

BIG DATA Summer 2020 - Assessment 2 Report - Group 14

Considering social, economic and political variables, can we predict whether a country is happy or not?

1.0 Introduction

The following report outlines the development and implementation of a model that predicts - in a binary format classification - the happiness of a country. There are a variety of surveys, and metrics used to define happiness. In this case a 'life ladder' variable is taken from the World Happiness Data issued by the United Nation's Sustainable Development Solutions Network (sourced from the Gallup world poll). The 'Life Ladder' value is derived from the Cantrail Ladder - obtained by asking participants to rate their current life using a scale on which ten is their best possible life and zero their worst. An average is then taken to give a score for 'Life Ladder'. If a country's score is above six, the population is deemed to be happy. This model could be useful in a range of situations; from diagnosing mental health conditions in a medical context, to making governmental policy changes and predicting the way in which a population will vote as a result.

2.0 Data Preparation

2.1 Data Information

The original dataset has 19 features and 1562 data entries. After visually analysing the data, it was discovered that NaN (not a number) values were spread amongst the data, rendering the data on those rows impractical. Two features had over 20% of its data as NaN values: gini of

household income reported in Gallup, by wp5-year: 357 NaNs, and GINI index by World Bank estimate: 979 NaNs). They were both removed, and the rest of the sparse NaNs were eliminated by removing the rows on which they featured. After sanitisation of the data, our dataset had 17 features and 1120 data entries. A new feature was then created: 'LifeLadBin', with binary values, a score above six denoted happiness = 1, a score below six denoted unhappiness = 0 (Appendix 6.2).

Since people from 164 countries are surveyed over the years we plotted (Figure 1) the happiness trend from 2005 to 2017 for the countries we grouped into 10 regions.

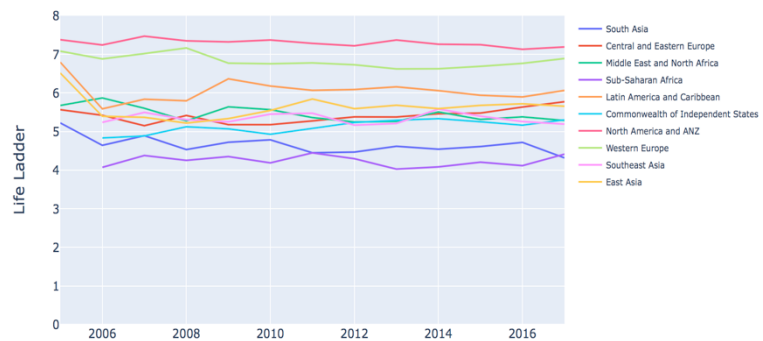


Figure 1 - Happiness trend fluctuation plot

The fluctuations were not very high and the happiness trends were mainly stable for each region. Therefore, we decided to place our focus on the countries without comparing it to the different years of observation.

2.2 Data Splitting

In order to build an accurate model that works on all data points without over-fitting, the original data set was split into three subsets; the training set contained 80% of the data, the validation and test sets contained 10% each. (Appendix 6.3).

2.3 Undersampling

A histogram was created (Figure 2) to display the distribution of the training set for descriptive analysis (Appendix 6.5 for the code).

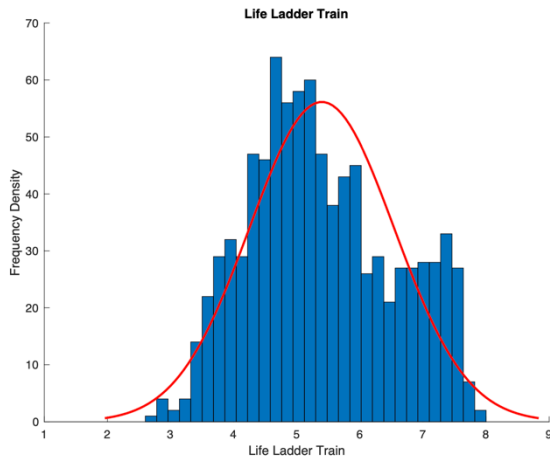


Figure 2– Descriptive analysis histogram

This showed that the training data for the ‘Life Ladder’ feature had a higher density at numbers lower than our threshold, and needed to be undersampled. ‘Unhappy people’ was the majority class, so this class was undersampled to balance the data. Before undersampling, there were 262 ‘happy’ data points and 634 ‘unhappy’ data points; after implementing the code (Appendix 6.4) there were 262 of both, meaning we had a balanced training set that would output a fair model. The same was done to the validation set.

