
ALCOHOL CONSUMPTION ON CORRELATING FEATURES

Jonathan Li

School of Engineering and Applied Sciences
University of Virginia
Charlottesville, VA 22904
jhl3xn@virginia.edu

Nicholas Talton

School of Engineering and Applied Sciences
University of Virginia
Charlottesville, VA 22904
nrt3xs@virginia.edu

Justin Zhang

School of Engineering and Applied Sciences
University of Virginia
Charlottesville, VA 22904
jmz8rm@virginia.edu

December 9, 2022

Contents

1	Abstract	3
2	Introduction	3
2.1	Overview	3
2.2	Hypothesis	3
2.3	Related Work	3
3	Method	4
3.1	Data	4
3.2	Models run	5
4	Experiments	5
5	Results	7
6	Discussion & Conclusion	8
6.1	Discussion	8
6.2	Conclusion	8
6.3	Limitations	8
7	Team Members and Contributions	8
7.1	Jonathan Li	8
7.2	Nicholas Talton	8
7.3	Justin Zhang	8
8	Citations and References	9

List of Figures

1	VABC Data	4
2	Census Data	5
3	Joined Table (VABC and Census)	5
4	Sklern models: Linear Regression and RandomForestRegressor	6
5	Architecture of the neural network we used	7
6	MSE of the model after 150 epochs	7
7	Performance of the neural network over the epochs	7

1 Abstract

Alcohol consumption has a significant impact on a person’s life and can have far-reaching consequences. Our goal was to create a machine learning model that relates socioeconomic factors in a particular region to the prevalence of alcohol consumption in that area, with the aim of helping to inform policy decisions on which aspects of a community should be targeted for the prevention of alcohol abuse and addiction. We used two datasets: Virginia socioeconomic data from the U.S. Census Bureau and VABC revenue by store in Virginia. After preprocessing, we tried various models we learned about in class for regression, such as linear, random forest, and neural networks. We found the most important factors to alcohol consumption to be being divorced or being widowed. Overall, we can treat the findings of our model as areas of interest, but our model is not a definitive authority on what factors affect alcohol consumption.

2 Introduction

2.1 Overview

Alcohol consumption has a significant impact on a person’s life and can have far-reaching consequences. Physical health, happiness, quality of relationships, and economic status are all variables that are directly or indirectly connected to alcohol use. On average, the higher the alcohol consumption in any given region, the more attention is needed to address the individual and collective effects of alcoholism. To this end, our goal was to create a model that relates socioeconomic factors in a particular region to the prevalence of alcohol consumption in that area, with the aim of helping to inform policy decisions on which aspects of a community should be targeted for the prevention of alcohol abuse and addiction. Such policies may include increased access to treatment facilities, increased public awareness campaigns, and increased access to social and medical services, among others. By focusing on the most influential factors, we can work towards a future of reduced alcohol consumption and improved quality of life for all.

We plan to quantitatively test correlations between socioeconomic factors and the alcohol consumption that influences the overall well-being of a community. From the patterns and correlations we find, we hope this data can be used to influence future policy around Virginia and make it more effective.

2.2 Hypothesis

By using machine learning we hope to gain a new interpretation of correlations with communities of people that directly relate to alcohol consumption. Using machine learning allows us to find new correlations that we perhaps did not realize were correlated with one another as if we were to use just manual analysis of this dataset. We do know that alcoholism is higher in groups such as the less educated, financially unstable, and socially/romantically unstable, so we expect the top factors to include at least some of these.

2.3 Related Work

[An overview of the genetic susceptibility to alcoholism³](#)

[Classification of EEG signals to detect alcoholism using machine learning techniques⁵](#)

[Using Machine Learning to Identify and Investigate Moderators of Alcohol Use Intervention Effects in Meta-Analyses⁴](#)

From our preliminary research, it seems that most research on alcoholism using machine learning has been done on an individual level, where certain features of an individual are used to predict a patient’s susceptibility to alcoholism. However, there has not been much research using machine learning to examine larger demographic trends in alcoholism to inform policy.

