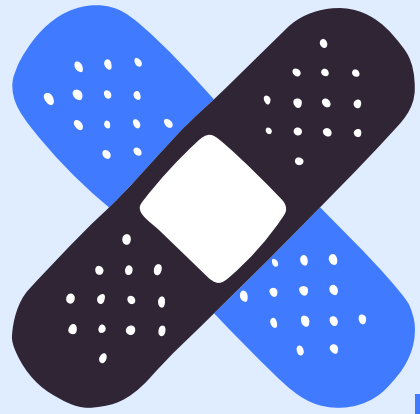# LONG COVID PREDICTION ALGORITHM

By Devin Gee and Juan Pablo Selame Fernandez

# Background Long COVID Prediction

## Long COVID

- COVID symptoms lasting more than 4-12 weeks
- No available diagnostic test or treatment options
- The risk factors for long COVID are poorly understood

## Behavioral Risk Factor Surveillance System

- Annual telephone survey conducted by the CDC for noninstitutionalized adults
- 2022 Data is the first year with COVID-19 data
- Previously used in machine learning projects to predict diabetes

## Previous Work

- Machine learning models have been employed to predict COVID mortality and severity of symptoms
- Other studies focused on more clinical variables such as medications taken and types of providers seen to predict long COVID

## Novelty

- We are specifically using behavioral and demographic data to identify vulnerable populations, and lifestyle risk factors
- Similar studies have focused primarily on chronic diseases such as diabetes and heart disease

# DATA DESCRIPTION

- Largely categorical and ordinal survey data
- 445,132 individuals surveyed total with 100s of questions asked. (Questions varied by participants)
- 110,877 individuals reported that they have contracted COVID
- 86,901 Individuals in the final dataset after cleaning missing values
- 29 features selected for initial dataset. Ranging from demographics (gender, income, race, education), lifestyle (exercise, sleep time), health (self reported mental and physical health), health history (diabetes, asthma, other chronic diseases), and more.

Label: Have an 3 month or longer covid symptoms?
Section Name: Long-term COVID Effects
Core Section Number: 17
Question Number: 2
Column: 266
Type of Variable: Num
SAS Variable Name: COVIDSMP
Question Prologue:
Question: Did you have any symptoms lasting 3 months or longer that you did not have prior to having coronavirus or COVID-19?

| Value | Value Label | Frequency | Percentage | Weighted Percentage |
|---|---|---|---|---|
| 1 | Yes | 26,783 | 21.56 | 21.42 |
| 2 | No—Go to next section | 94,596 | 76.16 | 76.63 |
| 7 | Don't know/Not Sure—Go to next section | 2,710 | 2.18 | 1.87 |
| 9 | Refused—Go to next section | 110 | 0.09 | 0.09 |
| BLANK | Not asked or Missing Notes: Section 17.01, COVIDPOS, is coded 2, 7, 9, or Missing | 320,933 | . | . |

Label: Frequency of Days Now Smoking
Section Name: Tobacco Use
Core Section Number: 12
Question Number: 2
Column: 224
Type of Variable: Num
SAS Variable Name: SMOKDAY2
Question Prologue:
Question: Do you now smoke cigarettes every day, some days, or not at all?

