

# IEOR 142: Final Report

## Group 10

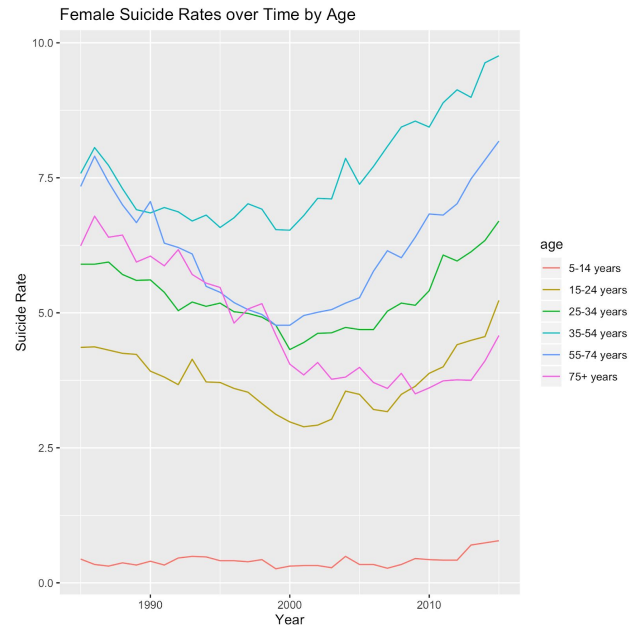
Dorothy Leung, Fanice Nyatigo, Haibin Lim, Kana Mishra, Ashleigh Purvis

### Motivation & Impact

The rate of people committing suicide every year in the United States is continually increasing under different circumstances and factors. The rate for women has been increasing for all age groups for the last 20 years, and the rate for men<sup>1</sup> for the last 15 years. Our final project focuses on predicting suicide rates for every 100,000 people in the United States by age group and gender over time. We are interested in learning the impact of various socio-economic factors in relation to suicide rates, in order to better ascertain which effects brought by the society and economy have the most impact on people prone to committing suicide. As a result, we hope to give more informed recommendations for public health interventions on significant factors based on our findings.

To accomplish this, the factors we decided to analyze were population size, HDI<sup>2</sup> (Human Development Index), GDP per capita, age group, gender, year, the depression percentage by gender, and the drug related death rate by age and gender. The GDP per capita and HDI are demonstrative of how much the government is investing in its people, with the population representing how large of a body is being governed and how the limited resources are being divided. The age group and gender are indicative of what individuals might be facing or experiencing in their lives (school for younger populations, work for middle aged, retirement for old age, and other life events). The year is used to find the trend over time. We were able to gather the depression rate over time by gender and average drug related death rates for every 100,000 people by age and gender over time from two additional data sources. We reasoned that the depression rate would give us a bit more insight into how individuals might be feeling on average, which we believe has a strong correlation with not just suicides, but also thoughts and attempts. The drug related death rate shed light on what individuals might be intaking to help cope with depression or other medical and non-medical issues. Overdosing can also be an indicator of a drug related death that was a suicide.

With our findings, we hope to be able to recommend changes be made to things like HDI and GDP per capita in order to help suicide rates decrease. In the best case, we would be able to see suicide rates decreasing and verify which factors are the most significant. This way, we could suggest that whatever changes caused the decreased be continued. The ultimate goal is to help save lives in the long run, whether it be for someone who has barely lived their life, lived most of it, or even anywhere in between. With more accessible open source data, our models could even be applied to other countries to help reduce suicide rates there. We could also further the impact by considering other features such as education level, food quality, and stress levels for all age groups by gender and country per year.



<sup>1</sup> The chart for males can be found in the appendix.

<sup>2</sup> For more information about HDI, visit <http://hdr.undp.org/en/content/human-development-index-hdi>.

