



MONASH
University

MONASH
BUSINESS
SCHOOL

Department of
Econometrics &
Business Statistics

☎ (03) 9905 2478
✉ BusEco-Econometrics@monash.edu

ABN: 12 377 614 012

Report on Happiness up to 2022

Zhixiang Yang
EBS Honours Student

Yiqi Wang
Master of BA Student

Xintong You
Master of BA Student

Report for
Group 07 ETC5513

26 May 2022



<Trained Models >	<Model description >
Multivariate Linear Model	Simple Linear regression with multiple variables
Support Vector Machine Model	Use multiple learning algorithms (resampling and tree) to give us better results.
Decision Tree Model	Binary tree model have control statement.
Random Forest Model	Use multiple learning algorithms (resampling and tree) to give us better results.

Table 1: *Model Description of our Trained Models*

1 Modelling

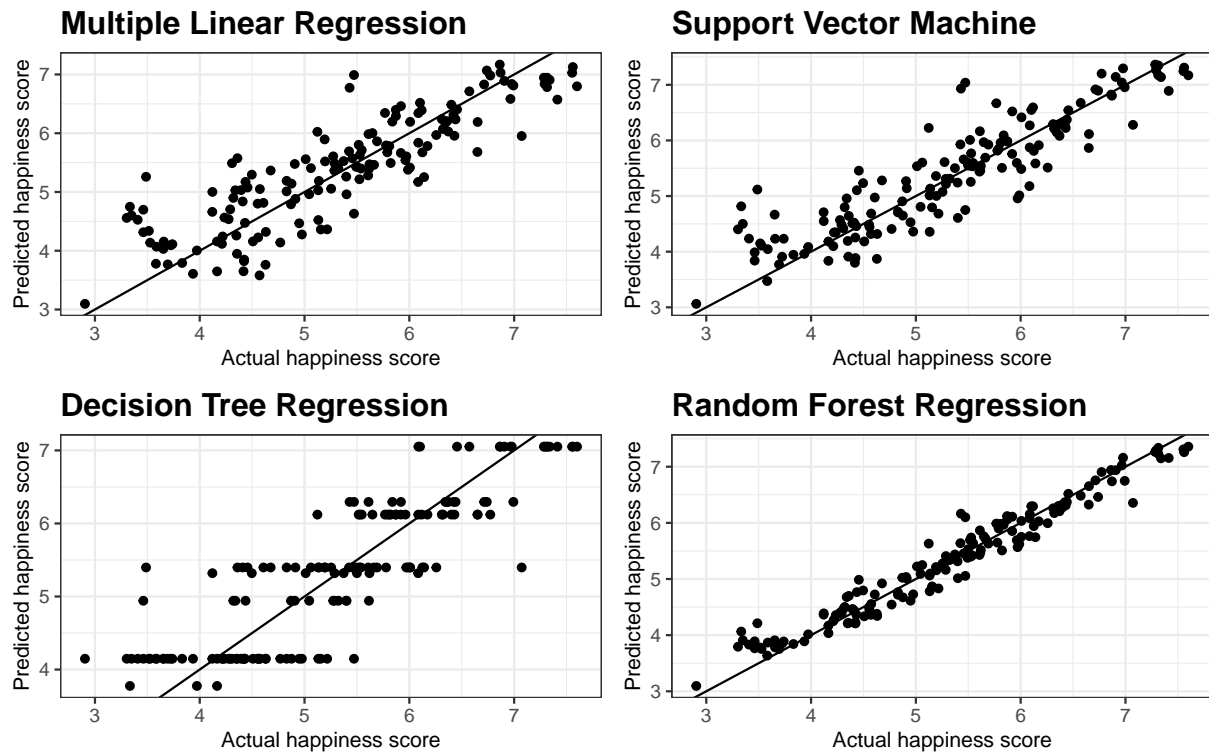
1.1 What is the most important variable to explain the happiness score differences across different countries in different years?

In this part, we only keep the common variables (Economy, Health, Generousity and Freedom) for these years to conduct our analysis. To explore the factors that could be contributing to the happiness score differences between each year, we first test few common models listed in Table1:

After trained our model based on the all historical data (from 2015 to 2022). We can see from the Figure ?? that the best model to fit the data is Random Forest model while the Multilinear and SVM performed similarly. The decision tree performed the worst because it changed the structure of the data. The result is similar when we have different training split ratios (see).