3.0 Models

3.1 Logistic Regression

If a logistic regression analysis was carried out on the original unbalanced dataset, it could predict ‘not happy’ every time and achieve 70.8% accuracy. This is a benchmark to compare the final unbalanced test set to. However, as the

training and validation set had been balanced, a higher accuracy needed to be achieved.

A backward selection model was also tested (Appendix 6.20). Variables were removed until there was no improvement to the accuracy of the model. However, when comparing the accuracies from the backward and forward selection models, the forward selection model gave better results and was used to select the features used in logistic regression (Appendix 6.7). The least absolute shrinkage and selection operator was also used with a L1-norm and changing the C value. The selection process (with $C=1e9$) was chosen to avoid having unnecessary features. This led to the selection of the 6 features listed below:

```
(['Healthy life expectancy at birth',  
 'Positive affect',  
 'Social support',  
 'Generosity',  
 'Standard deviation/Mean of ladder by country-year',  
 'Standard deviation of ladder by country-year'],
```

The model built using these features produced 78.8% accuracy in the validation set (Appendix 6.8) and 83.0% accuracy in the test set (Appendix 6.9).

3.2 Decision Trees

Decision tree was used as a greedy model that made a locally optimal decision at each step, to get a global optima at the end. The impact of depth and the impurity on the accuracy of the model was plotted for both the training and validation sets (Appendix 6.10). As the model was made more complex, the accuracy of the training set kept on increasing while the validation accuracy was the highest between 3 and 5 features (Figure 3). Many combinations of depth were explored and a depth of 3 was chosen as it provided the best accuracy without overfitting the data.

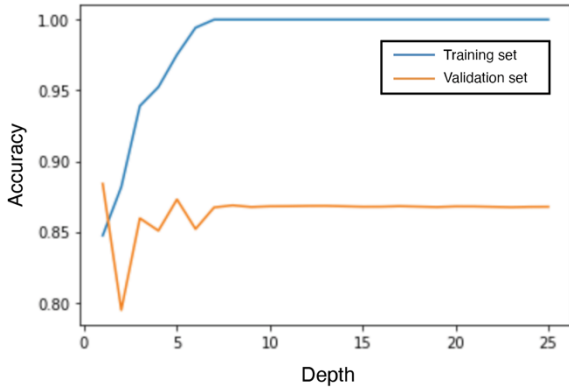


Figure 3 – Accuracy vs depth plot

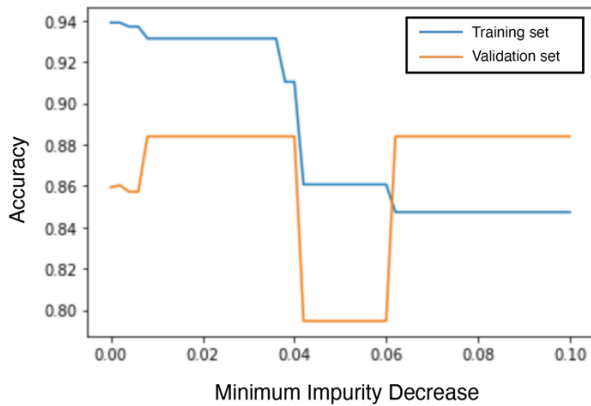


Figure 4 – Accuracy vs MID plot

The analysis started at the top of the tree and then going left or right according to the gini impurity that measures the probability of incorrect cases, this process gave us a “leaf node” with the lowest gini number (close to 0).

