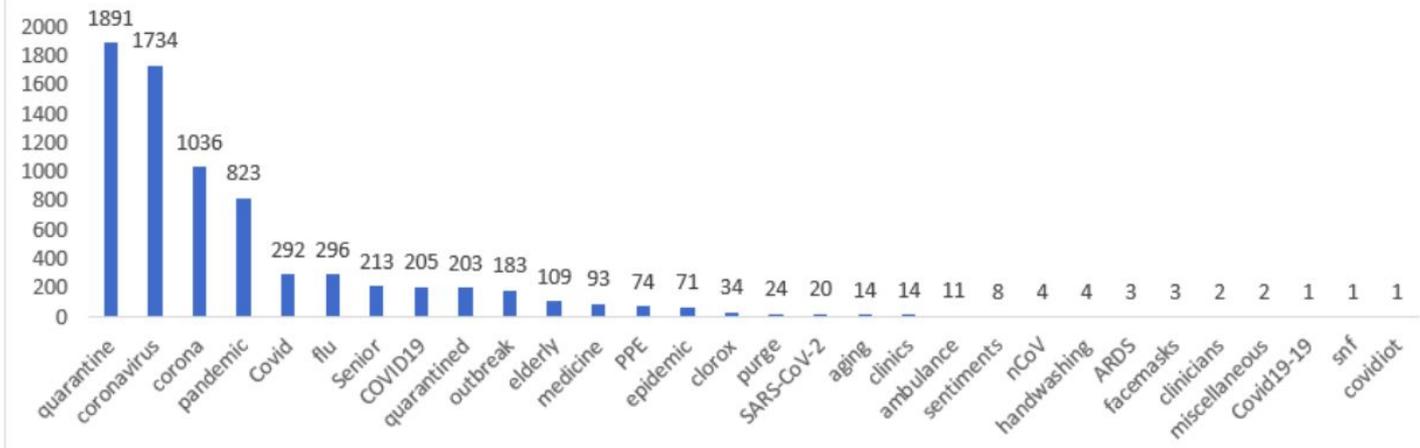
6:46 PM · Mar 13, 2020 · [Twitter for iPhone](#)



SWB + MCCERT

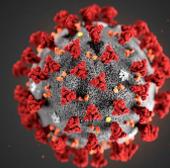
### Case 1: Exact Matches

Key Words Frequency







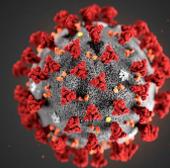


SWB + MCCERT



## Example of Methods and Models Applied

- Word Embeddings: TF-IDF, Word2vec, GLOVE, fastText
- Unsupervised Learning: K-Means Clustering, DBSCAN, Latent Dirichlet Allocation (Topic Modeling)
- Upsampling the minority class: SMOTE
- Transfer Learning: MERS → COVID-19
- Supervised Learning:
  - Naive Bayes, Logistic Regression, GLMNET, Support Vector Machines, ULMFIT, and XGBOOST,



SWB + MCCERT



## Deliverables to MMCERT

- Data acquisition pipeline
- Text preprocessing scripts
- Auditable model pipelines
- A collection of tweets over the course of the beginning of widespread awareness of the COVID-19 epidemic with emergency response relevance predictions

# Data Science For Social Good



Wrangling  
Visualizing  
Modeling  
Explaining  
Communicating



```
import pandas as pd

from pandas.io.json import json_normalize
from sklearn.metrics.pairwise import linear_kernel

pd.set_option('display.max_colwidth', -1)
pd.set_option('display.float_format', lambda x: '%f' % x)
```

```
In [526]: # Normalize JSON data to get the original urls
def normalize_url(data):

    import datetime

    print('normalize_url Start:', datetime.datetime.now())

    data['Ext URL'] = data['Ext URL'].fillna('').map(lambda x: x.strip())

    norm_url = pd.DataFrame()

    for i, row in data.iterrows():
        ext_url = row['Ext URL']
        if ext_url != '':
            ext_url = ext_url.replace("'url': '", "'url": ") \
                .replace("'", 'expanded_url': '", "'expanded_url": ') \
                .replace("'", 'expanded_url': '", "'expanded_url": ') \
                .replace("'", 'display_url': '", "'display_url": ') \
                .replace("'", 'display_url': '", "'display_url": ') \
                .replace("'", 'display_url': '", "'display_url": ') \
                .replace("'", 'indices"', "'indices") \
                .replace("'", 'indices"', "'indices")

            new_row = pd.DataFrame(json_normalize(json.loads(ext_url)))
            new_row['TweetID'] = row['TweetID']
            norm_url = norm_url.append(new_row, sort=False)

    norm_url['indices']
    norm_url.drop_duplicates(inplace=True)

    print('normalize_url End:', datetime.datetime.now())
```



