# Brexit

## Prajwal Amin

## 2023-04-28

In 2016, The UK had a national referendum to decide wether the country should leave or remain in the EU ('Brexit'). This was declared through an election where votes came from different wards comprising of certain proportion of people across England, Scotland and Wales.
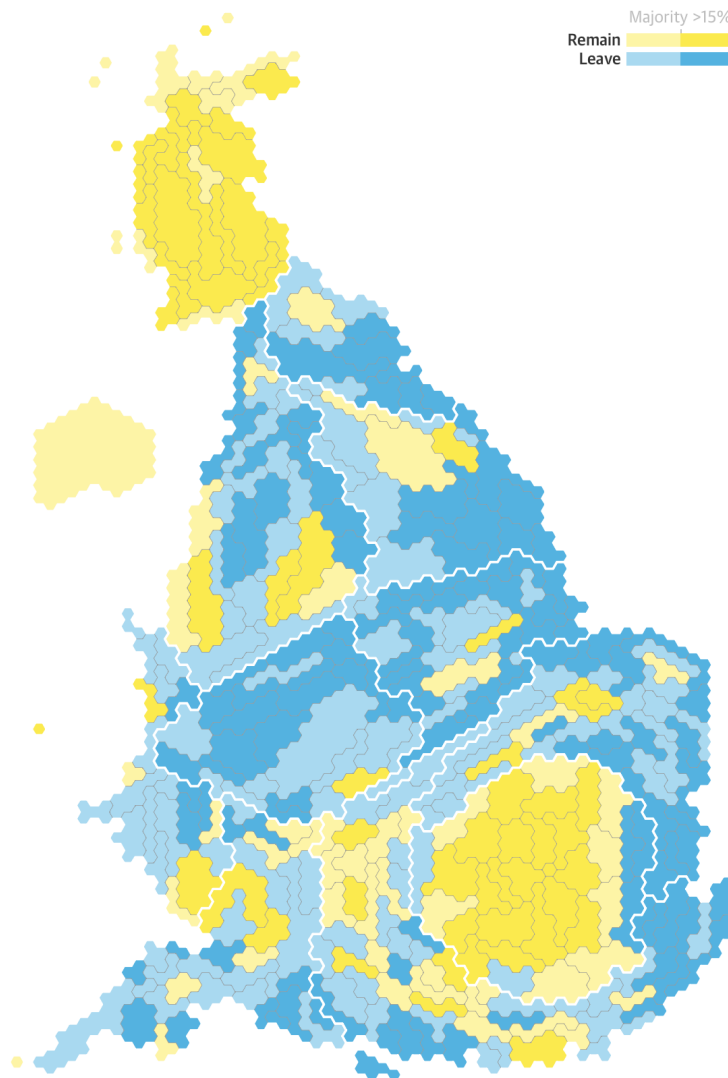


Figure 1: EU Referendum Results Map (The Guardian, 2016)

The Guardian newspaper presented some trends in modelling the outcome of voting. They have done the analysis but stopped at presenting the results graphically and commenting on the apparent patterns. Hence, we will use their analysis to compare our findings and perform some statistical analysis on the data.

The data contains 6 variables including the output variable 'voteBrexit which follows binary outcome. The variables are normalized so that the lowest value is 0 and highest is 1.

- abc1: proportion of individuals who are in the ABC1 social classes (middle to upper class)
- medianIncome: the median income of all residents
- medianAge: median age of residents
- withHigherEd: proportion of residents with any university-level education
- notBornUK: the proportion of residents who were born outside the UK

In the tasks below, we will fit a logistic regression model, explore the coefficients which have greater effect on the output, discuss the difficulties faced in interpreting these coefficients of the model and finally adapt an alternative approach to overcome these difficulties.

## Task 1

In this task we will fit a logistic regression model to the data to model the outcome of *voteBrexit* using all inputs. Through the summary of the model, will find the direction and magnitude of each of the inputs. Out of which, we will also identify the inputs having strong effects on the outcome. Finally, we will discuss about the findings of the model and compare them with the plots featured on the Guardian.

