



**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

01 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

1	Introduction	3
2	The Data	3
2.1	Data Preprocessing	4
2.1.1	Categorical features	5
2.1.2	Ordinal features	5
2.1.3	Exploratory Data Analysis	9
2.1.4	NA treatment	10
2.1.5	Feature Selection	11
2.2	Categorical Variables	12
2.2.1	The information Value	12
2.3	The Continous Variables	12
2.3.1	Decide which Continuous Variable to Use	12
2.4	Data Binning	12
2.4.1	The Categorical Variables	12
2.4.2	The Continuous variables	13
3	The Logistic Regression	13
4	The performance of the Model	13
5	Validation of the Model	13
5.1	Monte Carlo Cross Validation	14
6	The Challenger Models	14
6.1	Neural Network	14
6.2	Another logistic regression: logistic 2	14
7	Conclusion	15
8	Bibliography	15

1 Introduction

Tbd. . .

2 The Data

The data is downloaded from 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 in the dataset describing satisfaction level, 0 means *Not Applicable*.

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. We will perform quality checks on the data, such as outlier detection and treatment of not-available values (*NAs*). Lastly in this pre-modeling step we will look at the actual *information* carried by particular features to remove variables which are unlikely to have significant added value. The following section will be divided into parts corresponding to the data type of the features.

2.1.1 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.

2.1.2 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.