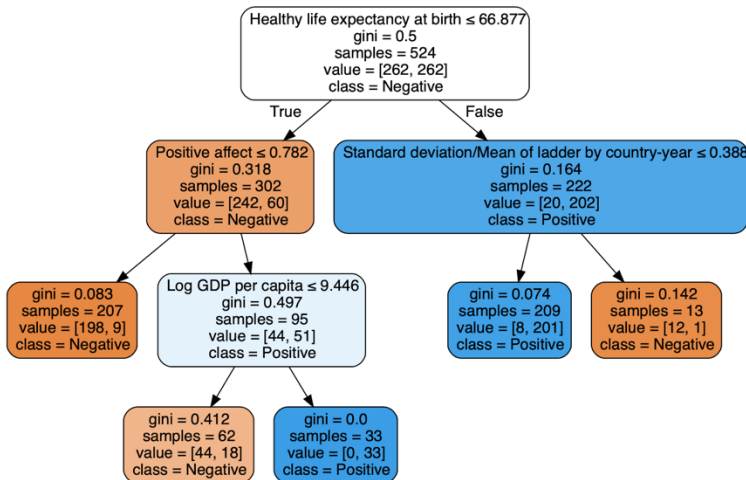


Figure 5 – Decision tree plot

3.3 Random Forests

To finally analyze the performance of our model, we ran the validation data sets with Random Forests. The Random Forests ran with a depth of 3 and with a minimum impurity decrease (MID) of 0.01 - many combinations of depth and MID were tested but the one which gave us the best accuracy in the end was kept. Our model achieved the following accuracies: 88.4% for the validation set and 87.5% for the test set (Appendix 6.16).

	Validation Set	Test Set
Precision	0.778	0.750
Accuracy	0.884	0.875
Recall	0.848	0.800

Table 1 – Results (more in Appendix 6.16)

3.4 Support Vector Machines

A support vector machine is a discriminative classifier, which allows us to validate the models we have used previously. Our validation data sets were run with SVM (using the hyper parameters C=1, gamma = 0.005 Appendix 6.17) and the following accuracies were achieved: 87.9% for the validation set (Appendix 6.18) and 78.6% for the test set (Appendix 6.19). These are similar values that our logistic regression model produced, validating our previous manual models.

4.0 Results and Analysis

Model	Validation set accuracy	Test set accuracy
Logistic Regression	78.8%	83.3%
Decision Tree	88.4%	58.9%
SVM	87.9%	78.6%
Random forest	88.4%	87.5%

Table 2 – Summary of findings

4.1 Logistic regression

Through forward selection of features, the logistic regression model showed that the following features were most influential in predicting happiness: 'Healthy life expectancy at birth', 'Positive affect', 'Social support', 'Generosity', 'Standard deviation/Mean of ladder by country-year' and 'Standard deviation of ladder country-year'. The difference in accuracy between the validation and test set was only 4.2% - this shows that the model has not overfitted. The accuracy of this model on the test set was 83.0% - this is significantly higher than the original 70.8 %, but not the highest accuracy of all our models

4.2 Decision Tree

Based on gini values, the decision tree selected the following features: 'Healthy life expectancy at birth', 'Positive affect', 'Standard deviation/Mean of ladder by country-year' and 'Log GDP per capita'. The first three features in this list are consistent with the logistic regression model. The training data was initially split between whether or not a country had a healthy life expectancy at birth. If that country did, the positive affect for that country was analysed, and a second split carried out. This process was continued for a depth of 3. The final ginis were in a range of 0 to 0.412 - there were two pure leaf nodes (ginis of 0.0 and 0.074) whereas the rest were impure. The accuracy on the validation set for the tree was high (88.4%), whereas the accuracy for the test set was relatively low (58.9%). This implies that the tree was potentially over-fitting, which is addressed in the implementation of Random Forests.

4.3 Random forests

The random forest model was created to assess the quality of the selection of features of the decision tree shown above. The random forest model runs a number of randomly created decision trees (the optimal number of tree number is 14 and was computed in the Appendix 6.14) and combines the results, to determine the most accurate prediction. Our random forest model gave an accuracy of 87.5% for the test set, which was higher than that of the decision tree. This is predictable due to the tendency of decision trees to over fit.

4.4 Analysis

Our models consistently show that having a healthy life expectancy at birth is the most significant feature (validated with a relatively low p-value: Appendix 6.25) when predicting whether or not a population will be happy. However, it is also interesting to consider the other features highlighted by forward selection, such as social support and generosity - the prominence of these features imply that physical and mental support in society play a larger role in the happiness of a population than other more concrete metrics, such as 'Log GDP per capita'. This information could prove useful for a government looking to make changes to improve the general happiness of its people.

