National
College *of*
Ireland

# Configuration Manual

MSc Research Project
Data Analytics

# Naumaan Mohammed Saeed Kazi

Student ID: x18130208

School of Computing
National College of Ireland

Supervisor:     Dr. Muhammad Iqbal

# National College of Ireland
## Project Submission Sheet
### School of Computing

| | |
|---|---|
| **Student Name:** | Naumaan Mohammed Saeed Kazi |
| **Student ID:** | x18130208 |
| **Programme:** | Data Analytics |
| **Year:** | 2019 |
| **Module:** | MSc Research Project |
| **Supervisor:** | Dr. Muhammad Iqbal |
| **Submission Due Date:** | 12/12/2019 |
| **Project Title:** | Configuration Manual |
| **Word Count:** | XXX |
| **Page Count:** | 15 |

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.

**ALL** 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.

| | |
|---|---|
| **Signature:** | |
| **Date:** | 12th December 2019 |

**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

Naumaan Mohammed Saeed Kazi

x18130208

MSc Research Project in Data Analytics

11th December 2019

# 1 Introduction

This configuration manual would present the software requirement, hardware requirement and system setup. Along with this the codes that have been used for programming that is written for the implementation of the research study:

"Using Machine Learning Models to Study Human Error Related Factors in Aviation Accidents and Incidents"

# 2 System Configuration

## 2.1 Hardware

Processor: Intel(R) Core i7-7200U CPU@2.6GHz GPU: NVIDIA GeForce GTX 1060Ti RAM:8GB Storage: 1 TB HDD; Operating system: Windows 10, 64-bit.

## 2.2 Software

\* Python using Jupyter notebook: Data analysis, data cleaning, pre-processing and manipulation. Feature selection and the implementation of machine learning algorithms and plots were done using libraries in Python.

\* Microsoft Excel:Used for saving of data, data exploration, and plots for explorations.

# 3 Project Development

Steps for project development are as follows: data exploration, data pre-processing (exporting data, handling mull values, removing unused attributes and again calculating features) and using hyper-parameter for model tuning. Number of codes have implemented during the many steps for analysing such as feature selection using Pearson Correlation, Carmer's V Rule and Random Forest for selection. K-fold validation using both approaches Stratified K-fold and 10-cross fold validation on all the models of machine learning used in this research. Generating confusion matrix, experimenting models and creating charts. In the later section codes for the study are shown with detailed information step by step.

## 3.1 Data Preparation & Pre-processing

The original data was downloaded from National Transportation Safety Board (NTSB) [1]. Which originally consist of a zip file, in which had the data file in MS Access database type. After which it was extracted to CS file file format to load in Python. After applying feature selection on the data, we had split the data into two categorizes test and train data for running the model. A systematic view of this data split is shown in the Figure 1



Figure 1: Data Collection and Splitting

The libraries used for data splitting is "sklearn.model_selection" from Python [2]. As 80% data is used for training purpose and the remaining 20% for testing. The data out is where the data that is no longer required for analysis i.e. discarded data. In data pre-processing contains label encoding as the values were "low" and "medium" which is shown in later. As the data required to be in 0's and 1' prior to model fitting. Data also needs to be normalized prior to fitting model, especially in Neural Network.

# 4  Code used for Machine Learning Models

The coding for this study has been performed using Jupyter Notebook. As the complete code was executed from correlation, cleaning and model in Python language only. The layout of this codes for explanation are designed as follows: feature selection, cross validation, class imbalance, label encoder and experiments on models.

---

[1] https://www.ntsb.gov/_layouts/ntsb.aviation/index.aspx
[2] https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html

## 4.1 Choosing Method for Feature Engineering

These three types are compared and contrasted against each other on depending which one best suit our classification techniques. As all these three methods are run on a single model which is XGBoost to see which gives the best accurate results. In Figure 2 we have used "SelectFromModel" function us used for random forest classifier to get the best feature from our data set[3]



Figure 2: Feature Selection Using Cramer's V Rule

Cramer's V Rule is most optimal for categorical data. In our dataset those variables which are categorical in nature are chosen and on them Carmer's V Rule is applied using the XGBoost model to check the accurate selection Figure 3. First, the cramer's function is defined using x and y. It used chi2 technique for feature selection[4].



