



**AGH UNIVERSITY OF SCIENCE AND TECHNOLOGY  
FACULTY OF APPLIED MATHEMATICS**

# The Selection of a Model for Airlines Customer Satisfaction

Joanna Krężel, Anna Matysek, Piotr Mikler, Adam Szczerba

02 January, 2022

## **Abstract**

The project aims to analyze the data about satisfaction of Investico Airlines passengers. Having customer-granular observations about the cruise and a reported satisfaction level for particular aspects of the flight we try to fit models which predict whether they are satisfied with the service or not. The binary classification models used during the project are  $\{ \} \{ \} \{ \}$ , out of which our recommendation is  $\{ \}$  based on  $\{ \}$ . This document describes the process we undertook and presents the results of data preprocessing, model selection and model validation.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Data</b>	<b>3</b>
2.1	Data Preprocessing . . . . .	4
2.1.1	Feature Encoding . . . . .	5
2.1.2	Missing data treatment . . . . .	7
2.2	Exploratory Data Analysis . . . . .	10
2.2.1	Feature Selection . . . . .	11
2.3	Categorical Variables . . . . .	12
2.3.1	The information Value . . . . .	12
2.4	The Continous Variables . . . . .	12
2.4.1	Decide which Continuous Variable to Use . . . . .	12
2.5	Data Binning . . . . .	12
2.5.1	The Categorical Variables . . . . .	12
2.5.2	The Continuous variables . . . . .	13
<b>3</b>	<b>The Logistic Regression</b>	<b>13</b>
<b>4</b>	<b>The performance of the Model</b>	<b>13</b>
<b>5</b>	<b>Validation of the Model</b>	<b>13</b>
5.1	Monte Carlo Cross Validation . . . . .	14
<b>6</b>	<b>The Challenger Models</b>	<b>14</b>
6.1	Neural Network . . . . .	14
6.2	Another logistic regression: logistic 2 . . . . .	14
<b>7</b>	<b>Conclusion</b>	<b>15</b>
<b>8</b>	<b>Bibliography</b>	<b>15</b>

# 1 Introduction

*Tbd...*

## 2 Data

The data is downloaded from [www.kaggle.com](http://www.kaggle.com) and delivered by an airline organization. The dataset consists of the details of customers who have already flown with them. The feedback of the customers on various context and their flight data has been consolidated. The main purpose of this dataset is to predict whether a future customer would be satisfied with their service given the details of the other parameters values. Also the airlines need to know on which aspect of the services offered by them have to be emphasized more to generate more satisfied customers. The data consists of 129880 rows and 23 columns.

Below we list all column names with explanations of the variables' meaning. For categorical variables describing satisfaction level, 0 means *Not Available* and reflects situation in which the passenger did not provide a rating.

Feature	Description	Values
Satisfaction	Airline satisfaction level	satisfied/dissatisfied
Gender	Gender of the passenger	male/female
Customer type	The customer type	loyal / disloyal customer
Age	The age of a passenger	[7; 85] years
Type of travel	Purpose of the flight	personal / business travel
Class	Travel class in the plane	business / eco / eco plus
Flight distance	The flight distance of the journey	[50; 6951] miles
Seat comfort	Satisfaction level of seat comfort	{0-5}
Departure/arrival	Satisfaction level of departure/arrival time	{0-5}
Food and drink	Satisfaction level of food and drink	{0-5}
Gate location	Satisfaction level of gate location	{0-5}

Feature	Description	Values
Inflight WiFi service	Satisfaction level of the inflight wifi service	{0-5}
Inflight entertainment	Satisfaction level of inflight entertainment	{0-5}
Online support	Satisfaction level of online support	{0-5}
Ease of online booking	Satisfaction level of online booking	{0-5}
On-board services	Satisfaction level of on-board service	{0-5}
Leg room service	Satisfaction level of leg room service	{0-5}
Baggage handling	Satisfaction level of baggage handling	{0-5}
Checkin service	Satisfaction level of check-in service	{0-5}
Cleanliness	Satisfaction level of cleanliness	{0-5}
Online boarding	Satisfaction level of online boarding	{0-5}
Departure delay in minutes	Delay upon departure	[0; 1592] minutes
Arrival delay in minutes	Delay upon arrival	[0; 1584] minutes

## 2.1 Data Preprocessing

As both academia and business point out, the data-related operations typically constitute about 80% of the whole effort of a modeling pipeline. The performance of any model is heavily driven by the quality of its inputs. It can be easily proven in a simple trial by combat that even a suboptimal model running on high quality data can oftentimes bring a sophisticated one with poor inputs to its knees. For that reason it is of utmost importance to pay extra care and attention to the data which is fed to the decision making models.