3 Method

3.1 Data

We used two datasets: Virginia socioeconomic data from the U.S. Census Bureau and VABC revenue by store in Virginia. We combined these datasets with relevant features that could be correlated with alcohol consumption and abuse. The links to our datasets are: [VABC sales](#)⁶, [Income](#)¹, [Age and sex](#)².

The demographic data we have is from the American Community Survey (ACS), a monthly survey (ran by the Census Bureau) of around 1% of homes in the US. The survey asks questions more detailed than the Census, so it would help us find more features that might be useful. As the Census website says itself, the ACS "provides local and national leaders with the information they need for understanding local issues and conditions," which reflects the purpose of this project.

The biggest issue we had is that our data does not have a 1-to-1 fit between the two data sets. Census data can be split into regions and counties, but the data from Virginia ABC stores are not split in the same manner. Because of this, we knew we needed to come up with a better way to pair each data with each other. Without the format of the rows of data matching as close as possible, the results may have not been accurate. We decided to split the data from the ACS and VABC stores into regions, as that both datasets could be sorted into regions. We cleaned the data and were able to get it into a Pandas format. We combined all of the datasets we have (ACS data and ABC data) into one CSV file with Pandas for convenience. The Virginia ABC revenue data had no missing values to impute, and because it's the thing we are trying to predict, we don't need to scale the data. From ACS data, we imported all four datasets we plan on using (tables DP02-DP05 on data.census.gov separated by counties in Virginia) into Pandas and combined all of the datasets into one table.

We found some missing values that needed to be imputed, so we imputed with the median where necessary. We did a one-hot encoding on all categorical data. Finally, we scaled all of our data (except our output labels) to have a mean of 0 and a variance of 1. Our labels, the values being predicted, were gallons of alcohol sold from ABC stores. We do realize that ABC stores are not the only sellers of alcohol but for this experiment, but we will assume the majority of alcohol sales are from ABC stores. Features we believed were extraneous and could either be captured by other variables or were irrelevant for our purpose of policy-making, such as use of computers, language spoken at home, place of birth, etc, were dropped.

A1	B1	C1	D1	E1	F1	G1	H1	I1	J1	K1	L1	M1	N1
1	District, Locality	Gallons sold	Gross Sales	Spirits & Wine % Net Sales	Cost of Goods \$	Gross Profit	Store Expenses	Net Store Profit	Allocation of Get Adjusted	Net Pric Rate of Return to Virginia			
2	Fairfax/Wellington	63027	4597360	760277	383680	2199926	1637055	843368	703065	219826	574057	29.82	
3	Southwest	94898	7285158	1209582	6075576	3504146	2571430	1486201	1085229	347971	737258	26.72	
4	Wytheville/Hillview	216340	16432729	2721500	13711229	7808875	5830354	2011743	3810811	791861	3026951	34.86	
5	Raeford/Blackob	276910	21155465	3486981	17856574	10146817	7511956	2891288	4820988	1010420	3010248	34.54	
6	Roseville	513687	43708487	7231275	36478212	21044780	15432452	4272871	11160580	2060770	9069810	37.29	
7	Staunton/Wayne	471487	38007870	6258350	31749520	18276357	13472163	4955204	8517959	1818485	6698474	34.09	
8	Winchester/Frost	287688	32699862	5388961	27276201	15730710	11548491	4015499	7531862	1561105	5975887	34.78	
9	Northern/Virginia	3450275	365133828	60201345	304632483	176912779	128819704	45347400	82572394	17491896	65060409	34.31	
10	Warrenton/Culpe	309161	26509247	4372369	22136878	12785047	9351831	3255736	6066095	1262812	4833283	34.73	
11	Charlottesville	470032	42584371	7023568	35500803	29576891	14983912	4829191	10953721	2017235	8025488	35.34	
12	Lynchburg	409334	32848132	5453557	27486965	15505377	11681188	4091088	7800088	1575006	6164293	36.89	
13	Danville/Martins	342424	28650951	4750016	23900935	13756867	10143968	3299040	6844828	1370304	5488624	35.67	
14	South/Easton	145618	11704252	1944032	97950226	5627931	4132088	1390740	2741949	556276	2162073	36.26	
15	Farmville	138310	10915613	1808448	9190565	5237112	3895452	1803912	1865541	520376	1443964	29.8	
16	Richmond	2135786	188732373	32851991	155880382	95866759	70013623	26211338	49802286	9528817	40273668	36.8	
17	Fredericksburg	564653	53652233	8887043	44785190	25888859	18916231	5748473	13167758	2567073	1050005	36.26	
18	Norfolk/Bedford	120675	9347380	1544348	7803442	4615411	3287631	1278031	2068889	448130	1553573	33.25	
19	West/Piedmont	202093	15544352	2585095	13059293	7568161	5491132	2101978	3389155	744777	2643377	33.43	
20	Petersburg/Hopk	336099	29756095	4941080	24815035	14332036	10482999	3231772	7251226	1429006	5821720	36.17	
21	Norfolk/Virginia	2591320	228229736	37770990	169456756	110540851	60417006	22176026	58237980	10963291	47254669	37.25	
22	Newport/News/H	1027325	96325590	15946445	60379145	45447740	33931405	10038887	23882917	4634320	19258198	36.55	
23	Eastern Shore	106825	8159161	1340418	6812743	3925587	2887155	1231700	1055455	391053	1263882	31.99	