Figure 3: Feature Selection Using Random Forest Selection

---

Correlation was accomplished using Pearson's correlation matrix. As the value above 0.8 is extreme correlation and above 0.4 is medium correlations Figure 20.



Figure 4: Feature Selection Using Correlation Matrices

## 4.2 Cross Validation

Cross validation (CV) is used in Zhang and Mahadevan (2019), the two approaches were opted first is the Stratified K-Fold and second is 10-fold CV. For all the models CV is used in the Figure 20 has shown for XGBoost. CV is helpful for problem like over-fitting & selection bias, as it gives meaningful insight about model. Using "Cross_val_score" function we implemented CV [5].



Figure 5: Stratified K-Fold for XGBoost

The CV which is very effective with model improvement is stratified k fold. The function "StratifiedKFold" function is used setting the n splits to 10. As the train and test would be split as per the CV Figure 20.

---

[5] https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.cross_val_score.html

Figure 6: 10-Fold Cross Validation for XGBoost

## 4.3 Experiment 1: ANN



Figure 7: Experiment on ANN

As observed from the Figure 7 before data is split and fitted into the model the data for ANN needs to be scaled. Using "StandarScaler" function this scaling is achieved in ANN [6]. We have used Keras library for this implementation as the optimizer chosen is Adam and the activation function would be 'relu'. ANN was applied in Burnett and Si (2017)

---

[6] https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html

Figure 8: Experiment on ANN

The epoch is run for 10, 50, 100, 150 Burnett and Si (2017). on which the optimal performance was achieved at 50 for non-hyper parameter ANN Figure 8.



Figure 9: Experiment on ANN

Classifier compile function is used for training the model of ANN as seen in Figure 9

## 4.4 Experiment 2: ANN Using Hyper-Parameter Tuning

Hyper-Parameter tuning is used to make the model better in terms of exploring more depths of the data and finding new hidden layers in the dataset. As observed from the Figure 10 we have defined only neurons for the Adam optimizer ANN function. It would run in the batches of 32 which the learning rate.
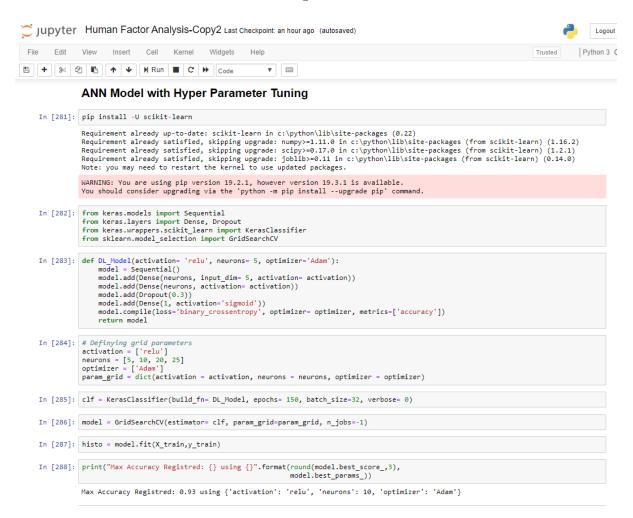


Figure 10: Experiment on ANN using Hyper-Parameter Tuning

## 4.5 Experiment 3: SVM

SVM was build using the "sklearn" library in python as this was the slowest performing algorithm Kumar et al. (2016). the function used here is the SVC() function were the kernel defined is linear Burnett and Si (2017). Based on this the prediction is done. Figure 11.

Figure 11: Experiment on SVM
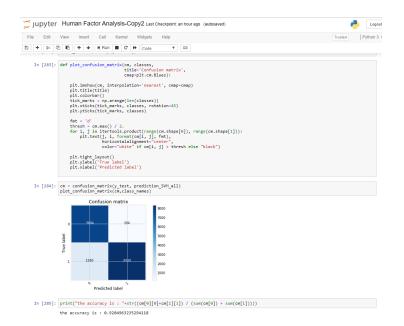
## 4.6 Experiment 4: SVM using Cross Validation