As the first step in our modeling pipeline we are going to look at the dataset

to gain insight about it's statistics and information it conveys. We'll refactor the feature names to something more manageable and represent accordingly different data types present in the dataset. Lastly we will perform quality checks on the data, such as outlier detection and treatment of not-available values (*NAs*).

### 2.1.1 Feature Encoding

Most machine learning algorithms require numerical inputs. Our data is mostly categorical and ordinal, hence we need to start with encoding those features.

**2.1.1.1 Categorical features** The dataset contains some binary categorical information such as *Male/Female*, *Loyal/Disloyal Customer*, etc. We are going to employ binary encoding for those features, that is: map values to 1 or 0 and rename the factors to `IsSatisfied`, `IsFemale`, `IsLoyal` for easier interpretation.

**2.1.1.2 Ordinal features** The main challenge of the data preparation in this dataset is the treatment of ordinal features. Take for example the `SeatNote` feature which is a customer note describing their satisfaction level with the seating arrangement. One could ask himself the following questions:

- What did the passenger have in mind? Satisfaction with seat location? Seat comfort? Possibility of choosing the seat?
- Does '*SeatNote*' = 3 imply a negative attitude towards a service? Or it's a moderate 'OK'?
- Is the satisfaction "*difference*" between notes 3 and 2 the same as between notes 5 and 4?
- Is a note '*SeatNote*' = 5 given '*Class*' = '*Eco*' the same as '*SeatNote*' = 5 given '*Class*' = '*Business*'?

The point is valid for any note-type variable in the dataset. The issue boils down to the problem that there is **no universal "unit" of satisfaction**. It is just as non-trivial to measure it as to predict it - simply because everyone perceives it in a subjective way. Our problem has an additional layer of complexity since we don't have information how precisely the survey questions have been described to the customers - so even if we *did* have some carefully

designed satisfaction unit, we cannot be sure if all respondents referred to the same aspects of service when filling out the survey.

Before discussing this further let's take a short detour to the options we have when dealing with ordinal variables for Machine Learning. Two most common approaches emerge: **Dummy encoding** and **Ordinal encoding** - both are valid, depending on what we're trying to achieve.

We could use ordinal encoding and assign numbers to each vote. This is pretty much what we already have in our "note" features. We could encode `Class` this way and assign a mapping like: {'Eco': 1, 'EcoPlus': 2, 'Business': 3}. This type of representation ensures the quality of the service is properly represented in the numeric data, but the question is whether this translates the same to the overall satisfaction? Yes, the business class is clearly more comfortable to travel in, but the *expectations* (the baseline) of business-class passengers will also be quite higher than the expectations of say, passengers in the economic class. This may result in a counterintuitive drop in the satisfaction level, simply because the sub-populations across business classes will perceive the service differently.

The second possibility we have for encoding ordinal variables is the *dummy encoding* which will split the feature `Class` into features: `Class.Eco`, `Class.EcoPlus` and `Class.Business` assigning ones and zeros in appropriate places. One of those features will be dropped to avoid perfect linear relationship (otherwise the sum of the new features would always be 1), but we'll not lose information. We only need  $n - 1$  features to encode full information about a factor with  $n$  possible levels.

We chose to employ dummy encoding to encode `Class` - to avoid making assumptions about baseline satisfaction criteria across different passenger classes. For `Note` features however this problem is non-existent, since a higher note should correspond to higher satisfaction for any rational passenger. Here to avoid inflating the dataset with additional  $4 * 14 - 14 = 42$  sparse binary columns we will stick to the original ordinal encoding. This choice nonetheless should be revisited and controlled once we reach the stage of model choice and model fitting.