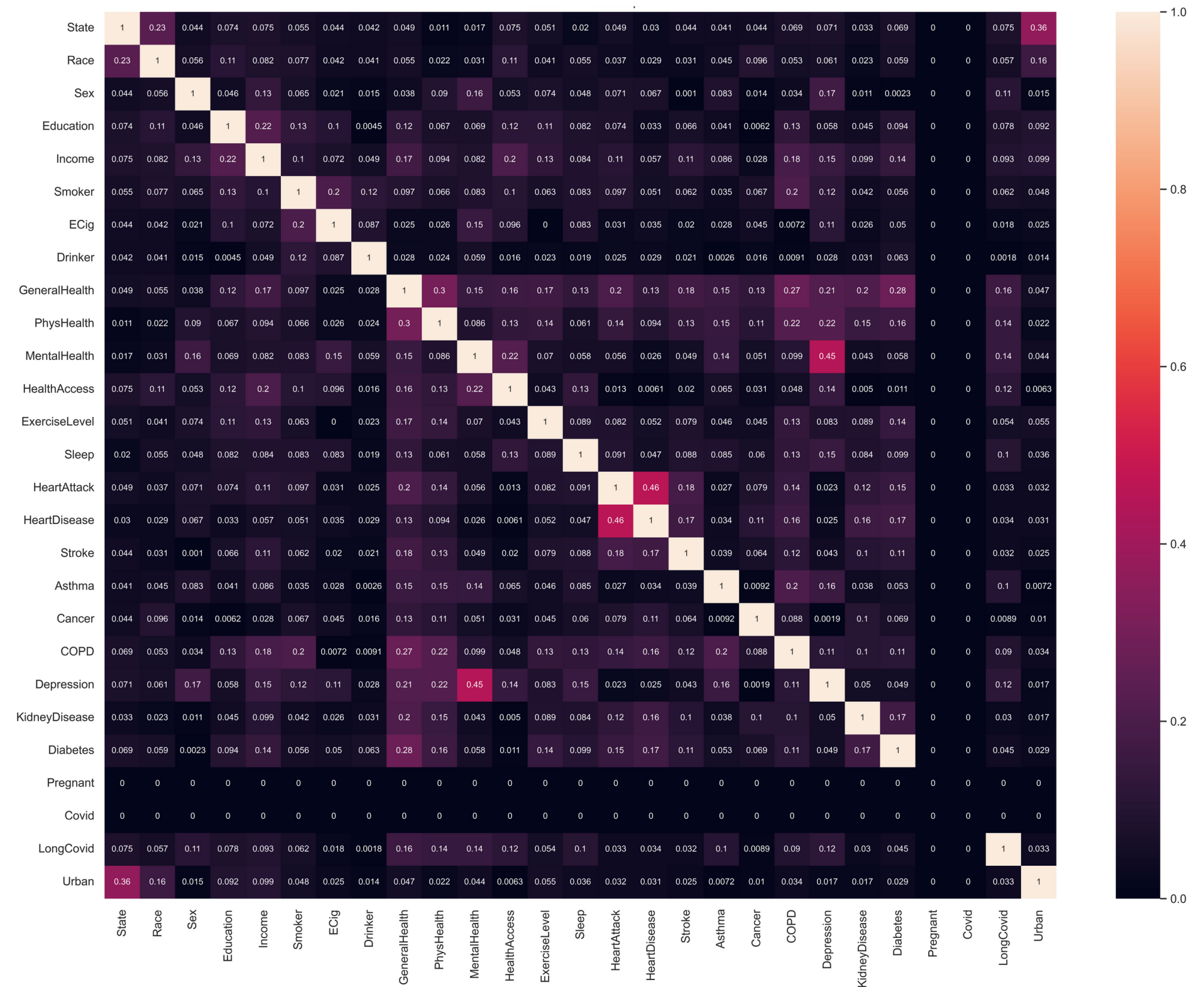| Value | Value Label | Frequency | Percentage | Weighted Percentage |
|---|---|---|---|---|
| 1 | Every day | 36,003 | 21.95 | 24.09 |
| 2 | Some days | 13,938 | 8.50 | 10.54 |
| 3 | Not at all | 113,774 | 69.35 | 65.15 |
| 7 | Don't Know/Not Sure | 165 | 0.10 | 0.12 |
| 9 | Refused | 173 | 0.11 | 0.10 |
| BLANK | Not asked or Missing Notes: Section 12.01, SMOKE100, is coded 2, 7, 9, or Missing | 281,079 | . | . |

| Stroke | Asthma | Cancer | COPD | Depression | KidneyDise | Diabetes |
|---|---|---|---|---|---|---|
| 2 | 2 | 2 | 2 | 2 | 2 | 0 |
| 2 | 1 | 2 | 1 | 2 | 2 | 1 |
| 2 | 2 | 2 | 2 | 2 | 2 | 1 |
| 2 | 2 | 2 | 2 | 2 | 2 | 0 |
| 2 | 2 | 2 | 2 | 2 | 2 | 1 |
| 2 | 2 | 2 | 2 | 2 | 2 | 1 |
| 2 | 2 | 2 | 2 | 2 | 2 | 0 |
| 2 | 2 | 1 | 2 | 2 | 2 | 0 |
| 2 | 2 | 2 | 2 | 2 | 2 | 0 |
| 2 | 2 | 2 | 2 | 2 | 2 | 1 |
| 1 | 2 | 2 | 1 | 1 | 2 | 1 |
| 2 | 2 | 2 | 2 | 2 | 2 | 0 |