The limitations of our model include the fact that it is only applicable to the countries and areas involved in the data collection - whilst the list is extensive, it is not quite worldwide (there are 195 countries in the world as opposed to the 164 in the dataset). The use of 19 specific features is also a limitation, as the results of our models can only predict happiness based on these metrics, so a wider set of features (like access to education, medical services...) in the

initial survey would have maybe improved our model.

Finally, using a category system other than the binary 'happy' or 'not happy' would have been more insightful, with our model predicting happiness on a scale, more conclusions could be drawn. For example, making the happiest country: Finland, a benchmark to compare it to other happy countries and seeing what social, economic and political features made the overall difference.

However, we were constrained by our research that provided us with the information that a Life Ladder score of 6 is the threshold above which people are deemed happy - for this reason we used a binary system, rather than a continuous.

5.0 Conclusion

After many trials, we have designed a model which gives improved predictions of whether a population is happy or not, and the factors that affect happiness most were uncovered. To get a tangible visualization of the life ladder distribution as well as its correlation with the socio-economic and geopolitical context, we have created a geographical visualization of happiness with a heat map with a gradient range between 3 and 8 (Figure 6).

Another heat map was plotted with the binary life ladder showing the countries that are happy in yellow and unhappy in blue - Appendix 6.21.

We have established a model that produces more false negatives and a lower recall. This result is better than getting more false positives because it is safer to predict that people are unhappy to push governments and the United Nations to act by reducing conflict and achieving peace for a better well-being of the population.

Another interesting application of this model could be uncovered as a result of recent events - in the wake of Covid-19, the dire economic situation that many countries face could lead to a severe decrease in happiness globally due to factors such as job loss and inability to travel. Our model could help governments to consider and include other factors to improve general happiness, and therefore maintain the population's wellbeing.

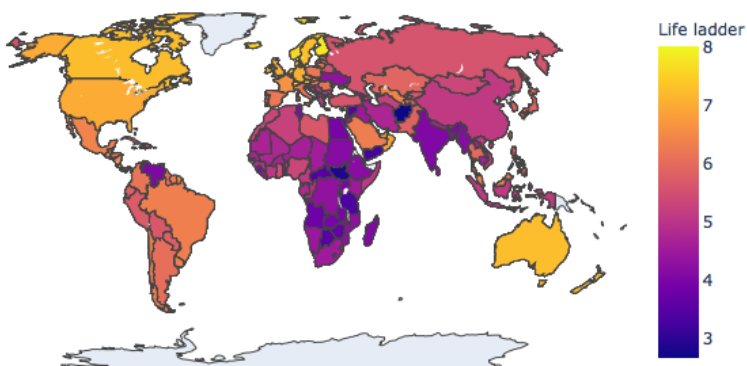
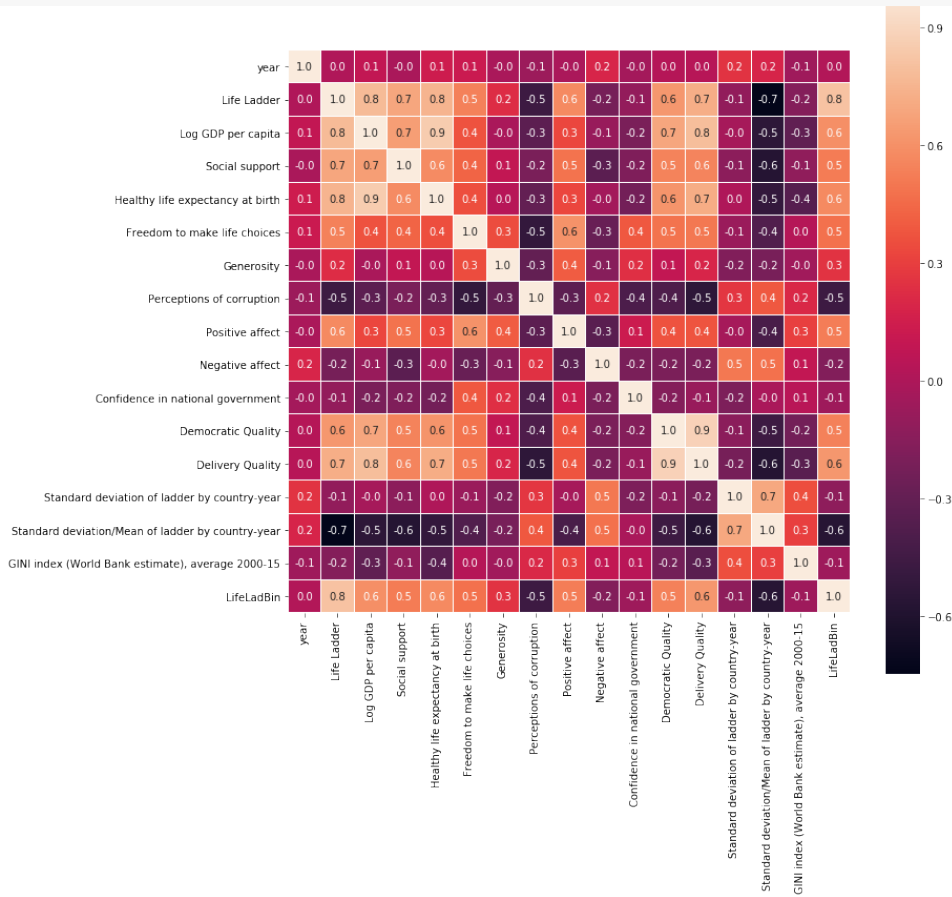