Figure 12: Experiment on SVM Confusion Matrix

Using the "Confusion_Matrix" function the plot is sets and the prediction has been achieved. The model very precisely makes a very good results on the confusion matrix. Figure 12 Using estimator as svc we are able to get set that CM for the SVM. There has been 10 folds for the SVM Burnett and Si (2017), Management et al. (2014). Figure 13
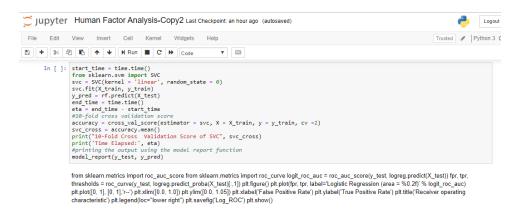
Figure 13: Experiment on SVM using Cross Validation

## 4.7 Experiment 5: Logistic Regression

The most simplest form of model it was implemented using the "sk.learn.linear_model" for the LR model Bazargan and Guzhva (2007), Bazargan and Guzhva (2011). Here the log.reg.predict is set for prediction. As observed from Figure 14. As the 10-Fold CV is implemented with cross val score function Mathur et al. (2017), Kharoufah et al. (2018).



Figure 14: Experiment on Logistic Regression and using 10 Fold Cross Validation

As observed from the Figure 15 we have implemented the Stratified fold with 10 CV using the function "StratifiedKFold" using the sciki learn [7]. It helps in exploring the data and finding better outputs Christopher and Appavu (2013). The SMOTE on data the model of Lr is run to handle the class imbalance using the library "imblearn.over_sampling". As the number of states for random synthetic sampling are restricted to 42Hofmann (2019).
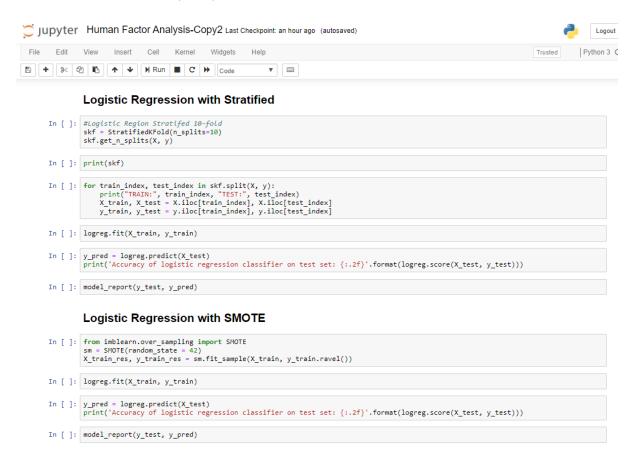


Figure 15: Experiment on Logistic Regression Stratified K-Fold and SMOTE

## 4.8   Experiment 6: XGBoost

It is one of the best performing models as the implementation was done using the "XG-BClassifier" function [8]. So were the SimpleImputer function being imported from sklearn.impute as it handles if any missing values, which in this study has been deal in the initial stages of this research. Prediction were rounded of using the round() of the values for the predictors. As observed in the Figure 16.

---

[7]https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.StratifiedKFold.html

[8]https://xgboost.readthedocs.io/en/latest/python/python_api.html

Figure 16: Experiment on XGBoost

The evaluation of CV in XGBoost is done using the function cross_val_score which is imported from scikit-learn. The folds are set to 10 and for experiment purpose it was set to 5 folds and 15 folds too. In which the best output was received with 10-folds. The function and model fit is seen in theFigure 17.
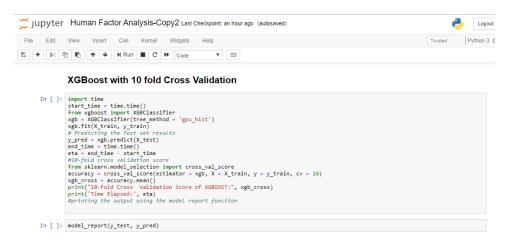


Figure 17: Experiment on XGBoost with 10 Fold Cross Validation

A similar approach has been taken for the stratified as the other models to explore the data, using stratified k fold function with the splits of 10 as observed in the Figure 18 XGBoost is one of the best performing models using CV we managed to get a very accurate result. The data is spilt and the trained and tested.