We will first split the data into training and testing using the ratio of 80:20. Then, will apply *train* data while training the model and *test* data for evaluation.

```
#Splitting the data (80:20 ratio).
set.seed(2)
split <- sample.split(brexit, SplitRatio = 0.8)

train <- subset(brexit, split == "TRUE")
test <- subset(brexit, split == "FALSE")
```

We will be using **glm** command in R to fit the data to the model. The "." (dot) indicates that all variables except *voteBrexit* to be used as inputs. Additionaly, *family = "binomial"* specifies the model to perform logistic regression.

```
# Training the model
model <- glm(voteBrexit ~ ., family = binomial, data = train)
summary(model)
```

```
##
## Call:
## glm(formula = voteBrexit ~ ., family = binomial, data = train)
##
## Deviance Residuals:
##      Min        1Q    Median        3Q       Max
## -2.57568  -0.07988   0.32642   0.62410   1.93422
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -0.2414     0.9876  -0.244  0.80689
```

```
## abc1            15.7788      3.6001    4.383 1.17e-05 ***
## notBornUK        6.1156      2.2322    2.740  0.00615 **
## medianIncome    -5.8385      2.2713   -2.571  0.01015 *
## medianAge        6.2603      1.6714    3.746  0.00018 ***
## withHigherEd   -25.0804      4.4940   -5.581 2.39e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 279.28  on 229  degrees of freedom
## Residual deviance: 165.45  on 224  degrees of freedom
## AIC: 177.45
##
## Number of Fisher Scoring iterations: 6
```

The estimates of the coefficients given by the model can be used to determine the magnitude and direction for each of the inputs. The magnitude indicates the strength of the impact a variable has on the outcome, whereas direction is whether an impact is positive or negative. Hence, given below are the inputs ranked on decreasing order of magnitude.

**1. withHigherEd:** This input variable has the highest magnitude compared to others. However, this has a negative impact on the outcome. As there is one unit increase in *withHigherEd*, the log-odds of voting for Brexit (TRUE) decreases by 25.08.

**2. abc1:** This variable has both higher magnitude as well as positive impact on the outcome. A single unit increase in this variable results in an increase of 15.78 in the log-odds of voting for Brexit.

**3. medianAge:** medianAge is a positive coefficient however the magnitude is not drastically high. An increase in one unit of this variable results in the increase of log-odds of voting for Brexit by 6.26.

**4. notBornUK:** This variable is also similar to *medianAge* considering its effect on the outcome. This variable has an increasing effect on the outcome by 6.12.

**5. medianIncome:** This variable has the least magnitude among all others, and the impact too is negative. A unit increase in medianIncome will decrease the log-odds of voting for Brexit by 5.84.

Apart from this, we can also observe that **p-value** for all of the inputs are below $< 0.05$ which suggests that the inputs are statistically significant. Although, the best predictor in modelling the outcome is **withHigherEd** (proportion of individuals with higher degree). The magnitude of this variable is exceptionally high and the p-value is quite low too.

## Evaluation

In order to evaluate the model we will first need to perform predictions using the trained model. Note that, while making predictions *test* data will be used.

```
# Making predictions
pred <- predict.glm(model,test, type = c("response"))
```