Figure 6 – Happiness heat map

! " ? \$ % & ' () * + , - . / : ; < = > ?

? - \$ B ^ 0 2 . 8 ^ h . / * 1 (\$ / 1 \$ 0 F 1 H \$ / F ' \$ 4 1 3 3 ' 8 . / * 1 (\$ C ' / H ' ' (\$ / F ' \$) * D D ' 3 ' (/ \$ D ' . / 2 3 ' 0 \$. () \$ 6 ' / \$. \$ D ^ 3 0 / \$ * (0 * 6 F / \$ 1 (/ 1 \$ H F * 4 F \$ 1 (' 0 \$. D D ' 4 / \$ / F ' \$ F . & & * (' 0 0 * () ' + "\$

```
import seaborn as sns
f,ax = plt.subplots(figsize = (12, 12))
sns.heatmap(happiness.corr(), annot = True, linewidths = 0.1, fmt = '.1f', ax = ax, square = True)
```



! " ? ^ \$ b 1) ' \$ / 1 \$ & 8 1 / \$ T * 6 2 3 ' \$, \$ K \$ 9 . & & * (' 0 0 \$ / 3 ' () \$ 8 2 4 / 2 . / * 1 (\$

```
layout = go.Layout(title="Happiness Trend over the years", font=dict(size=18),
xaxis=dict(title='Year', titlefont=dict(size=18),
tickfont=dict(size=14)),
yaxis=dict(range=[0, 8], title='Life Ladder',
titlefont=dict(size=18), tickfont=dict(size=14)),
legend=dict(font=dict(size=10)))
fig = {'data': [{'x': df_allyear[df_allyear['region'] == region].groupby('year')
.agg({'happiness': 'mean'}).reset_index()['year'],
'y': df_allyear[df_allyear['region'] == region].groupby('year')
.agg({'happiness': 'mean'}).reset_index()['happiness']},
'name': region, 'mode': 'lines', } for region in df_allyear['region'].unique()],
'layout': layout}
py.iplot(fig)
```

! " ? X \$ Y K B . 8 2 ' \$. () \$ ' 0 4 3 * & / * B ' \$ B . 8 2 ' 0 \$

	coefficient	std	p-value
intercept	-1.986	4.097	0.025
Healthy life expectancy at birth	0.050	0.040	10.105
Positive affect	0.744	2.885	0.208
Social support	-0.372	3.323	-0.029
Generosity	1.044	1.514	-4.972
Standard deviation/Mean of ladder by country-year	-1.949	4.994	-6.957
Standard deviation of ladder by country-year	-0.797	1.223	-11.846
			-3.220