Table 2: Variable importance for Random Forest model

	IncNodePurity
Economy	348.66900
Health	327.08463
Freedom	169.32991
Generosity	97.79738



1.1.1 Random Forest regression

Random Forest Model have the variable selecting system (via bootstrapping) to decide the most significant tree and can reduce overwriting compared with decision tree. With that said, random forests are a strong modeling technique and much more robust comparing with many different methods (Liberian 2017). We can see from the plot that this model have captured the data well in the past few years for various training set (see).

To get which are the most important variables we then check their loadings in RF model (see Table 2).

In the Table 2, we conclude that the most important variables on explaining the happiness scores will be the Happiness and Health due to the loading for them are significant higher than others.

Table 3: Linear regression model for happiness scores without new data

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	2.437442	0.0752984	32.370444	0
Health	1.223204	0.1392052	8.787054	0
Economy	1.454316	0.0841520	17.282015	0
Freedom	1.372180	0.1095802	12.522156	0
Generosity	1.170207	0.1396671	8.378545	0

1.2 Multivariate Linear Model Analysis.

One advantage of multivariate linear regression is that it can allow us to analyse the relationship between different variables in a statistical coherent way. For example, marginal effects and percentage changes.

In our classic linear model, which is

Classical Linear Model :

$$\text{Happiness score} = \text{Economy} + \text{Health} + \text{Generosity} + \text{Freedom}$$

We can see from Table 3, in our classic linear model, they have similar loadings and all of them are significant, which we cannot tease out the important variables out of this model.

In order to better analyse the relationships, I add two new variables, which are CPI values and the population size for each country. However, due to the limitation of the new dataset, we can only conduct our analysis based on the 2020 data.

$$\begin{aligned} \log(\text{score}) = & -4.6000 - 0.0008 \text{ cpi} + 0.2591 \log(\text{economy}) \\ & -0.0026 \log(\text{population}) + 0.0114 \log(\text{health}) + 0.0032 \log(\text{year}) \end{aligned}$$

```
##
```

```
## Call:
```

```
## lm(formula = log_score ~ cpi + log_eco + log_population + log_health +
```

```
##      year, data = lognarm_hapall)
```

```
##
```

```
## Coefficients:
```

```
##      (Intercept)          cpi      log_eco log_population  log_health
```

```
##      -4.5998547    -0.0007945    0.2590779    -0.0025806    0.0113946
```

```
##              year
```

```
##              0.0031566
```

We can see that the total proportion of variance explained by the model with these variables are 60.49%. For the 4 predictors, the economy status(GDP per capita) contribute most to the happiness scores than other variables.

