

# **Configuration Manual**

MSc Research Project MSc. in Data Analytics

Richard Burke Student ID: 15034097

School of Computing National College of Ireland

Supervisor: Dr. Catherine Mulwa

#### National College of Ireland



#### **MSc Project Submission Sheet**

	School of Computing	
rke		

	Richard Burke	
Student Name:		
	15034097	
Student ID:		
	MSc. In Data Analytics	2020
Programme:		Year:
	Research Project	
Module:		
	Catherine Mulwa	
Lecturer:		
Submission	17/08/2020	
Due Date:		
	The Significance of External Factors on an I	individual's Health
<b>Project Title:</b>	Determination	
-		
	1027 11	L
Word Count:	Page Count:	

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

<u>ALL</u> internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

	Richard Burke
Signature:	
	17/08/2020
Date:	

#### PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

Attach a completed copy of this sheet to each project (including multiple copies)	
Attach a Moodle submission receipt of the online project	
submission, to each project (including multiple copies).	
You must ensure that you retain a HARD COPY of the project,	
both for your own reference and in case a project is lost or mislaid. It is	
not sufficient to keep a copy on computer.	

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

# **Configuration Manual**

Richard Burke 15034097

# **1** Introduction

The objective of this configuration manual is to offer the reader an understanding of the software and hardware requirements required for the delivery of the research project: 'Extraneous Factors in an Individual's Health Determination'. Furthermore, the reader will be able to implement independently post completion of this document.

# 2 System Configuration

### 2.1 Hardware

Model: Dell Latitude 7480, 2018, 256 GB Processor: 2.8 GHz Intel Core i7-7600U RAM: 16 GB 2904 MHz Graphics: Intel HD Graphics 620 Dynamic ram allocation Operating System: Microsoft Windows 10

#### 2.2 Software

R Studio: Version 1.2.1335 R: R version 3.6.1 (2019-07-05)

# **3** Project Development

Following a Cross Industry Standard Process for Data Mining (CRISP- DM) Methodology, the data collection, data transformation, data processing, and data visualisation are outlined below.

## 3.1 Datasets

Deprivation index data received from Jonathan Pratschke. Ireland shapefile and CSO data downloaded from the Irish Governmental Open Data site.<sup>1</sup> Data sets:

- Playgrounds (2011)
- Streetlights (2010)
- Parks (2011)
- Schools (2011)
- Health Centres (2011)
- Leisure Sites (2009)
- Trees (2011)
- Census 2016 Small Area Population Statistics
- Ireland Census 2016 Boundary File

#### Jonathan Pratschke & Trutz Hasse (RIP):

• SAP Level Deprivation Index data – 2016<sup>2</sup>

## 3.2 Data Preparation

Data preparation is completed in R using R Studio. To generate a meaningful response variable, the self-assessed health determination from the census was ranked based on the proportion of total residents considering their health to be fair, bad, or very bad.

These were then split into five equal-sized groups based on their ratio and assigned categories.



Figure 1: Assigned Categories in R Studio

<sup>&</sup>lt;sup>1</sup> <u>https://data.gov.ie/organization/fingal-county-council</u>

 $<sup>^{2} \</sup>underline{http://trutzhaase.eu/deprivation-index/the-2016-pobal-hp-deprivation-index-for-small-areas/}$ 

Every facility or amenity needed to be mapped to a small area to generate aggregate statistics. These are completed in R Studio using the spatial packages (sp and rgdal). This was completed for each amenity separately before being amalgamated. Example of transformation for play areas included in Figure 2.



Figure 2: Overlay play areas and small areas in R Studio

The streetlight data set was mapped using easting and northing rather the longitude and latitude. This required a spatial transformation to generate their respective longitude and latitude references.



Figure 3: Reprojection of streetlights to latitude longitude

The points for all amenities were aggregated by each small area in Fingal.

10)	🔊 🔐 🖸 Source on Sale   🔍 🥕 🕴 📄	🔿 Run	24	le Source	•
1	library(dplyr)				
2 3 4 5	play_count <- play_geo %b% count(SMALL_AREA) colnames(play_count)[2] <- 'Play_Area'				
6 7 8	streetlight_count <- streetlight_geo %>% count(SMALL_AREA) colnames(streetlight_count)[2] <- 'Streetlight_Area'				
9 10 11	parks_count <- parks_geo M>K count(SWALL_REA) colnames(parks_count)[2] <- 'Parks_Area'				
12 13 14	school_count <- schools_geo %>% count(SWALL_AREA) colmames(school_count)[2] <- 'school_Area'				
15 16 17	health_count <- health_geo M>% count(SMALL_AREA) colnames(health_count)[2] <- 'Health_Area'				
18 19 20	leisure_count <- leisure_geo %>% count(SMALL_AREA) colnames(leisure_count)[2] <- 'Leisure_Area'				
21 22 23	tree_count <- tree_geo %>% count(SMALL_AREA) colnames(tree_count)[2] <- 'Tree_Area'				
24 25 26 27	Fingal_Areas <- as.data.frame(Fing)SWALL_AREA) colnames(Fingal_Areas)[1] <- 'SMALL_AREA'				

Figure 4: Aggregation to small area

Due to the sparse nature of the data sets (apart from trees and streetlights), secondary factors were created to represent the nearest distance from each polygon (small area) to the amenity.