Size	Kind
53.4 MB	comma...values
53.5 MB	comma...values
53.8 MB	comma...values

Share Comments

Lightning bolt icon Ideas Sensitivity

Yes No

M	N	O	P

connected\_graphs

Conditional Formatting Format as Table Cell Styles

J	K	L	M

J	K	L	M

EXPLORER

- OPEN EDITORS
  - parse\_json.py src\data 9+
  - string\_match.py src\data
  - get\_data.py src\data
  - visualize.py src\data
- SWB-COVID-19-TWITTER-REPO
  - misc-resources
  - pipeline
  - src
    - data
    - features
    - models
    - readme U
    - .env
    - .gitignore U
    - README.md



```

src > data > parse_json.py > accuracy_score
27 import gensim
28 from gensim.models import word2vec
29 import nltk
30 from sklearn.decomposition import PCA
31 from sklearn.model_selection import train_test_split
32 from nltk.stem import PorterStemmer
33 from sklearn.linear_model import LogisticRegression
34 from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
35 from sklearn.metrics import accuracy_score, confusion_matrix
36 import matplotlib.pyplot as plt
37 import csv
38 from datetime import datetime
39 from gensim.models import KeyedVectors, Word2Vec
40 import sys
41 from tweepy.streaming import StreamListener
42 import multiprocessing
43 from pandas.io.json import json_normalize
44 import tweepy
45 import nltk
46 from nltk import wordpunct_tokenize
47 import re
48 from nltk.corpus import stopwords
49 import boto3
50 from nltk.tokenize import TweetTokenizer
51 from nltk.stem import WordNetLemmatizer
52 import spacy
53 import string
54 import pandas as pd
55 from nltk.tokenize import word_tokenize
56 import boto3
57 import numpy as np
58 from pathlib import Path
59 from dotenv import find_dotenv, load_dotenv
60 import time
61 import ijson
62 import sqlalchemy
63 import preprocessor as p
64 from gensim.models.doc2vec import TaggedDocument
65 from sqlalchemy import create_engine
66
67 tqdm.pandas(desc="progress-bar")
68
69 class SymptomFinder():
70
71     def __init__(self, client_csv_fname, oov_csv_fname, broadened_client_csv_fname, mers_csv_fname):
72         self.nlp = spacy.load("en_core_web_sm")

```

TERMINAL

```

C:\Users\rachel> cd src\data
C:\Users\rachel> python parse_json.py

```

DEBUG OUTPUT

```

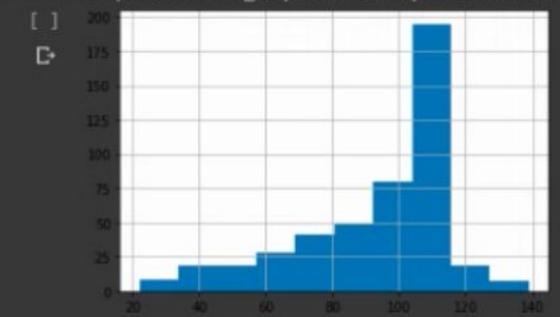
2020-05-21 10:00:00,000 DEBUG: Loading model from C:\Users\rachel\AppData\Local\Microsoft\Windows\Temporary Internet Files\Content.IE5\...

```

Files

Connecting to a runtime to enable file browsing.