Figure 1: VABC Data

The image shows a screenshot of a spreadsheet application displaying a large table of census data. The table has many columns, with headers in a smaller font at the top. The data is organized into rows, with some rows highlighted in light blue. The table appears to be a detailed demographic and economic dataset, likely from the US Census Bureau. The interface includes standard spreadsheet controls like zoom, pan, and search tools.

Figure 2: Census Data

The image shows a screenshot of a spreadsheet application displaying a joined table. The table has a wide range of columns, including demographic and economic indicators. The data is organized into rows, with some rows highlighted in light blue. The table appears to be a detailed demographic and economic dataset, likely from the US Census Bureau. The interface includes standard spreadsheet controls like zoom, pan, and search tools.

Figure 3: Joined Table (VABC and Census)

3.2 Models run

Since we are trying to predict alcohol sales/consumption with features, this is a regression problem and we needed to train regression models. We planned on using linear regression, random forest, and neural networks for this problem. Because of the limited number of regions in Virginia (<10), we need to prevent overfitting. To do this, we should reduce the complexity of the models while looking for the same performance, and conduct cross-fold validation or a train-test split, depending on how long our models take to train and run.

4 Experiments

In order to gauge the baseline performance of our model, the simpler models of linear regression and random forest regression were trained using 12-fold cross-validation on the Joined Table data. We chose these simpler models first as opposed to other more advanced models because they are relatively simpler models and do not require many hyper-parameters to tune. This gave us a general understanding of how these simple models fit the data.

```
[33] 1 from sklearn.linear_model import LinearRegression
      2 from sklearn.model_selection import cross_val_score
      3
      4 regr = LinearRegression()
      5 regr.fit(data_x_sklearn, data_y_sklearn)
      6 lin_scores = cross_val_score(regr, data_x_sklearn, data_y_sklearn,
      7                             scoring="neg_mean_squared_error", cv=12)
      8 lin_rmse_scores = np.sqrt(-lin_scores)
      9 display_scores(lin_rmse_scores)
```

Scores: [0.10355929 0.03123252 0.02089243 0.07005791 0.48675726 0.04317711
0.0518543 0.01625493 0.03236278 0.10101525 0.37247457 0.10207441]
Mean: 0.11930939612050494
Standard deviation: 0.1438708105097129

```
1 from sklearn.ensemble import RandomForestRegressor
2
3 regr = RandomForestRegressor()
4 regr.fit(data_x_sklearn, data_y_sklearn)
5 rf_scores = cross_val_score(regr, data_x_sklearn, data_y_sklearn,
6                             scoring="neg_mean_squared_error", cv=12)
7 rf_rmse_scores = np.sqrt(-rf_scores)
8 display_scores(rf_rmse_scores)
9
```

Scores: [0.10102256 0.01019727 0.03633315 0.04607457 1.19152649 0.05325424
0.13119278 0.12141895 0.15456234 0.03214606 0.55380402 0.31134998]
Mean: 0.22857353511934161
Standard deviation: 0.32530604238381244

Figure 4: Sklearn models: Linear Regression and RandomForestRegressor

After these preliminary experiments, moving forward, we used more advanced models that can capture more patterns from our features. We took our dataset, marked a section of our dataset as labels, and split the test data into training and testing.