Figure 5: Derive minimum distance between points and polygons

Finally, combine all the data derived data sets ready for processing and visualisation.

```
🕈 Run 🛛 🏞 📑 Source 🔻
             🔊 🗧 🗌 Source on Save 🛛 🔍 🎢 🗸 📗
    1
           #Merge data sets
           HP_Index <- read.csv("C:/Users/burkeri/Desktop/Project_2/Copy of HP Index 2006-2016 HP Index Scores
   2
   3
           HP_Index$Small_Area_New <- ifelse(HP_Index$SmallArea2017 == ''</pre>
   4
   5
                                                                                                   as.character(HP_Index$SmallArea2011),
   6
                                                                                                    as.character(HP_Index$SmallArea2017))
   8
          Fingal_Areas <- merge(Fingal_Areas, play_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, play_count, by = 'SMALL_AREA', all.x = IRUE)
Fingal_Areas <- merge(Fingal_Areas, streetlight_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, parks_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, school_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, health_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, health_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, leisure_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
Fingal_Areas <- merge(Fingal_Areas, tree_count, by = 'SMALL_AREA', all.x = TRUE)
15 Fingal_Areas <- merge(Fingal_Areas, dist_to_health, by = 'SMALL_AREA')
16 Fingal_Areas$row_id <- NULL
17
           Fingal_Areas <- merge(Fingal_Areas, dist_to_schools, by = 'SMALL_AREA')</pre>
18 Fingal_Areas$row_id <- NULL
           Fingal_Areas <- merge(Fingal_Areas, dist_to_playg, by = 'SMALL_AREA')</pre>
 19
20 Fingal_Areas$row_id <- NULL
21
           Fingal_Areas <- merge(Fingal_Areas, dist_to_leisure, by = 'SMALL_AREA')
22
           Fingal_Areas$row_id <- NULL</pre>
23
           Fingal_Areas <- merge(Fingal_Areas, dist_to_park, by = 'SMALL_AREA')</pre>
24
           Fingal_Areas$row_id <- NULL</pre>
           Fingal_Areas <- merge(Fingal_Areas, Fing[,c(1, 32)], by = 'SMALL_AREA', all.x = TRUE)
25
26 Fingal_Areas[is.na(Fingal_Areas)] <- 0
27 Fingal_Areas_2 <- merge(Fingal_Areas, HP_Factors[,c(1:3)], by = "SMALL_AREA", all.x = TRUE)
```

Figure 6: Merge the derived points and aggregations

## **4** Visualisation and Inspection

The aggregated data sets were inspected to produce visual and tabular depictions showing the dispersion of areas with below-average health determination across Fingal. The polygons must be fortified to produce the shapes and merged back to the derived data.

The plots are built using ggplot2 and colour coding polygons based on their health categorisation.

```
Source on Save  Sav
```

Figure 7: Depict visually using ggplot2



Figure 8: Dispersion of health categorisation across Fingal

The data was further aggregated by Assigned category and Electoral district to view of the underlying data points for each entity.

Figure 9: Aggregate to ED and Category

Assigned	EDNAME	n	Total	Percentage
Poor	Blanchardstown-Corduff	10	12	83.33
Poor	Balbriggan Urban	15	28	53.57
Poor	Blanchardstown-Roselawn	3	6	50.00
Poor	Blanchardstown-Tyrrelstown	3	6	50.00
Poor	Kilsallaghan	3	6	50.00

#### Table 1: Lowest 5 by categorisation

Visually comparing each ED to its overall mean in ggplot2 by transposing the wide data sets to long format and generating group means to add as points.

<sup>159
160</sup> summed\_tables <- aggregate(.~Assigned, Fingal\_Areas, sum)
161 summed\_ED\_tables <- aggregate(.~EDNAME, Fingal\_Areas, sum)
162 mean\_tables <- aggregate(.~Assigned, Fingal\_Areas, mean)
163 mean\_ED\_tables <- aggregate(.~EDNAME, Fingal\_Areas, mean)
164
165</pre>







Figure 11: Lowest 6 by self-determined health

# 5 Modelling

The package Caret in R was used to implement all the models (excluding Auto ML) and to develop optimised parameters for XGBoost and Random Forest. Classification accuracy, precision, and recall were used as measures of a model's performance on a random 25% holdout set.

### 5.1 Auto ML

Specifying the max number of model's limits run time and ensures model selection can be completed.



Figure 12: AutoML implementation

The output includes a set of models as a part of a leader board and allows the user to select which model to implement or allow default selection. Predictions of the testing data set can be made within the code, and a confusion matrix demonstrated using caret's in-built function.