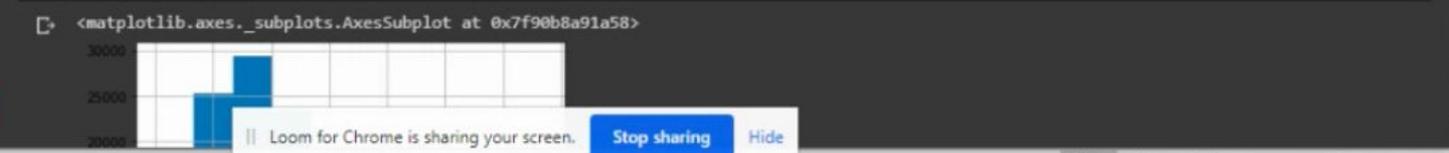
+ Code + Text Reconnect Editing



```
[ ] # checking for shortened tweet_text  
dat[dat._preprocessed_text.apply(lambda x: '\\.\\.' in x)]
```

```
id FINAL_LABEL status_obj text hashtags url expanded_url label_mod preprocessed_text
```

```
[ ] palo_tweets.Processed_Text.apply(len).hist()
```



Loom for Chrome is sharing your screen. Stop sharing Hide



... [X] [Pause] [Checkmark]

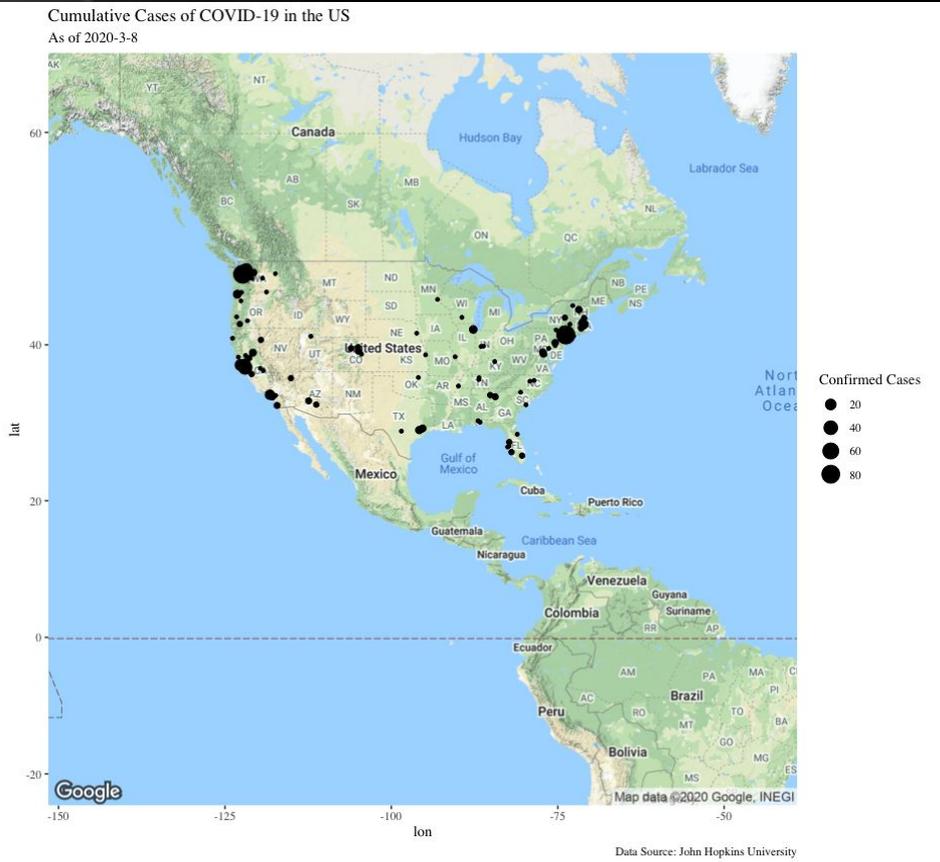