19163 Positive for Long COVID (~22%)
67738 Negative for Long COVID
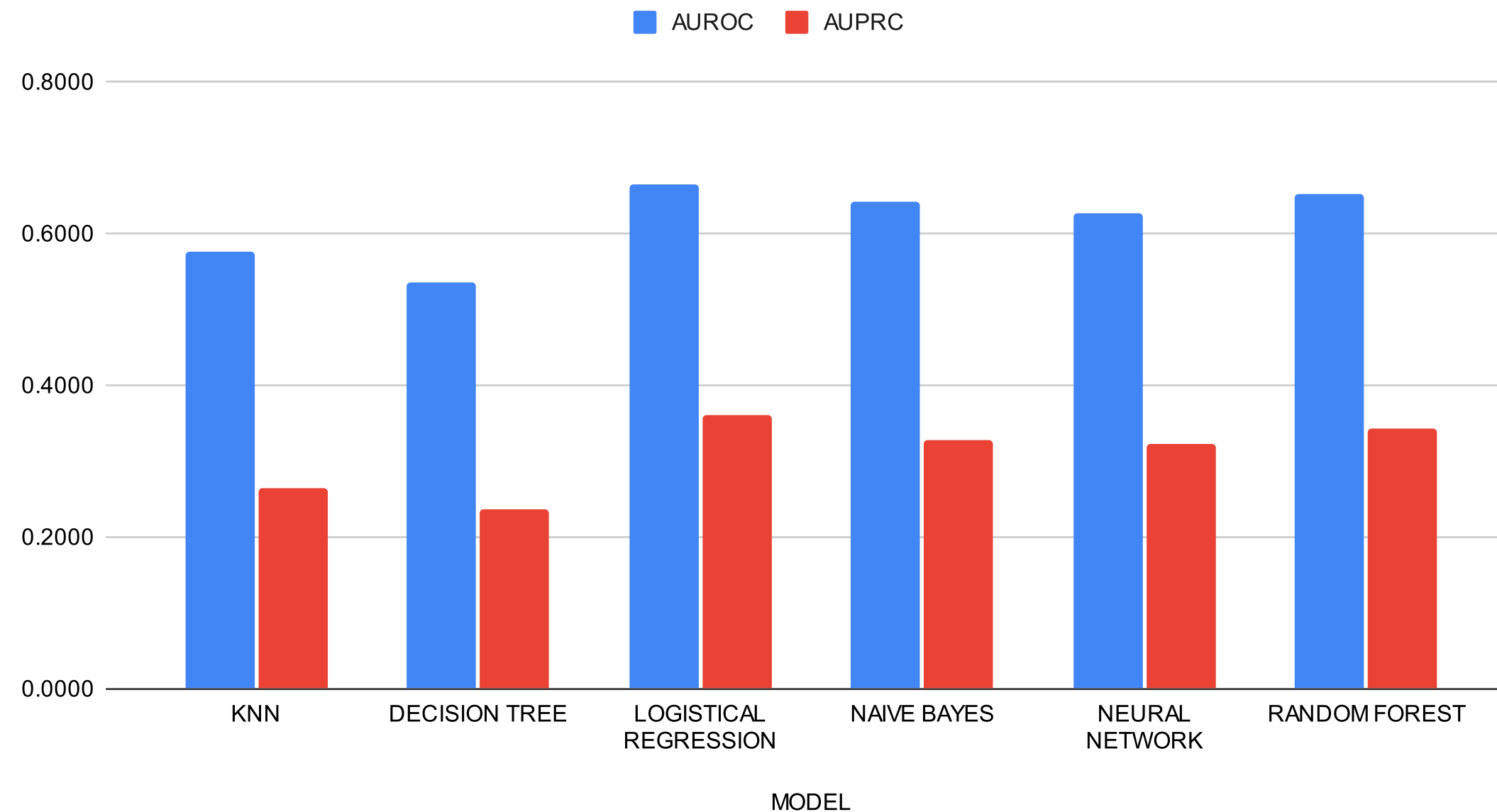Imbalanced Dataset

# FEATURE SELECTION

**Cramer's V Correlation Test**

- Extremely weak correlations filtered from data set (<0.1)
- No features were very highly correlated
- Pregnancy and COVID features were removed due to only 1 value across the dataset

# MODEL CREATION

AUROC and AUPRC



70/30 Test Split
AUPRC as metric because of dataset is imbalanced
Previous work on predicting diabetes ~ .90 AUROC

# HYPER PARAMETER TUNNING

- KNN - k =30
- DECISION TREE - max_depth 10, min_samples_leaf = 1, min_samples_split = 10
- LOGISTICAL REGRESSION - C=0.001, max_iter = 100, penalty=l2, solver= saga
- NAIVE BAYES - alpha = 0.01, nb_type = gaussian
- NEURAL NETWORK - activation = relu, alpha = 0.01, batch_size = auto, learning_rate=adaptive, verbose = True, validation_fraction=0.1
- RANDOM FOREST - max_depth = 10, min_samples_leaf = 1, min_samples_split = 5, n_estimators = 200

**BEFORE TUNING**

| MODEL | AUROC | AUPRC |
|---|---|---|
| KNN | 0.575 | 0.266 |
| DECISION TREE | 0.535 | 0.237 |
| LOGISTICAL REGRESSION | 0.666 | 0.362 |
| NAIVE BAYES | 0.643 | 0.328 |
| NEURAL NETWORK | 0.626 | 0.324 |
| RANDOM FOREST | 0.651 | 0.344 |

**AFTER TUNING**

| MODEL | AUROC | AUPRC |
|---|---|---|
| KNN | 0.629 | 0.321 |
| DECISION TREE | 0.624 | 0.316 |
| LOGISTICAL REGRESSION | 0.669 | 0.366 |
| NAIVE BAYES | 0.643 | 0.328 |
| NEURAL NETWORK | 0.632 | 0.327 |
| RANDOM FOREST | 0.670 | 0.365 |

# DISCUSSION
# &
# FUTURE DIRECTION

Difficulty increasing AUPRC is likely due to the nature of the problem. Long COVID is most likely not strongly associated with the included sociobehavioral, or demographic features. Other biological features such as viral strain and viral load are probably better indicators.

Future Directions:
- Add more features (COVID vaccination, former smoking status)
- Test new model (gradient boosting binary classifier)
- Shift question focus (COVID vs Long COVID)