We will evaluate using **confusion matrix** metric from *caret* package which also provides the accuracy for the model.
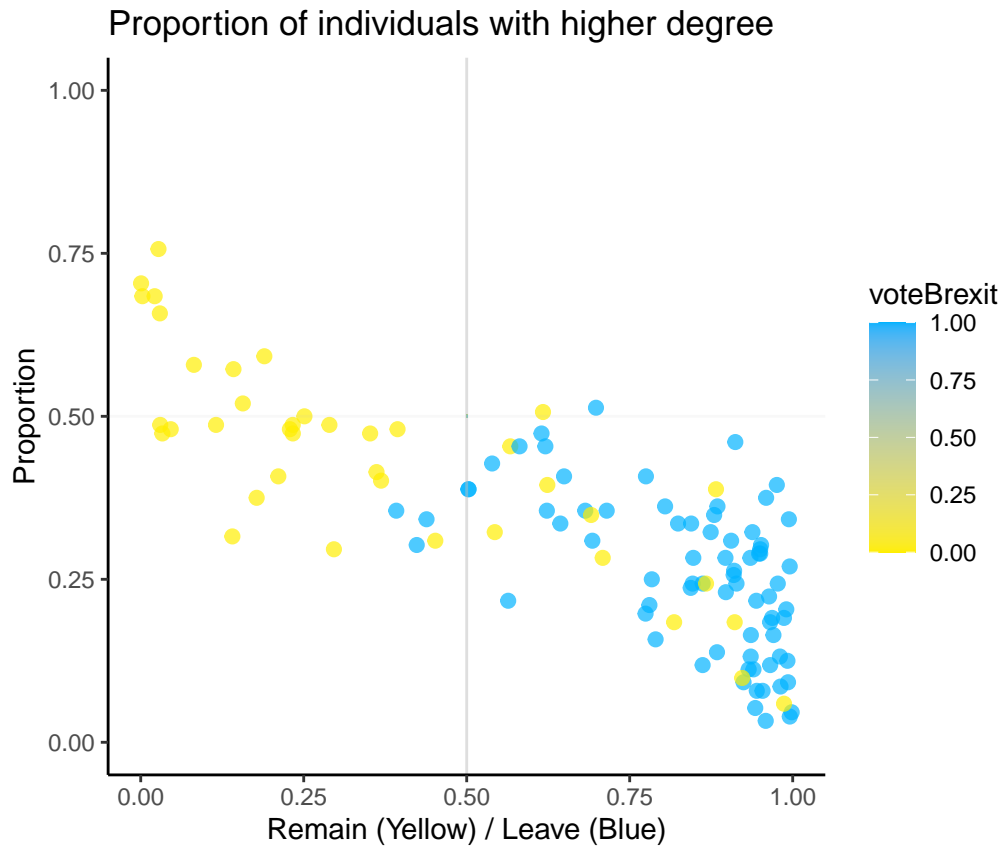
```
# Create a confusion matrix
confusionMatrix(table(test$voteBrexit, prediction_probs))
```

```
## Confusion Matrix and Statistics
##
##    prediction_probs
##      0  1
##   0 27 12
##   1  3 72
##
##                Accuracy : 0.8684
##                  95% CI : (0.7923, 0.9244)
##     No Information Rate : 0.7368
##     P-Value [Acc > NIR] : 0.0005181
##
##                   Kappa : 0.6906
##
##  Mcnemar's Test P-Value : 0.0388671
##
##             Sensitivity : 0.9000
##             Specificity : 0.8571
##          Pos Pred Value : 0.6923
##          Neg Pred Value : 0.9600
##              Prevalence : 0.2632
##          Detection Rate : 0.2368
##    Detection Prevalence : 0.3421
##       Balanced Accuracy : 0.8786
##
##        'Positive' Class : 0
##
```
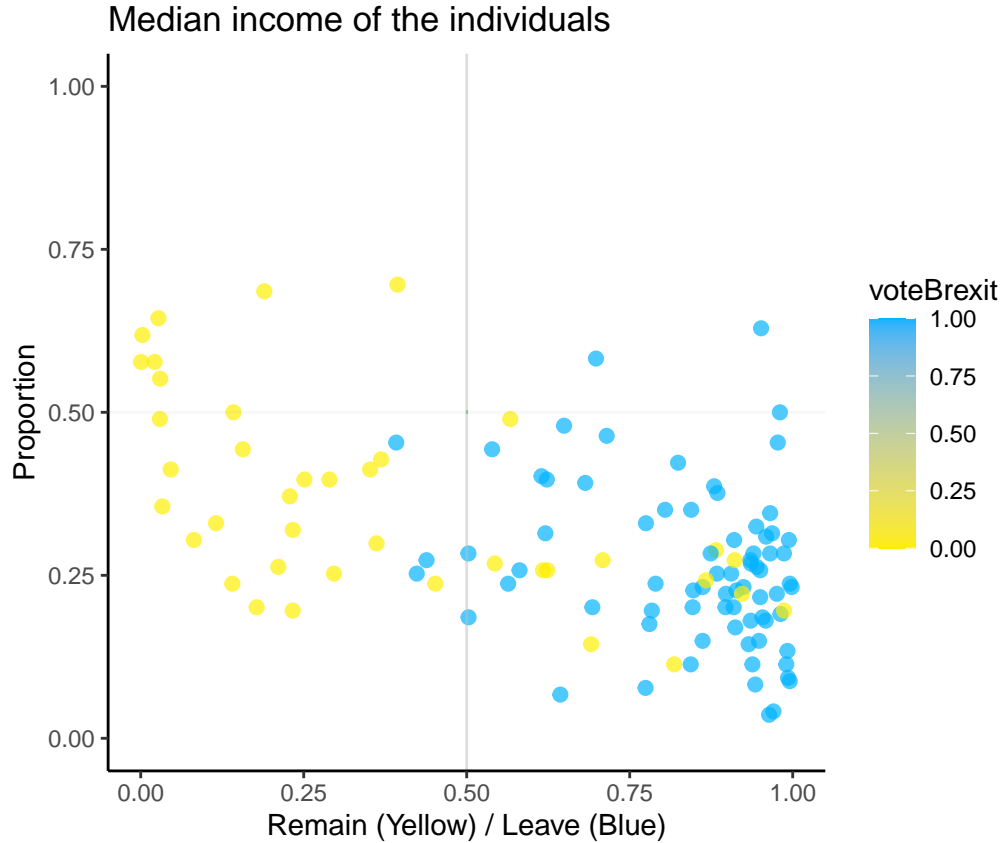
The **TP** and **TN** rate is high and the accuracy is $\approx 0.86$ which is quite good.