AutoML using the package h2o requires pre-processing of the data to execute. This involves recoding the textual response variable to numeric. For this problem, all other variables were set as numeric with a 75%/25% split between training and testing data. The h2o cluster is initialised, and the data sets uploaded.

```
2
3
4 # AutoML Leaderboard
5 lb = aml@leaderboard
6
9 # prediction result on test data
9 # prediction = h2o.predict(aml@leader, test_h[,-5]) %>%
1 as.data.frame()
2
3 prediction$predict
4 #prediction$predict <- as.factor(round(prediction$predict,0))
5 # create a confusion matrix
6 confusion_matrix_h20 <- caret::confusionMatrix(testingXG2$Assigned, prediction$predict)
7
8 # close h2o connection
9 h2o.shutdown(prompt = F)
0</pre>
```



#### 5.2 Naïve Bayes

Naive Bayes was implemented using caret and its in-built functions. The default grid search and parameter optimisation were selected by caret. As with AutoML, the training set was split into a test and train set with a 25% and 75% split respectively.



Figure 14: Naïve Bayes implementation

#### 5.3 Random Forest

Additional parameters were specified for the Random Forest model to expand the tuning grid and parameter selection. This was implemented using caret on the same test and train data sets.

230	
231	############## Random Forest Tune
232	
233	start.time <- Sys.time()
234	
235	# Specifying cross validation - 10 folds.
236	<pre>trctrl &lt;- trainControl(method = "cv", number = 10)</pre>
237	
238	#SPecify randomly slected predictors
239	<pre>mtry &lt;- sqrt(ncol(trainingXG2))</pre>
240	
241	<pre>tunegrid &lt;- expand.grid(.mtry=mtry)</pre>
242	
243	# fit model
244	
245	rf_fit <- train(Assigned ~., data = trainingXG2, method = "rf",
246	metric='Accuracy',
247	trcontrol=control,
248	tuneLength = 13
249	)
250	# Variable importance
251	of Verture ( of fit)
252	ri_vapimP <- varimp(ri_iit)
203	
234	# prodict on pay data
255	# predict of new data
250	tort prodict $rf < prodict(rf fit torting/c2)$
258	cost_predict_fri < predict(fr_fri, costingxaz)
250	#conflucion matrix
260	
261	confusion matrix rf <= confusionWatrix(test predict rf as factor(testingxG2\$Assigned))
262	
263	end_time <- Svs.time()
264	
265	time.taken_rf <- end.time - start.time
266	

Figure 15: Random Forest implementation

## 5.4 XGBoost

XGBoost (gradient boosted decision trees implementation) was developed using carets manual tuning options to create an expanded grid search for the optimised parameters.



Figure 16: XGBoost implementation

#### 5.5 Multinomial Regression

Multinomial regression was implemented utilising the same parameter optimisation features. However, decay is the critical parameters, and an appropriate search grid was selected. The goal of the decay parameter is to discourage overfitting on data sets and increase reproducibility.

```
466 #### MLM
467
468 tuneGrid_mnl <- expand.grid(decay = seq(0, 1, by = 0.1))</pre>
469
470 mlm_fit <- train(Assigned ~., data = trainingXG2, method = "multinom",
471
                   trControl=trctrl
                      ,tuneGrid = tuneGrid_mnl
472
473
                      ,tuneLength = 13
474 )
475
476
477 mlm_var_imp <- varImp(mlm_fit)
478
479 mlm_test_predict <- predict(mlm_fit, testingXG2)
480
481 summary(mlm_test_predict)
482
483 confuson_matrix_mlm <- confusionMatrix(mlm_test_predict, as.factor(testingXG2$Assigned))</pre>
484 confuson_matrix_mlm
485
```

Figure 17: Multinomial Regression Implementation

# 6 Model Results

The variable importance (where available) is computed for each model in the processing steps. They are aggregated and visualised. Each model produces predictions on the test data set and the creation of a confusion matrix to allow comparisons of the accuracy, sensitivity, and specificity of each model.

🕈 Run 🛛 🕈 📑 Source

```
Source on Save Q V | G | Content of Con
```

Figure 18: Variable Importance

The confusion matrices were amalgamated to show the differential in performance across the accuracy, sensitivity, specificity, and kappa.

```
boo
661 results <- rbind(confusion_matrix_rf§overall,
662 confuson_matrix_mlmfoverall,
663 confuson_matrix_xgb2soverall,
664 confusion_matrix_hD2soverall,
665 confusion_matrix_hD2soverall)
666 results <- as.data.frame(results)
667
668 colnames(Names)[1] <- "Model"
669 Names <- as.data.frame(c(as.character('Random Forest'),
670 as.character('Kandom Forest'),
671 as.character('Naive Bayes'),
672 as.character('Naive Bayes'),
673 [s.character('H2o')))
674 Results <- cbind(Names, results)
676
```

Figure 19: Confusion Matrices