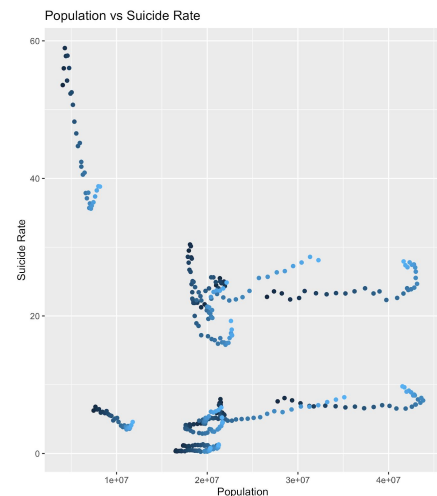
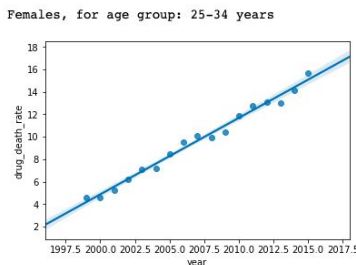
### Data<sup>3</sup>

The data used in this study was originally collected from Kaggle<sup>4</sup> with the original topic of “find signals correlated to increased suicide rates among different cohorts globally.” We extended the scope of our analysis by focusing on patterns particularly in the US of suicide rates by year, country, age group, etc. To fill in the NA values for the Kaggle dataset, we first made sure that the data was growing in a linear trend, and then we were able to interpolate by averaging the previous and next existing data points to fill the middle missing values.

To better analyze the overall impact of factors, we searched for additional data to add to our features. We found data by gender for depression<sup>5</sup>, which is the most common mental illness that people suffer before having a suicide attempt and is also often undiagnosed and untreated<sup>6</sup>. We were also able to find data for drug related death rates<sup>7</sup> by gender and age group per year. In order to merge the depression data which was not at the same level of granularity, we had to fill all age groups of the gender for the year with the same value, which could lead to limitations in finding patterns in the data. There were also some issues with NA values since the data we found did not fully span across the years we had from the Kaggle dataset. To fill in the NA values from the new datasets we performed EDA to study the general trend of the interval wrapping the missing values. We looked at each gender and age group pairing individually, and calculated best fit polynomials for each one.

Using these best fit predictions (courtesy of `numpy.polyfit`<sup>8</sup>), we were able to backcast<sup>9</sup> values to fill the early years we could not find data for. We chose not to overfit and most of these were filled in using linear or constant models, based on what our EDA showed us. However, some age groups seemed to follow no coherent trends, so we decided to simply average the next existing 5 years of values to find the value for that year.

**EDA:** Upon exploring the data<sup>10</sup> through our preliminary analysis, we found that the range of the suicide rate per 100,000 people ranged from 26% to about 59%. The sex most likely to commit suicide is male, and the generation with the highest suicide rate is the Silent generation. Interestingly, the age group that contains the highest rate was 75 and over (according to our data, this age group is mainly encompassed within the Silent and G.I. generations). We also discovered that population and suicide rate are slightly negatively correlated: suicide rates are highest in smaller populations. Suicide rates also seem to fluctuate the most around population of sizes 2 billion, which was an interesting find and may allude to a certain area with this size population having a particularly high suicide rate – we could further investigate this if we had additional data pertaining to individual regions within the U.S. We also found that there is a positive correlation between population and drug death rate of about 68%.



<sup>3</sup> Note: We changed our project completely a few weeks ago, so we had to start over with a new data set.

<sup>4</sup> Data obtained from <https://www.kaggle.com/russellyates88/suicide-rates-overview-1985-to-2016>.

<sup>5</sup> <https://www.statista.com/statistics/979898/percentage-of-people-with-depression-us-by-gender/>.