## Plotting the data

The predictions contain the probability of each electoral ward voting for Brexit. Hence, using the predictions, we can visualize and model the output for each of the inputs. First we will use *withHigherEd* variable against the model predictions to identify any patterns that is evident. We will also imagine a decision boundary so that if probability is $> 0.5$, the individual is likely to vote and if it is $< 0.5$ he will not.

**Proportion of individuals with higher degree**

The graph clearly shows that a ward with less than 50% of people with higher degree, is more likely to vote for leave. The correlation of the variables is strongly negative ($\approx -0.79$).

**Median income of the individuals**

Overall, *withHigherEd* seems like a good predictor for *voteBrexit*. In comparison to the findings featured on the Guardian website, it appears that their plots agree with the plots produced in this report and the difference is insignificant. The possible reasons for the patterns to vary might be due to the data that has been used to generate the plots or adapting different approaches in training the model. Further, the data points plotted on the website are classified much accurately compared to this report.

## Task 2

Since we have used a logistic regression model, there are certain factors that affect the interpretability of the coefficients fitted to the model, these are:

**1. Linearity** Similar to linear models, logistic regression assumes that there exists some linearity in the relationship between the variables. However, if non-linear then the model might not capture true effect of the variables, which might be misleading while interpreting the coefficients.

**2. Sample Size** It is always a good idea to perform predictions using large sample size, the results obtained using small sized data sets may sometimes produce inaccurate results. The standard errors of the coefficients might be high, thus leading the coefficients to be unstable.

**3. Normalization** Normalization rescales the input variables so that they have similar scales, which can help in interpreting the coefficients. But, sometimes interpretation might be difficult especially when variables with non-linear relationship have been normalized. Moreover, different variables will be measured in different scales which will get affected due to this process.

**4. High correlation** This results when the correlation between the input and output variables are high, which might make it hard to determine the individual effect of each variable and finding out which variable is relevant for the output.