Figure 18: Experiment on XGBoost with Stratified K-Fold

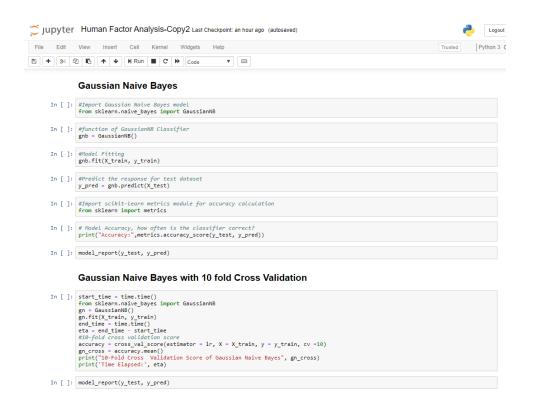## 4.9 Experiment 7: Gaussian Naïve Bayes



Figure 19: Experiment on Gaussian Naïve Bayes and with Using 10-Fold Cross Validation

in Figure 19 the model is a linear classifier which is best fit for supervised learning. As our data is big it is very much the best consideration. Using the function "GaussianNB()"

function for predictions. The library is from sklearn. And also the 10-Fold CV is implemented in this research using the cross val score() and also the implementation of stratified fold Figure 20.



Figure 20: Experiment on Gaussian Naïve Bayes and with Using Stratified K-Fold

## 4.10 Experiment 8: Random Forest

One of the model common and widely used models is the random forest as observed from the Figure 21 the model uses the ensemble approach which is used from the "sklearn.ensemble" to get the "RandomForestClassifier" function as the number of classifier decision trees and it would take an average in order to improve the accuracy and have control over the model being over-fit [9]. We have not defined any max depth as we wanted the model to learn on all the basis and not be restrictive in nature Burnett and Si (2017).



Figure 21: Experiment on Random Forest

---

[9]https://scikit-learn.org/stable/modules/ensemble.html

As seen in the Figure 22 stratified fold would create a test set in way that each of the class would have same class distribution and would be close connection. It is invariant to the class Li (2014). Using both the approaches of fold i.e. 10-Fold CV and stratified fold.



Figure 22: Experiment on Random Forest using Stratified and 10-Fold Cross Validation

# References

Bazargan, M. and Guzhva, V. S. (2007). Factors contributing to fatalities in General Aviation accidents, *World Review of Intermodal Transportation Research* **1**(2): 170.

Bazargan, M. and Guzhva, V. S. (2011). Impact of gender , age and experience of pilots on general aviation accidents, *Accident Analysis and Prevention* **43**(3): 962–970.
**URL:** *http://dx.doi.org/10.1016/j.aap.2010.11.023*

Burnett, R. A. and Si, D. (2017). Prediction of Injuries and Fatalities in Aviation Accidents through Machine Learning, pp. 60–68.

Christopher, A. B. and Appavu, S. (2013). Data mining approaches for aircraft accidents prediction: An empirical study on Turkey airline, *2013 IEEE International Conference on Emerging Trends in Computing, Communication and Nanotechnology, ICE-CCN 2013* (Iceccn): 739–745.

Hofmann, M. (2019). Analysis of aviation accidents data, (October 2018).

Kharoufah, H., Murray, J., Baxter, G. and Wild, G. (2018). Progress in Aerospace Sciences A review of human factors causations in commercial air transport accidents and incidents : From to 2000 – 2016, **99**(March): 1–13.

Kumar, P., Gupta, S., Agarwal, M. and Singh, U. (2016). Categorization and standardization of accidental risk-criticality levels of human error to develop risk and safety management policy, *Safety Science* **85**: 88–98.
**URL:** *http://dx.doi.org/10.1016/j.ssci.2016.01.007*

Li, G. (2014). Pilot-Related Factors in Aircraft Crashes : A Review of Epidemiologic Studies, (November 1994).

Management, E., College, S. E. and Force, A. (2014). Study on the Aviation Accidents Due to Human Factors Based on Improved Support Vector Machine 2 The selection of key indicator based on Analytic Hierarchy Process ( AHP ), pp. 278–283.

Mathur, P., Khatri, S. K. and Sharma, M. (2017). Prediction of Aviation Accidents using Logistic Regression Model, pp. 1–4.

Zhang, X. and Mahadevan, S. (2019). Ensemble machine learning models for aviation incident risk prediction, *Decision Support Systems* **116**(October 2018): 48–63.
**URL:** *https://doi.org/10.1016/j.dss.2018.10.009*