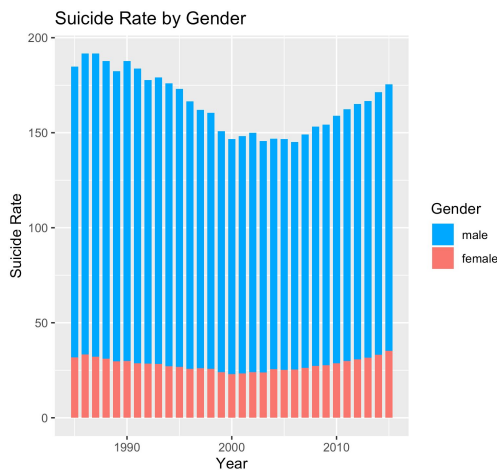
<sup>6</sup> <https://afsp.org/about-suicide/risk-factors-and-warning-signs/>.

<sup>7</sup> [https://www.cdc.gov/nchs/hsr/contents2018.htm#Table\\_008](https://www.cdc.gov/nchs/hsr/contents2018.htm#Table_008).

<sup>8</sup> <https://docs.scipy.org/doc/numpy/reference/generated/numpy.polyfit.html>.

<sup>9</sup> <https://stats.stackexchange.com/questions/59271/what-is-the-proper-name-for-a-backward-forecast>.

<sup>10</sup> Final version of our dataset: [https://github.com/kanam12/ieor142finalproject/blob/master/us\\_suicides\\_merged\\_no\\_na.csv](https://github.com/kanam12/ieor142finalproject/blob/master/us_suicides_merged_no_na.csv).



There is a correlation of about 29% between suicide rate and drug death rate.

Lastly, we determined that while the suicide rate declined from the late 1990s until about 2009, it has since increased and has steadily been increasing in recent years. In order to show the alarming increase of committed suicides the chart to the left shows the trend over time colored by gender. This further affirms the social magnitude of our project and the potential impact it can have.

Our final dataframe after merging and cleaning contains 14 columns and 372 rows. The original Kaggle dataset had 12 columns to which we added the 2 more mentioned earlier.

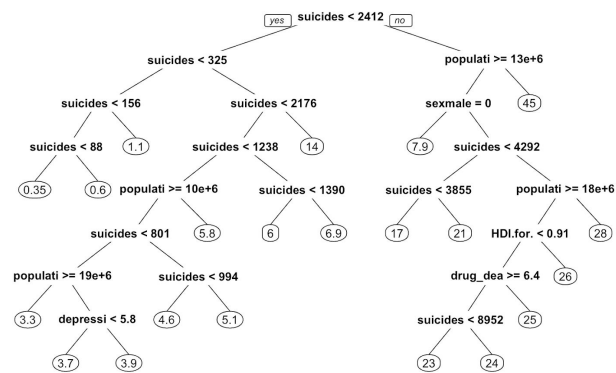
## Analytics Models

We focused on implementing our models with linear regression, CART, random forest, boosting, and time series. We also provided methods of cross validation to train models for an in-depth analysis. The metrics we are using are RMSE and out-of-sample  $R^2$  across the different models. Assuming our variables are independent and identically distributed, we performed a normalized random 70:30 split for training:test set ratio.

We initially used a baseline approach that always predicted the average suicide rate of 13.82 and obtained a fairly low  $OSR^2$  of 0.00003278162. We realized this naive baseline approach is not the right indication for us to reference our results for other upcoming models. Since suicides rate is a continuous variable, we opted to first start with a linear regression model to explore the relationship between suicides rate and some of the other variables. The first exploratory model included all the independent variables, after which we computed the Variance Inflation Factor (VIF) to determine multicollinearity<sup>11</sup>. From the linear model equation, it appears that the number of suicides, rate of death from drugs and generations X (part of intercept), Boomers, Silent and G.I all have a positive linear relationship with the rate of suicides. Generations Z, millenials and the population have a negative linear relationship. After 5 training rounds, the final linear regression model had an  $R^2$  value of 0.8014 and an  $OSR^2$  value of 0.7357. This shows that the model did not overfit too much because the model performed decently on the test set, relatively similar to the results of the training set. With CART, random forest, boosting, and time series, we are hoping to get even better results.

For CART, we built a regression tree using a 10-fold cross validation method with a tuning parameter,  $cp$  that ranged from 0 to 0.1, with a step size of 0.001.