Although, it is discussed in *Task 1* that the coefficients are interpreted in such a way that as one one unit increase in that variable increases the log-odds of the *output* by *coefficient value* when all other variables are held constant. This should not be considered as the only decision rule in determining the relationship between the input and output variables, as there might be other demographic factors that influence the decision to vote for or against Brexit.

Taking into account the factors mentioned above, it is not always reliable to determine the input variables based on their decreasing effect on the output. If done so there might be some uncertainty in the ordering. However, one ideal approach in determining the relevant features would be to consider the purpose of the task or the nature of the data along with the magnitudes of the model coefficients. Hence, taking such factors, the ordering of the relevant input features would be:

**1. withHigherEd** An individual with a higher degree will be well equipped with the knowledge to respond to any political affairs or national referendums.

**2. abc1** The individuals belonging to different socioeconomic groups possess different political attitudes or behaviour.

**3. medianIncome** This variable having the lowest magnitude can be considered as a good predictor since income level is also a driving factor that affects various voting behaviour.

**4. medianAge** Age cannot be said as a strong predictor for the output, hence the importance of this variable depends on the context of voting.

**5. notBornUK** An individual whether born inside or outside of UK may or may not affect the decision to vote.

As always these observations are subjective and the actual ordering varies in different scenarios. Additionally, factors such as sample size, linearity, normalization and high correlation could also affect in determining the order of the relevant input variables.

## Task 3

An alternative approach to carry out the analysis for task 1 would be to use other models such as random forests or decision trees. Since random forests provide better results in most situations compared to decision trees, we will use decision trees because of their structure which can be visualized. Decision trees create a tree-like structure where each *leaf* node represents the outcome for a combination of inputs.

We will split the data once again using the same split (80:20).

```
#Splitting the data into training & testing (80:20)
split_tree <- sample.split(brexit, SplitRatio = 0.8)

train_tree <- subset(brexit_tree, split == "TRUE")
test_tree <- subset(brexit_tree, split == "FALSE")
```
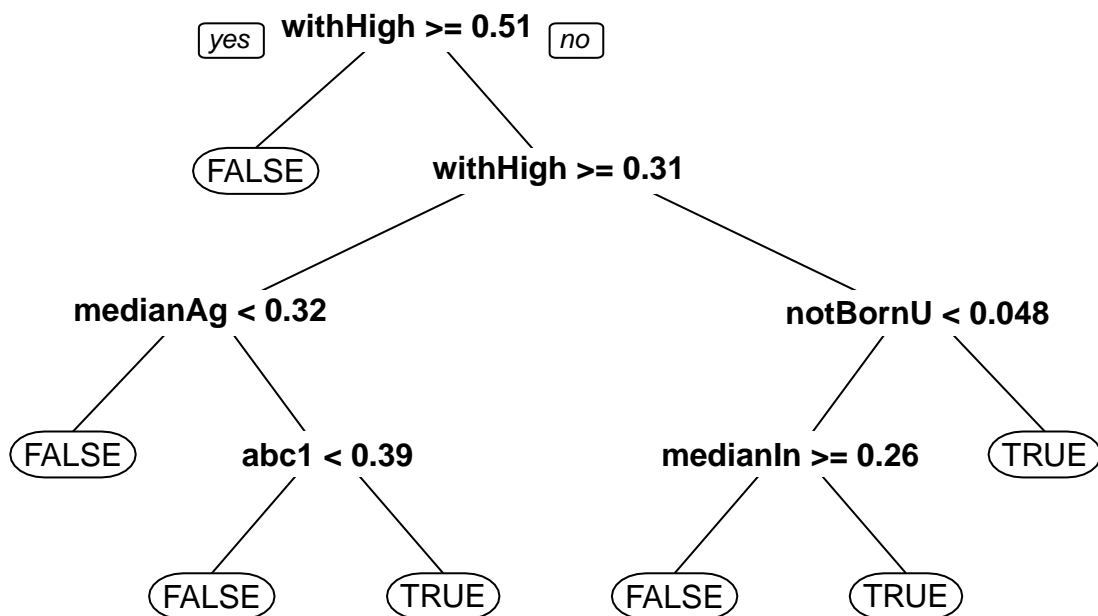
To fit a decision tree, we need to use **rpart** function from *rpart* package in R. Rest of the syntax is similar to that of **glm**.

```
# Training a decision tree
tree = rpart(voteBrexit ~ ., data = train_tree, method = 'class')
```