We ran a more complex model, a regression artificial neural network (ANN). We ran this experiment in hopes to find a better correlation to work with for analysis. Our design for the architecture of the neural network consisted of taking in all features initially and then consolidating the features into 120 neurons, then 60, then into a single one-dimensional output. We arrived at this architecture by choosing an estimate for the optimal number of nodes to be $2/3$ the input size. We hoped that this result would be able to predict a better model than our preliminary experiments and better show any features that stand out from other ones.

```

1 import tensorflow as tf
2 from tensorflow import keras
3 from sklearn.metrics import confusion_matrix, accuracy_score
4
5 model = keras.Sequential([
6     keras.layers.Dense(units=448),
7     keras.layers.Dense(120, activation='relu'),
8     keras.layers.Dense(60, activation='relu'),
9     keras.layers.Dense(1, activation='linear'),
10 ])
11
12 model.compile(loss='mean_squared_error', metrics=['mean_squared_error'], optimizer=keras.optimizers.Adam(learning_rate=0.0001),)
13 history = model.fit(X_train, y_train, epochs=150, validation_data=(X_val, y_val))

```

Figure 5: Architecture of the neural network we used

5 Results

After we ran the model for 150 epochs we found the model was returning a very small mean squared error. The data we had was scaled so it makes sense for the mean squared error to be very small as a result (MSE = 0.0168). As it turns out, the ANN model seemed to perform rather well compared to the previous models we have tried. Because of this, we believe that ANNs would be good at predicting alcohol consumption from demographic data.

Because neural networks introduce nonlinearity into the data, feature importance is less defined. However, anyone can use a goal demographic and use that as an input into the model and see how that would affect the predicted alcohol consumption (ex. feed the ANN data with demographics with better education and see if predicted alcohol consumption goes up or down).

```

171 [-----] - 0s 40ms/step - loss: 0.0037e-08 - mean_squared_error: 0.0037e-08 - val_loss: 0.0168 - val_mean_squared_error: 0.0168
Epoch 146/150
1/1 [-----] - 0s 38ms/step - loss: 3.0206e-08 - mean_squared_error: 3.0206e-08 - val_loss: 0.0168 - val_mean_squared_error: 0.0168
Epoch 147/150
1/1 [-----] - 0s 46ms/step - loss: 1.9815e-08 - mean_squared_error: 1.9815e-08 - val_loss: 0.0168 - val_mean_squared_error: 0.0168
Epoch 148/150
1/1 [-----] - 0s 40ms/step - loss: 3.7671e-08 - mean_squared_error: 3.7671e-08 - val_loss: 0.0168 - val_mean_squared_error: 0.0168
Epoch 149/150
1/1 [-----] - 0s 43ms/step - loss: 2.6357e-08 - mean_squared_error: 2.6357e-08 - val_loss: 0.0168 - val_mean_squared_error: 0.0168
Epoch 150/150
1/1 [-----] - 0s 36ms/step - loss: 3.5426e-08 - mean_squared_error: 3.5426e-08 - val_loss: 0.0168 - val_mean_squared_error: 0.0168