We obtained the highest  $R^2$  with the lowest RMSE at a  $cp$  value of 0. The regression tree yields the following results, as summarized in the diagram. The reason we picked CART as one of our models is because we do not need to make any implicit assumptions about the underlying relationships of the features we are using, allowing us to capture nonlinear trends. From our EDA, we found that depression seemed to be nonlinearly correlated with suicide rates, seeming to increase only when the average depression rates were high or



<sup>11</sup> See appendix for details about VIF procedure and linear equation.

nonexistent and decreasing otherwise. Our goal was for CART to be able to capture and describe these nuances. As a result, our CART model had an  $R^2$  of 0.9398040 and an  $OSR^2$  of 0.8777108.

Along with CART, we created a random forest model in search for a possible model similar to CART but with a better  $OSR^2$ . We first created a tree with 5-fold cross validation with `mtry` values ranging from 1 to 5, and obtained 5 as the optimal fitting `mtry`. The final model for this `mtry` value gave us an  $OSR^2$  of 0.9645385. To explore further, we trained another set of random forest with the same 5-fold cross validation except with `mtry` values ranging from 1 to 10. From this training set, we obtained best tuned `mtry` of 10. Our final model on the test set with `mtry` = 10 had an  $OSR^2$  of 0.9928406 which was higher than that of CART.

The next model we ran was a boosting model under a gaussian distribution with `n.trees` = 1000, `shrinkage` = 0.001, and `interaction.depth` = 2. However, this model surprisingly underperformed and only had an  $OSR^2$  of 0.7628486. Some of the most influential features were `suicides_no`, `population`, and `sex`. We also tried running different versions of this gradient boosting model with cross validation, but for some reason, all of those attempts<sup>12</sup> resulted in an  $OSR^2$  of -0.2342707.

The last models we tried was the time series models: random walk, auto-regressive (AR), and random forest. The random walk model had an  $R^2$  of 0.9965203. However, this seemed to be unreliable for predicting drastic fluctuations. The one term AR model had an  $R^2$  of 0.9967 and an  $OSR^2$  of 0.9097873. The two term AR model had a slightly higher  $R^2$  of 0.9969. Both the two term AR model and RF model performed similarly<sup>13</sup>, with the two term AR model having a slightly higher  $OSR^2$  of 0.9090414 as opposed to the RF model's  $OSR^2$  value of 0.8945664.

We also tried to run a neural net model, but unfortunately ran into some issues that we (with the help of the professor) nor Google were able to resolve. We decided to not run unsupervised models like k-means since we were trying to predict continuous data instead of cluster things into groups or classify things. For this same reason, we also did not run confidence intervals for metrics such as accuracy, TPR, and FPR since those metrics did not apply to our project.

Models	R2	OSR2
Baseline	-	3.27816e-05
Linear Regression	0.8014	0.7357
CART	0.939804	0.877711
RF (mtry = 5)	0.973565	0.964538
RF (mtry = 10)	0.986766	0.992841
Boosting	-	0.762849
Time Series (Random Walk)	0.99652	-
Time Series (one term AR)	0.9967	0.909787
Time Series (two term AR)	0.9969	0.909041
Time Series (RF)	-	0.894566

Above is a chart detailing all of the different models we tried along with their performances<sup>14</sup>. The final model we would propose based on our  $OSR^2$  values would be the random forest with `mtry` = 10. We could try to cross validate the random forest model more and potentially use a larger test set to gain more confidence in its performance, or we could also use a one or two term AR model.

<sup>12</sup> The code is in the `cartandrf.Rmd` file but is commented out at the end.