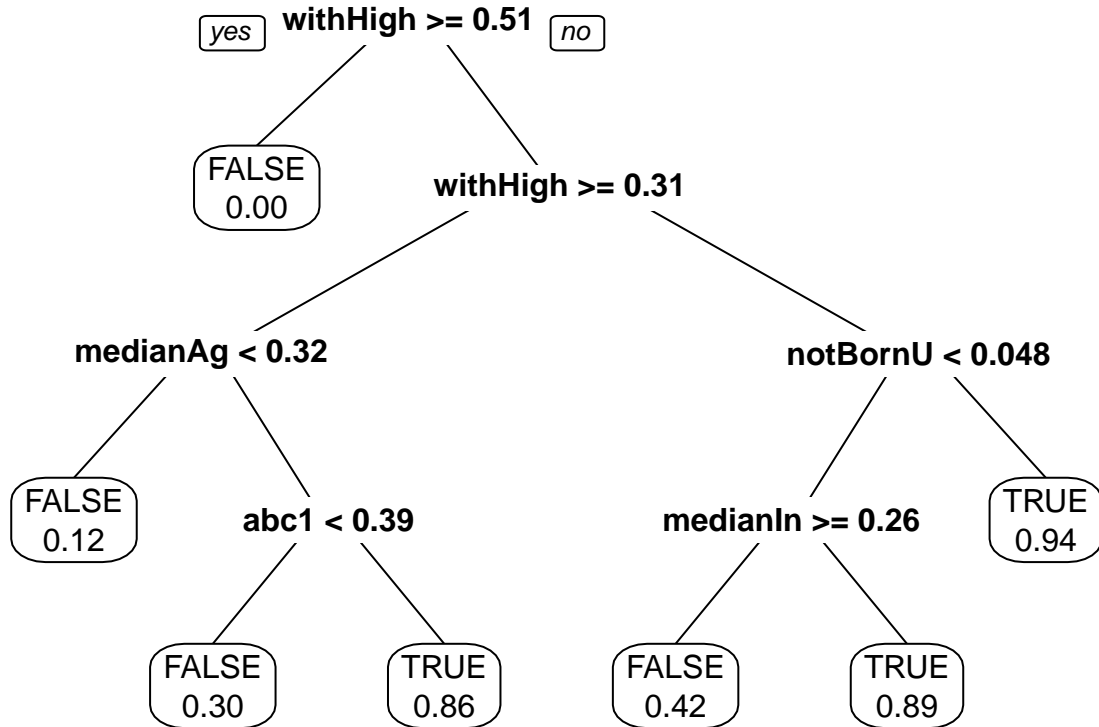Now, we can visualize the tree using the **prp** function

```
# Visualizing the tree,
prp(tree)
```

This gives a clear decision boundary at each split of the tree, thus improving the interpretation of the variables unlike the case in logistic regression.

We can also determine the probability of *voteBrexit = 'TRUE'* for every *leaf* node using an additional argument *extra=6.*

```
# Probability of the outcome
prp(tree, extra=6)
```

This approach of implementing decision trees comes with certain advantages as well as disadvantages which are given below:

**Advantages**

- The structure of the tree can be interpreted visually.
- Unlike logistic regression models, decision trees are capable in handling non-linear relationships between the input and output variables.
- Gives clear importance of the variables in modelling the outcome.

**Disadvantages**

- As the problem complexity increases, more decision points occur which makes it really hard to identify the structure of the tree.
- Decision trees are very unstable, a slight change in the data, the model produces different trees.
- This is the least accurate model when compared to other models such as random forests.