The resulting, encoded dataframe looks the following way:

```
## Rows: 129,880
## Columns: 24
```

```

## $ Class.Business      <fct> 0, 1, 0, 0, 0, 0, 0...
## $ Class.EcoPlus       <fct> 0, 0, 0, 0, 0, 0, 0...
## $ Age                 <dbl> 65, 47, 15, 60, 70,...
## $ FlightDistance      <dbl> 265, 2464, 2138, 62...
## $ SeatNote            <fct> NA, NA, NA, NA, NA,...
## $ ScheduleNote        <fct> NA, NA, NA, NA, NA,...
## $ FoodNote            <fct> NA, NA, NA, NA, NA,...
## $ GateNote            <fct> 2, 3, 3, 3, 3, 3, 3...
## $ WifiNote            <fct> 2, NA, 2, 3, 4, 2, ...
## $ EntertainmentNote   <fct> 4, 2, NA, 4, 3, NA,...
## $ eSupportNote        <fct> 2, 2, 2, 3, 4, 2, 5...
## $ eBookingNote        <fct> 3, 3, 2, 1, 2, 2, 5...
## $ ServiceNote         <fct> 3, 4, 3, 1, 2, 5, 5...
## $ LegRoomNote         <fct> NA, 4, 3, NA, NA, 4...
## $ BaggageNote         <fct> 3, 4, 4, 1, 2, 5, 5...
## $ CheckInNote         <fct> 5, 2, 4, 4, 4, 5, 5...
## $ CleanNote           <fct> 3, 3, 4, 1, 2, 4, 5...
## $ eBoardingNote       <fct> 2, 2, 2, 3, 5, 2, 3...
## $ DepartureDelay      <dbl> 0, 310, 0, 0, 0, 0, 0...
## $ ArrivalDelay        <dbl> 0, 305, 0, 0, 0, 0, 0...
## $ IsSatisfied         <fct> 1, 1, 1, 1, 1, 1, 1...
## $ IsFemale            <fct> 1, 0, 1, 1, 1, 0, 1...
## $ IsLoyal             <fct> 1, 1, 1, 1, 1, 1, 1...
## $ IsPersonalTravel    <fct> 1, 1, 1, 1, 1, 1, 1...

```

### 2.1.2 Missing data treatment

In the dataset we have *NAs* for several features. We see proper *NAs* in the `ArrivalDelay` column, but there are also some “hidden” *NAs* represented by zeros, corresponding to a missing customer note. Let’s tackle that issue in this short section.

Generally speaking we don’t have any critical issue related to missing values in our data. Yes, there are *NAs* present in 14 variables, but they constitute a minuscule portion of a very large dataset (see fig. 1).

Therefore we stand before three feasible choices:

- impute the missing values

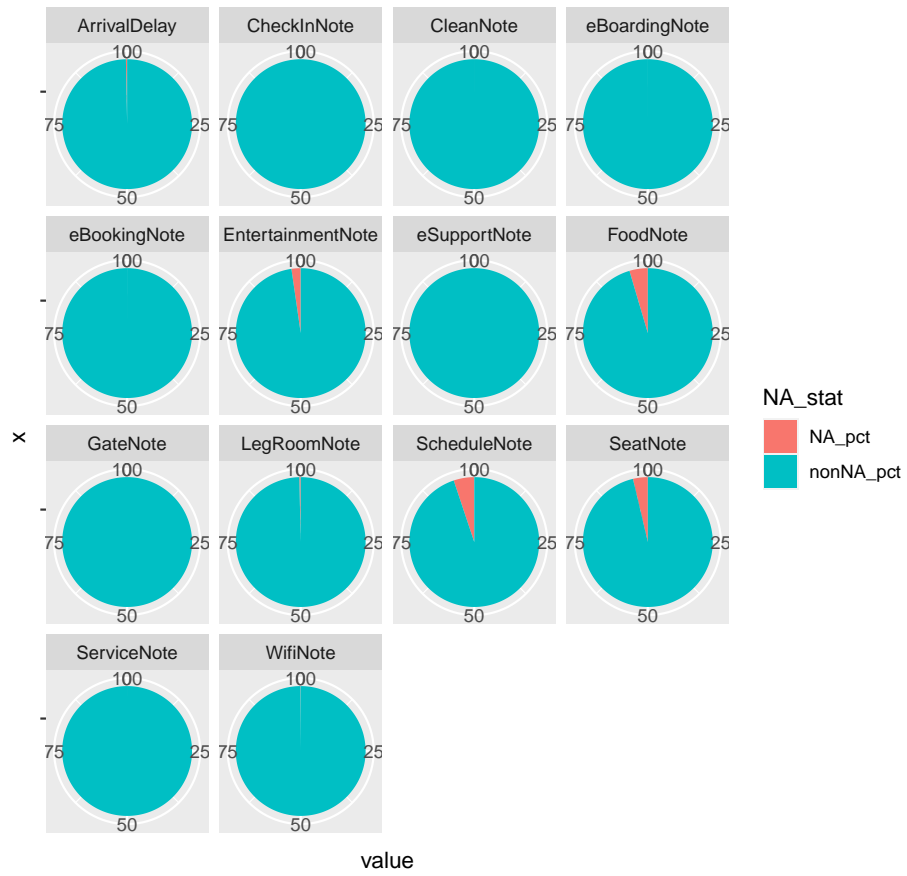


Figure 1: NA values constitute small percentage of the dataset



- drop the rows from the dataset
- discard the feature from the dataset

Handling missingness of `ArrivalDelay` turns out to be very straightforward, due to very strong linear relationship with `DepartureDelay` (see fig. 2). Looking at their scatterplot we see that we could easily impute missing values in `ArrivalDelay` by regressing it on `DepartureDelay`. It is also a very reasonable relationship, as intuitively the airplane departure delay should translate to delay in it's arrival roughly linearly. Given that the actual linear model beta is 0.9788, this imputation would be easy to justify and defend.

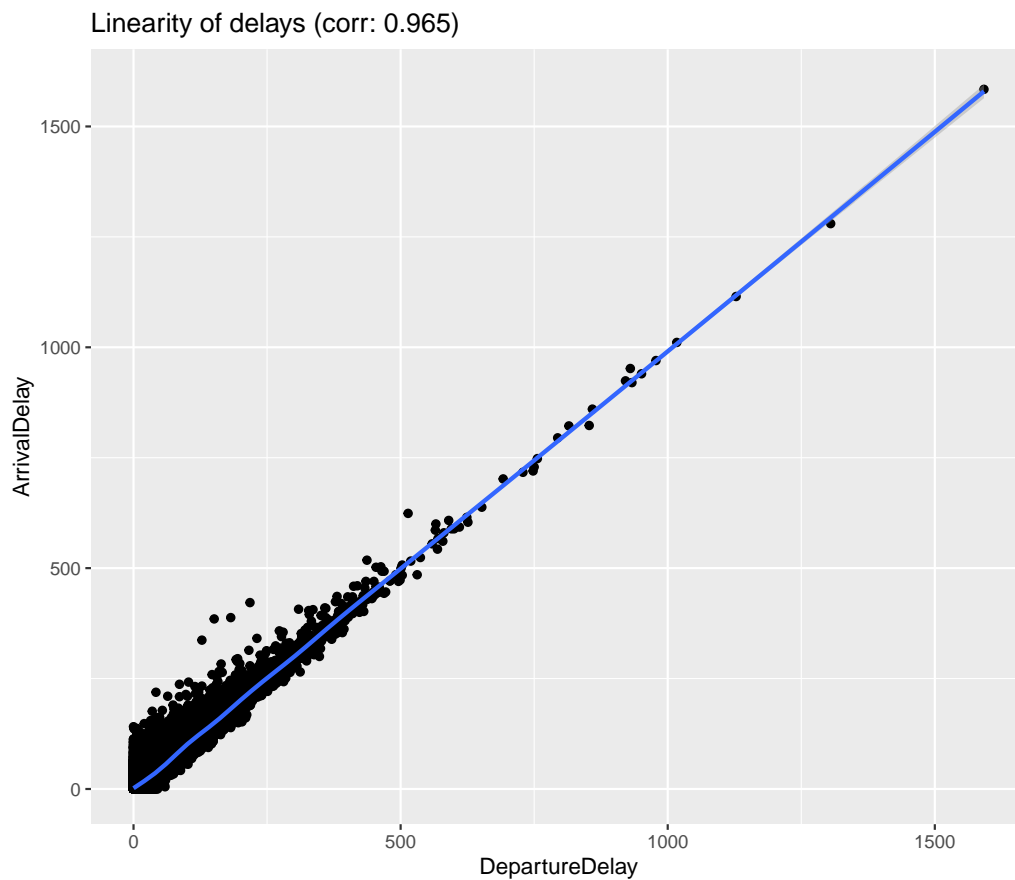


Figure 2: Strong linear relationship between departure delay and arrival delay allows to use one for imputing missing values in the other.

Therefore we could *technically* easily impute the values by regressing

**ArrivalDelay** on **DepartureDelay** - however given the high correlation of those variables (0.965 pearson correlation coefficient) one of them is bound to be dropped during multicollinearity analysis. For this reason we are not going to bother imputing the missing values, but will simply drop **ArrivalDelay** from the features at the stage of feature filtering.