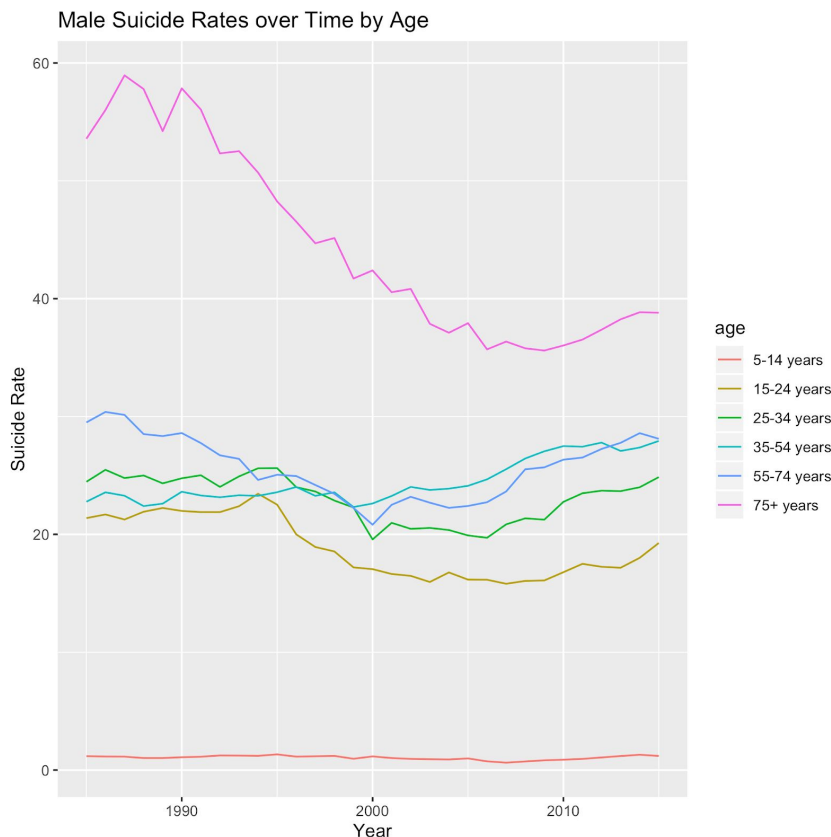
<sup>13</sup> See appendix for the plot with both models.

<sup>14</sup> The code for all our models can be found on our GitHub.

## Appendix

All of the project code can be found at <https://github.com/kanam12/ieor142finalproject>. Below are the elements from the footnotes, and following this is the pdf version of the Rmd file for our EDA.

### 1. Male Suicide Rate Trend over Time by Age (ref: page 1, Motivation & Impact)



3. See <https://github.com/kanam12/ieor142finalproject/tree/master/old> for the old dataset we were working with for previous rendition of the project. (ref: page 2, Data)

11. Not surprisingly, `gdp_per_capita ($)` had the highest VIF value ( $\text{gdp\_per\_capita } (\$) = \text{gdp\_for\_year } (\$) / \text{population}$ ), so it was removed and the model re-trained. We sequentially removed variables with the highest VIF and re-run the models until we were left with variables having little linear correlation with each other (VIF less than or equal to 5). The p-values of the variables removed (p-values greater than our cut-off of 0.05) were also not significant, and their removal only changed the model's  $R^2$  by 0.0024. Next, we re-trained the model by sequentially removing variables with the highest p-value above our 0.05 cut-off value, until we remained with those whose p-value was statistically significant. As with the VIF case, we constantly monitored the  $R^2$  value to ensure that it did not decrease significantly (the  $R^2$  only dropped by 0.0025). After 5 training rounds, the final linear regression model has an  $R^2$  value of 0.8014 and upon testing it on the test set, an  $\text{OSR}^2$  of 0.7357. The resulting equation is as follows:

```
suicides/100k pop = 3.815e+01 + 2.110e-03 suicides_no - 7.623e-07 population + 5.602e+00  
generationBoomers + 3.905e+00 generationSilent + 9.855e+00 generationG.I.  
Generation -1.682e+00 generationMillenials - 4.237e+00 generationGeneration Z -  
3.834e+00 depression_percentage + 2.253e-01 drug_death_rate
```

13. The red line represents the two term AR model and the green represents the RF model, with the black line representing the data.