2 Endogeneity and Sample Selection Bias .

Country included in our dataset

Red: Countries are not Included; Blue: Countries are included

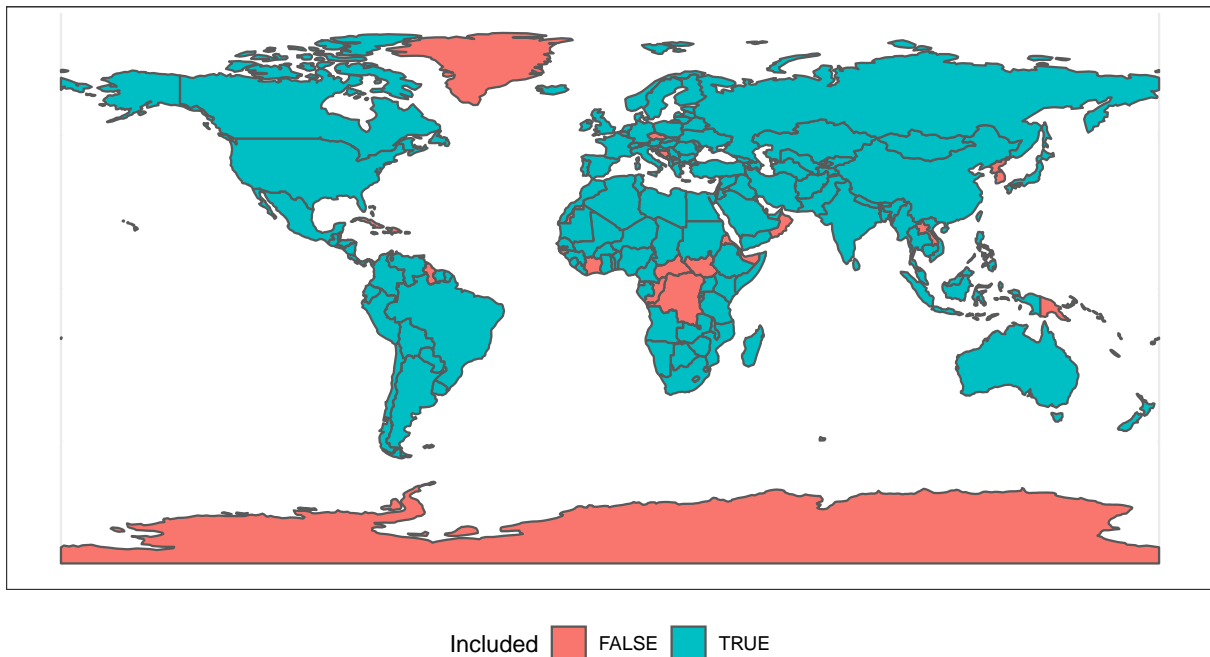


Figure 1: Colour the country that have included in Worldmap

Apparently, our model have some bias in selecting our sample. In our report, there are 34 are not included, which have been coloured in red from Figure 1

3 Potential Problems

The diagnostic of our regression model is shown in Figure 3. From the Residual and Fitted plot, we can see it has non-constant variance across the fitted value, which means the presence of Heteroskedasticity. Moreover, both the error distribution (see Figure 2) and the Q-Q plot Figure (see Figure 4) also suggest the density of our model is close to a normal distribution but seems there are some influences by outliers. In the