Imputing the missing values in “note” variables would require much more effort though. We could take an impute-by-model approach, but that would require selecting, fitting and evaluating a multilabel response model (like multinomial or even ordinal logistic regression). We could alternatively impute by some selected data statistic, like the median. They have their pros and cons but overall, given that top *NA* percentage in a feature is 0.0513089005235602, we settled on simply dropping rows containing *NAs* if necessary. They constitute roughly 0.08 of all observations, so we would still have 119255 observations left to work with. That should be enough.

However we will hold off with the actual *NA* dropping until feature selection is finalized. Take for example the feature **ScheduleNote** - out of all the rows that would be dropped due to data missingness, almost half is caused by this feature only. That is, if **ScheduleNote** turns out to be a redundant feature and we discard it, then instead of dropping 8% of all rows, we would be dropping only 5%.

## 2.2 Exploratory Data Analysis

Describe the ideas of this section. To be done, not urgent...

Data summary

##	Age	Class	FlightDistance	SeatNote	ScheduleNote
##	Min. : 7.00	Eco :58309	Min. : 50	1 :20949	1 :20828
##	1st Qu.:27.00	Business:62160	1st Qu.:1359	4 :28398	2 :22794
##	Median :40.00	EcoPlus : 9411	Median :1925	5 :17827	3 :23184
##	Mean :39.43		Mean :1981	2 :28726	4 :29593
##	3rd Qu.:51.00		3rd Qu.:2544	3 :29183	5 :26817
##	Max. :85.00		Max. :6951	NA's: 4797	NA's: 6664
##					
##	FoodNote	GateNote	WifiNote	EntertainmentNote	eSupportNote
##	1 :21076	2 :24518	2 :27045	4 :41879	2 :17260
##	2 :27146	3 :33546	3 :27602	2 :19183	3 :21609

```

## 3 :28150 4 :30088 4 :31560 3 :24200 4 :41510
## 4 :27216 1 :22565 5 :28830 5 :29831 5 :35563
## 5 :20347 5 :19161 1 :14711 1 :11809 1 :13937
## NA's: 5945 NA's: 2 NA's: 132 NA's: 2978 NA's: 1
##
## eBookingNote ServiceNote LegRoomNote BaggageNote CheckInNote CleanNote
## 3 :22418 3 :27037 4 :39698 3:24485 5 :27005 3 :23984
## 2 :19951 4 :40675 3 :22467 4:48240 2 :15486 4 :48795
## 1 :13436 1 :13265 2 :21745 1: 7975 4 :36481 1 : 7768
## 5 :34137 2 :17174 5 :34385 2:13432 3 :35538 2 :13412
## 4 :39920 5 :31724 1 :11141 5:35748 1 :15369 5 :35916
## NA's: 18 NA's: 5 NA's: 444 NA's: 1 NA's: 5
##
## eBoardingNote DepartureDelay ArrivalDelay IsSatisfied IsFemale
## 2 :18573 Min. : 0.00 Min. : 0.00 0:58793 0:63981
## 3 :30780 1st Qu.: 0.00 1st Qu.: 0.00 1:71087 1:65899
## 5 :29973 Median : 0.00 Median : 0.00
## 4 :35181 Mean : 14.71 Mean : 15.09
## 1 :15359 3rd Qu.: 12.00 3rd Qu.: 13.00
## NA's: 14 Max. :1592.00 Max. :1584.00
## NA's :393
## IsLoyal IsPersonalTravel
## 0: 23780 0:89693
## 1:106100 1:40187
##
##
##
##
##

```

### 2.2.1 Feature Selection

The more is not always the better. Every model has a certain computational complexity that increases with the number of additional explanatory variables. The feature selection in a pre-modeling environment serves identifying groups of variables which carry repeated or very similar informational value. *Filtering feature selection methods* allow one to discard redundant features in a model independent way. By reducing the number of variables they simplify the model

and increase it's interpretability. It is also a step which tackles multicollinearity (high linear codependency of explanatory variables) which kills stability and predictive power of some models.

## **2.3 Categorical Variables**

### **2.3.1 The information Value**

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

## **2.4 The Continous Variables**

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

### **2.4.1 Decide which Continuous Variable to Use**

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

## **2.5 Data Binning**

### **2.5.1 The Categorical Variables**

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis

nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

### **2.5.2 The Continuous variables**

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

## **3 The Logistic Regression**

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

## **4 The performance of the Model**

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

## **5 Validation of the Model**

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis

nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

## **5.1 Monte Carlo Cross Validation**

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

# **6 The Challenger Models**

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

## **6.1 Neural Network**

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

## **6.2 Another logistic regression: logistic 2**

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat.

Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

## **7 Conclusion**

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

## **8 Bibliography**