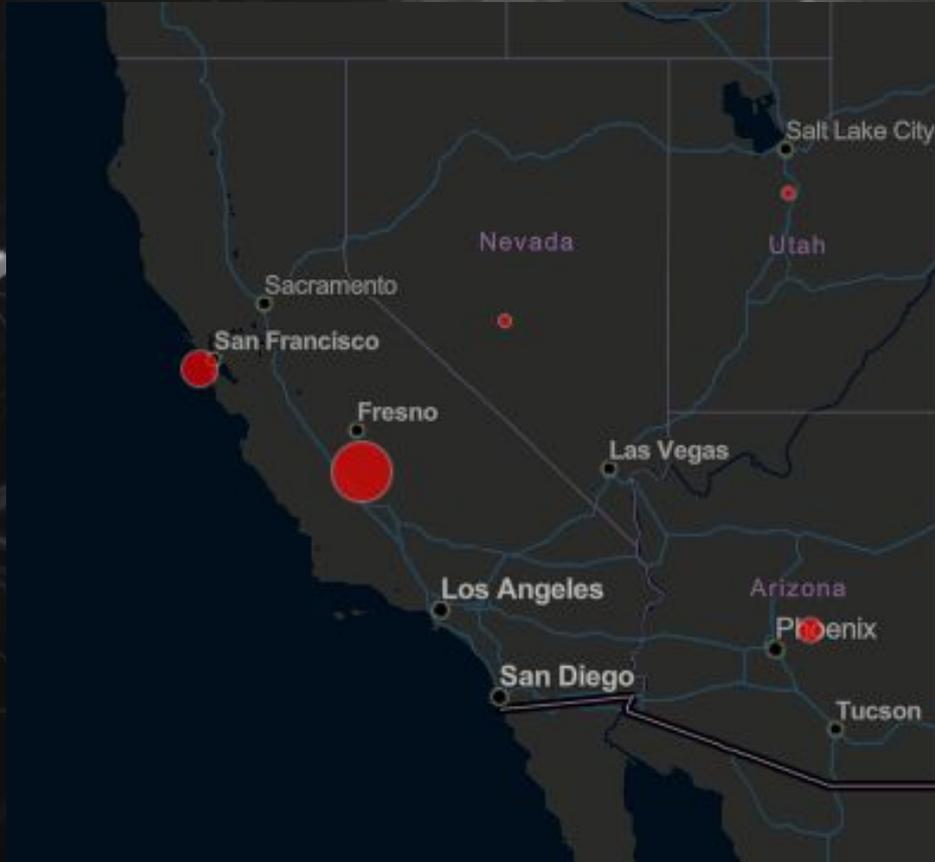
“Not only does data wrangling consume most of an analyst’s workday, it also represents much of the analyst’s professional process: it captures activities like understanding what data is available; **choosing** what data to use and at what level of detail; **understanding** how to **meaningfully combine(ing)** multiple sources of data; and **deciding** how to **distill(ing)** the results to a size and shape that can **drive(ing)** downstream analysis.”