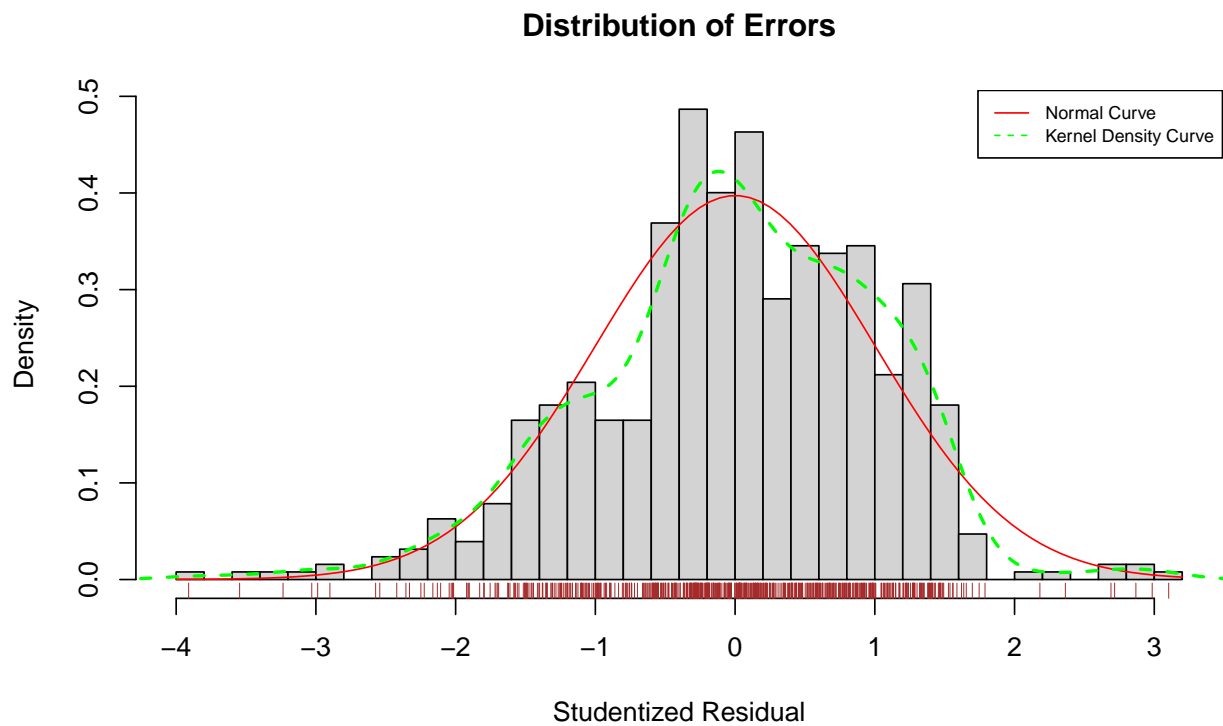


Figure 2: *Distribution of the residual term*

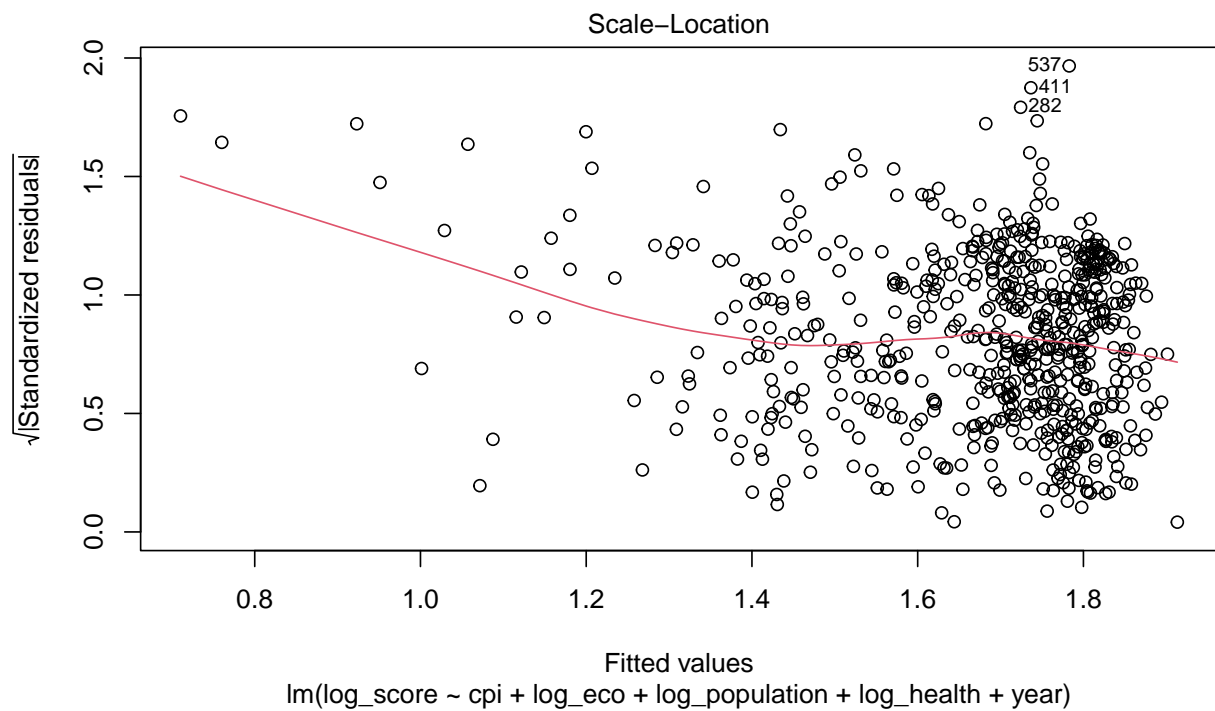


Figure 3: *Residual vs Fitted value plot (with standardised residuals)*

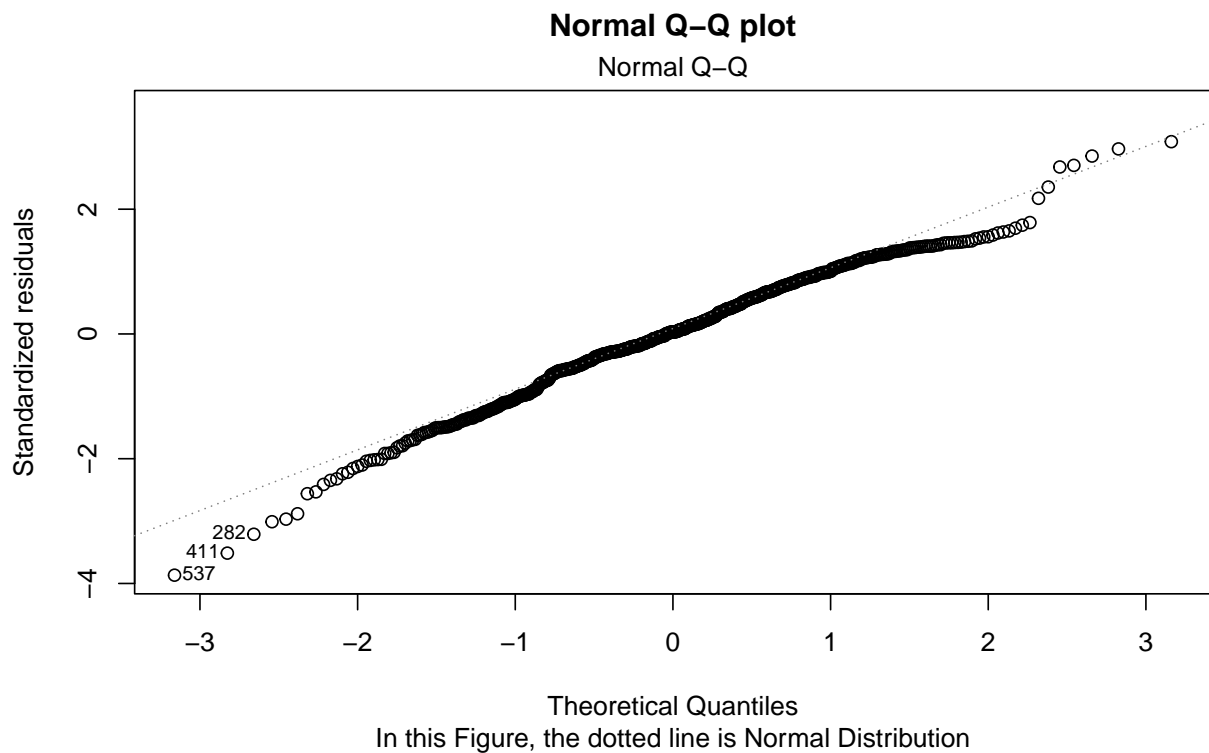