```

Figure 6: MSE of the model after 150 epochs

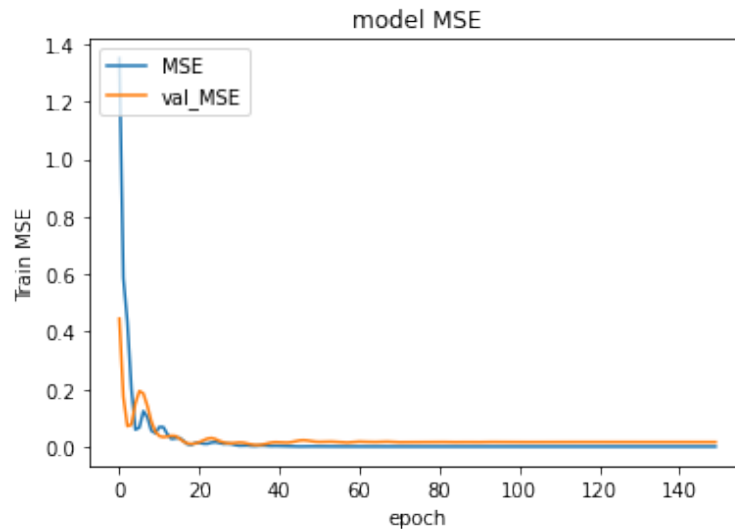


Figure 7: Performance of the neural network over the epochs

6 Discussion & Conclusion

6.1 Discussion

We found the most important factors to alcohol consumption to be being divorced or being widowed. This was expected and matched our hypothesis, as these have one thing in common: they are hardships. Hardships lead to a higher chance of alcoholism or just alcohol consumption in general.

This result shows that even with our assumptions, our model was still able to pick up on these trends, which means that our assumptions are not that problematic. However, part of the problem with our model was the low number of data points. The default way the VABC store data was organized was into huge regions, of which there were less than 10. Because of this, we had very limited data. However, overfitting didn't seem to be that much of a problem, because our validation error still went down as our training error went down, so we still feel that we can analyze our results confidently.

6.2 Conclusion

Overall, we can treat the findings of our model as areas of interest, but our model is not a definitive authority on what factors affect alcohol consumption. In the future, we could try to make our model more robust with more data points by reorganizing the VABC store data or generalizing our model with ABC stores around the country to inform our analysis on Virginia alcohol consumption.

6.3 Limitations

This research problem is a complex problem to solve, as predicting alcohol consumption involves a variety of factors such as cultural norms, economic factors, and demographic variables. This is possible to do, however, it would require a lot of data and perhaps multiple machine learning models to tackle. Additionally, to make the model be more precise in the future we could include more variables from other sources, such as local laws related to alcohol consumption and the availability of alcohol in stores and restaurants.

7 Team Members and Contributions

7.1 Jonathan Li

Data Preprocessing, program script writing, video recording and editing, and final report (abstract, preliminary experimentation, figures).

7.2 Nicholas Talton

Data preprocessing, video recording, final report (methods and next steps, figures, LaTeX outline).

7.3 Justin Zhang

Model training, acquisition of VABC and Census data, importing data to jupyter, final report (conclusions, abstract), Video recording.

8 Citations and References

- 1.) Bureau, U. S. C. (n.d.). S1901INCOME IN THE PAST 12 MONTHS (IN 2021 INFLATION-ADJUSTED DOLLARS). Explore census data. Retrieved December 7, 2022, from <https://data.census.gov/table?q=virginia>
- 2.) Bureau, U. S. C. (2020). S0101AGE AND SEX. Explore census data. Retrieved December 7, 2022, from <https://data.census.gov/table?q=virginia>
- 3.) Buscemi L, Turchi C. An overview of the genetic susceptibility to alcoholism. *Medicine, Science and the Law*. 2011;51(1_suppl):2-6. doi:10.1258/msl.2010.010054
- 4.) Parr, N. J., Loan, C. M., & Tanner-Smith, E. E. (n.d.). Using Machine Learning to Identify and Investigate Moderators of Alcohol Use Intervention Effects in Meta-Analyses. *Academic.oup.com*. Retrieved December 7, 2022, from <https://academic.oup.com/alcalc/article/57/1/26/6272648?login=true>
- 5.) Rodrigues, J. das C. et al. (2019, April 22). Classification of EEG signals to detect alcoholism using machine learning techniques. *Pattern Recognition Letters*. Retrieved December 7, 2022, from https://www.sciencedirect.com/science/article/pii/S0167865519301266?casa_token=RC8w7FsFSdgAAAAA%3Amd-JnjXkGogIJP9ud3fJxoovbfRAV2EEpNwDiMykr5bcjCilgyDKWlumsMELBcRvgKJOnFHHvIag
- 6.) Virginia Alcoholic Beverage Control Authority. (2021). Virginia ABC. Retrieved December 7, 2022, from <https://www.abc.virginia.gov/>