- **Principles of Data Wrangling**  
**Rattenbury et al. (2017)**

“Not only does data wrangling consume most of an analyst’s workday, it also represents much of the analyst’s professional process: it captures activities like **understanding** what data is available; **choosing** what data to use and at what level of detail; **understanding** how to meaningfully combine multiple sources of data; and **deciding** how to distill the results to a size and shape that can drive downstream analysis.”

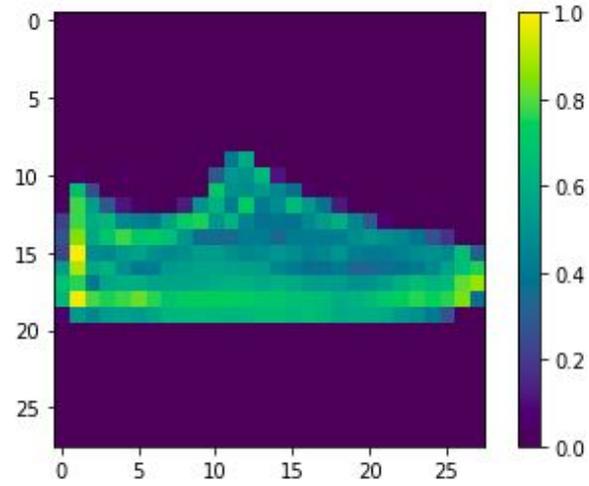
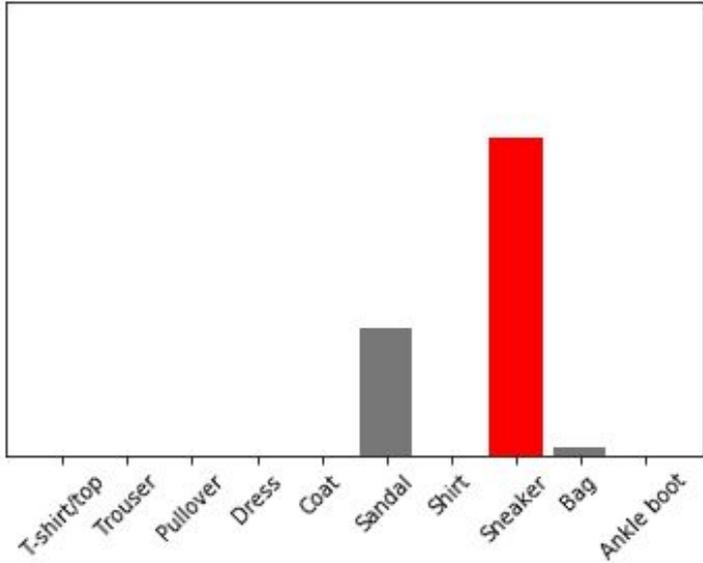
- Principles of Data Wrangling  
Rattenbury et al. (2017)

# JHU Public Dashboard vs. Using JHU Data (2020-3-8)





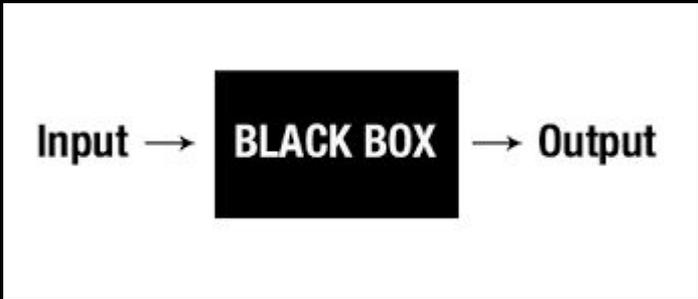
kinestry.io



ML/AI predictions are probabilistic.

# Machine Learning Interpretability (MLI)

More complex models take more work to explain, but **may no longer be** “black boxes”.



# Machine Learning Interpretability (MLI)

- **Goal 1a: Task Performance**
  - **Goal 1b: Understand the model (what's driving predictions?)**
  - **Goal 1c: Privacy, Fairness, and Provide the Right to Explanation**
- 
- **Tools that help:**
  - **Global Variable Importance**
    - What is the weighting of variable contributions to predictions, on average?
    - In NLP: Which words in which contexts contribute most to positive predictions?
  - **Local Variable Importance**
    - What is the weighting of variable contributions to specific observations?
  - **Surrogate Decision Trees**
    - Share a model of the prediction rules by outcome class

# Machine Learning Interpretability (MLI)

## - Sensitivity Analysis

- Vary the inputs; make small changes
- How does this influence predictions?
- What small changes would “push observations (or people) over the threshold”?
- This may inform subsequent iterations in data collection

# Fairness

- “...unfairness and discrimination are pervasive when decisions are being made by humans, which, unfortunately, are not automatically solved, and can even be amplified, when machines are put in control.” - Zhang and Bareinboim (2017)
- “Fairness in machine learning is an emerging topic with the overarching aim to critically assess algorithms (predictive and classification models) whether their results reinforce existing social biases.” - Kozodoi and Varga (2020)

## - General Approaches

### - Disparate Impact Analysis

- ex) Accuracy Parity... [performance metric] by group relative to the reference group

### - Root Cause Analysis

- Do we know whether **protected features** influenced the prediction?