Figure 4: Q-Q plot fitted model residual density vs normal distribution residual density

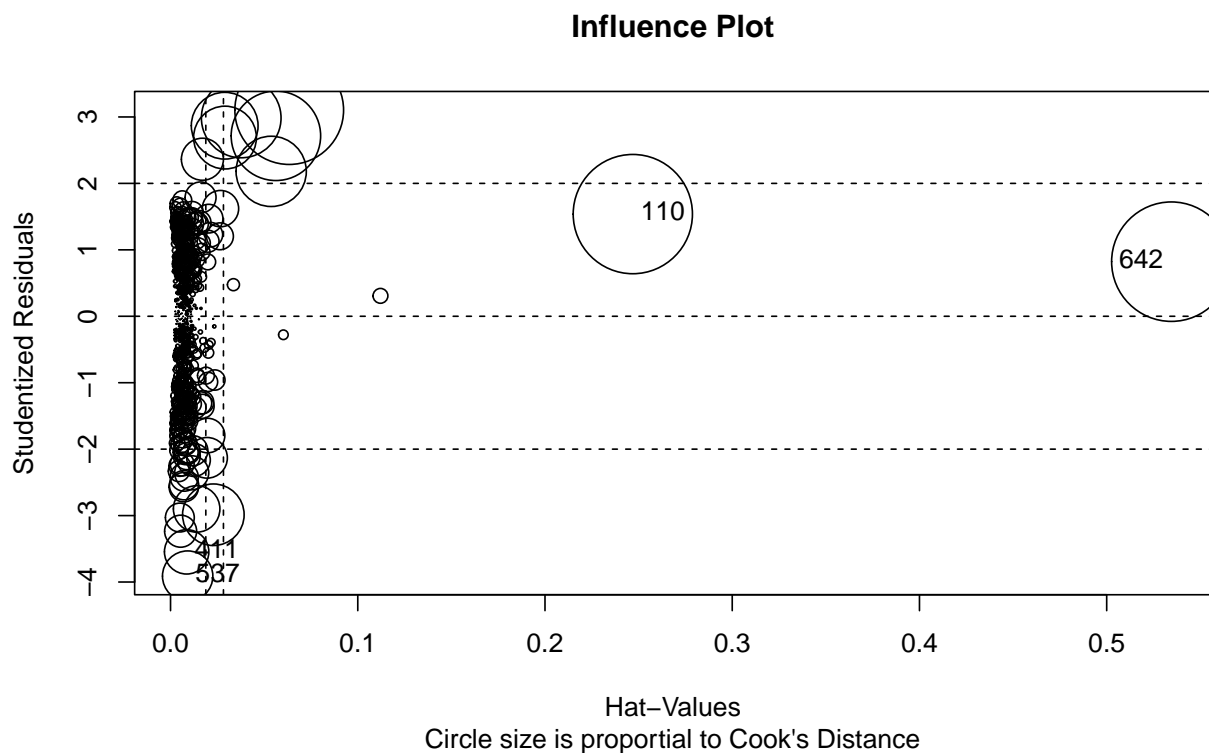


Figure 5: Influence Plot, the bigger circle means outliers

References

Liberman, N (2017). Decision trees and random forests. *Towards Data Science*.