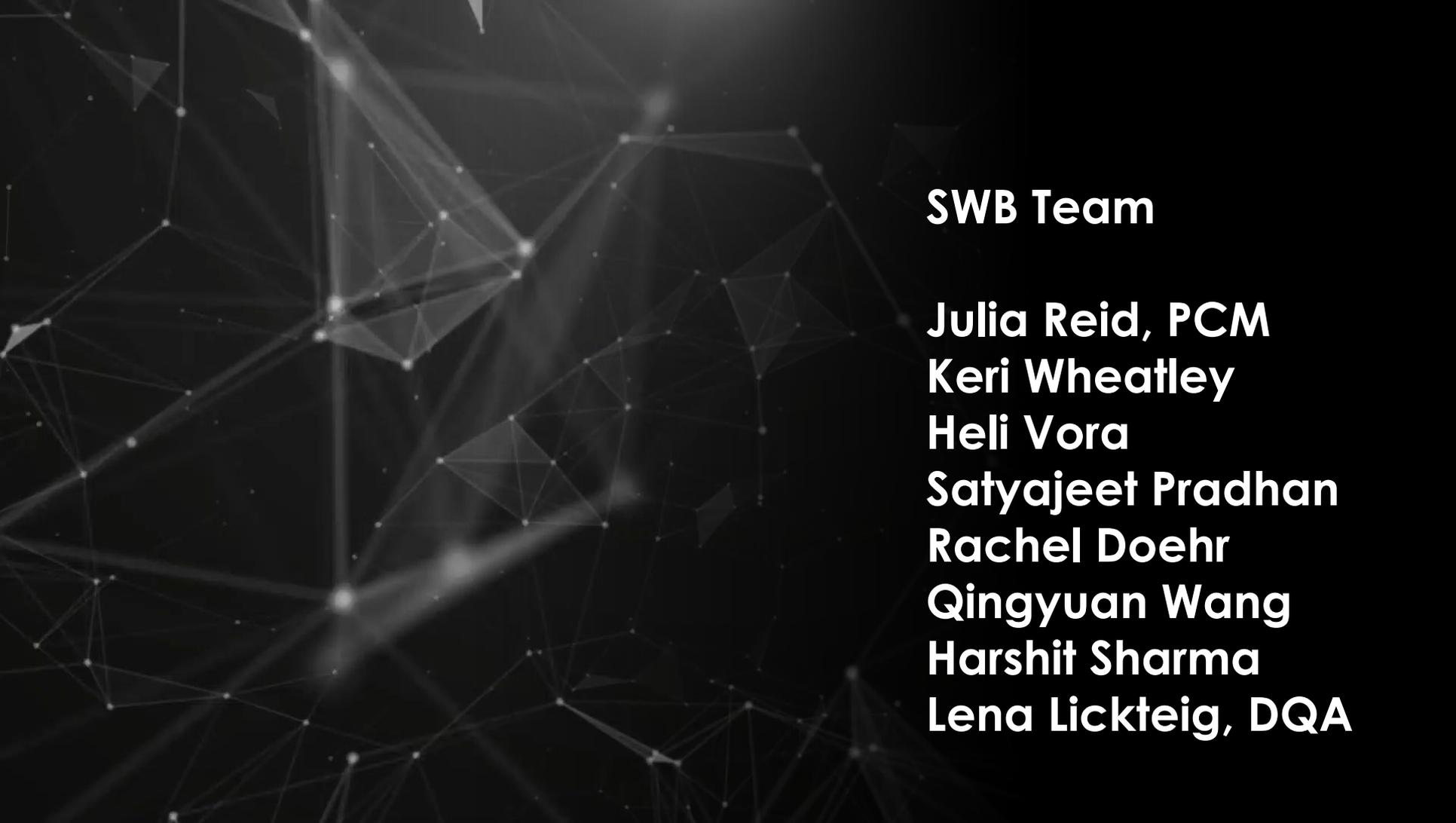
Now that we understand a model, do we trust it?

- What features did our ML or AI learn from?
  - Which of these features should it have learned from?
  - Which of these features **shouldn't** be learned from?
- Is this a model for social good?
  - Who does the model serve?
  - Who **doesn't** the model serve?



Ultimately, our goal is to do data science for **social good**.

This is why solving problems in ways that we can explain and trust is essential.



## **SWB Team**

**Julia Reid, PCM**

**Keri Wheatley**

**Heli Vora**

**Satyajeet Pradhan**

**Rachel Doehr**

**Qingyuan Wang**

**Harshit Sharma**

**Lena Lickteig, DQA**