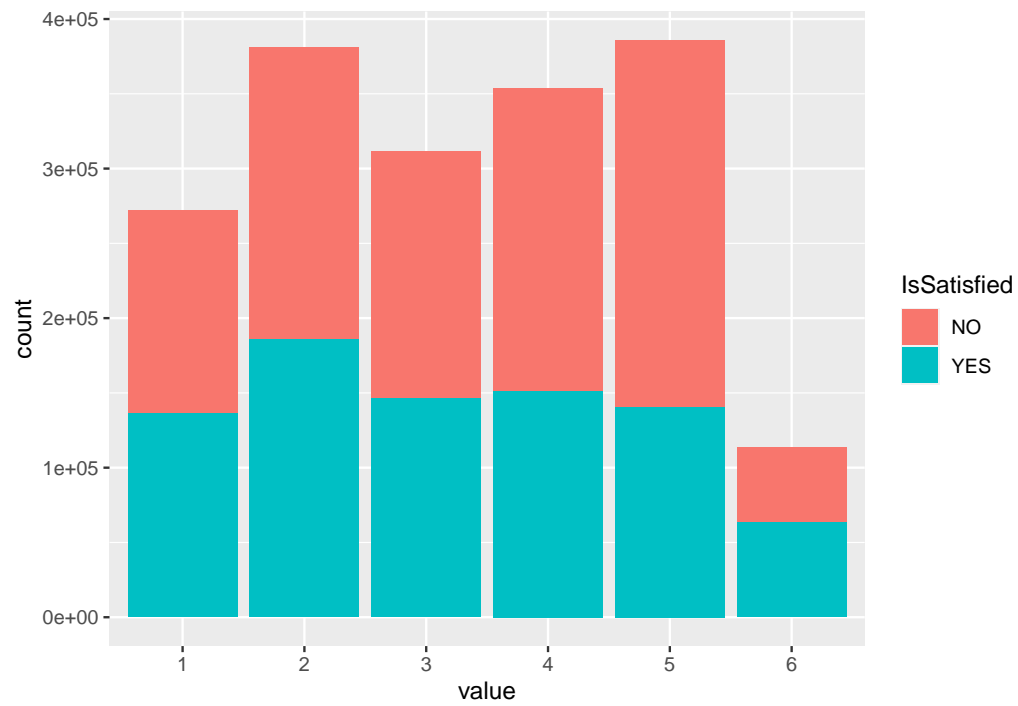
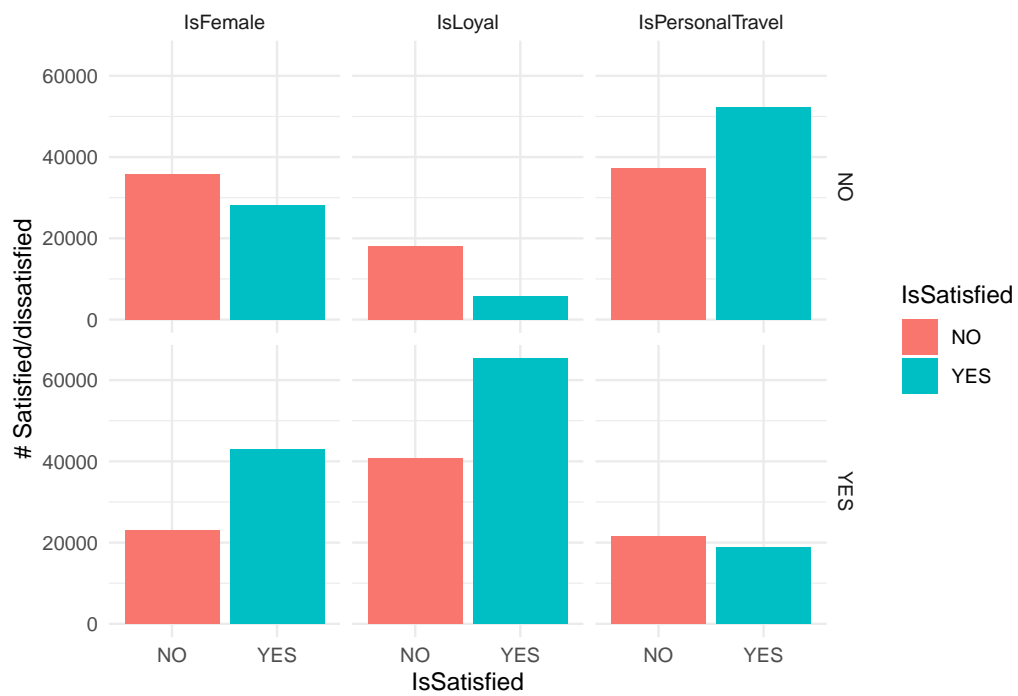
The second possibility we have for encoding ordinal variables is the *dummy encoding* which will split the feature `Class` into features: `Class.Eco`, `Class.EcoPlus` and `Class.Business` assigning ones and zeros in appropriate places. One of those features will be dropped to avoid perfect linear relationship (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, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0...
```

## \$ Class.EcoPlus	<fct> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0...
## \$ Age	<dbl> 65, 47, 15, 60, 70, 30, 66, 10, 56, 22, 58, 34, 6...
## \$ FlightDistance	<dbl> 265, 2464, 2138, 623, 354, 1894, 227, 1812, 73, 1...
## \$ SeatNote	<fct> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0...
## \$ ScheduleNote	<fct> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1...
## \$ FoodNote	<fct> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0...
## \$ GateNote	<fct> 2, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 4, 4, 1, 1, 1, 2...
## \$ WifiNote	<fct> 2, 0, 2, 3, 4, 2, 2, 2, 5, 2, 3, 2, 5, 4, 5, 1, 4...
## \$ EntertainmentNote	<fct> 4, 2, 0, 4, 3, 0, 5, 0, 3, 0, 3, 0, 0, 0, 2, 0, 0, 0...
## \$ eSupportNote	<fct> 2, 2, 2, 3, 4, 2, 5, 2, 5, 2, 3, 2, 5, 4, 1, 1, 4...
## \$ eBookingNote	<fct> 3, 3, 2, 1, 2, 2, 5, 2, 4, 2, 3, 2, 5, 4, 5, 1, 4...
## \$ ServiceNote	<fct> 3, 4, 3, 1, 2, 5, 5, 3, 4, 2, 3, 3, 1, 3, 5, 3, 4...
## \$ LegRoomNote	<fct> 0, 4, 3, 0, 0, 4, 0, 3, 0, 4, 0, 2, 3, 5, 0, 4, 4...
## \$ BaggageNote	<fct> 3, 4, 4, 1, 2, 5, 5, 4, 1, 5, 1, 5, 2, 2, 5, 1, 1...
## \$ CheckInNote	<fct> 5, 2, 4, 4, 4, 5, 5, 5, 5, 3, 2, 2, 2, 3, 2, 4, 3...
## \$ CleanNote	<fct> 3, 3, 4, 1, 2, 4, 5, 4, 4, 4, 3, 5, 4, 2, 5, 2, 1...
## \$ eBoardingNote	<fct> 2, 2, 2, 3, 5, 2, 3, 2, 4, 2, 5, 2, 5, 4, 2, 1, 4...
## \$ DepartureDelay	<dbl> 0, 310, 0, 0, 0, 0, 17, 0, 0, 30, 47, 0, 0, 0, 40...
## \$ ArrivalDelay	<dbl> 0, 305, 0, 0, 0, 0, 15, 0, 0, 26, 48, 0, 0, 0, 48...
## \$ IsSatisfied	<fct> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1...
## \$ IsFemale	<fct> 1, 0, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1...
## \$ IsLoyal	<fct> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1...
## \$ IsPersonalTravel	<fct> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1...



2.1.3 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      0: 4797      0: 6664
##  1st Qu.:27.00      Business:62160      1st Qu.:1359      1:20949      1:20828
##  Median :40.00      EcoPlus : 9411      Median :1925      4:28398      2:22794
##  Mean    :39.43                                Mean    :1981      5:17827      3:23184
##  3rd Qu.:51.00                                3rd Qu.:2544      2:28726      4:29593
##  Max.    :85.00                                Max.    :6951      3:29183      5:26817
##
##  FoodNote  GateNote  WifiNote  EntertainmentNote eSupportNote eBookingNote
##  0: 5945    2:24518    2:27045    4:41879                2:17260        3:22418
##  1:21076    3:33546    0: 132     2:19183                3:21609        2:19951
##  2:27146    4:30088    3:27602    0: 2978                4:41510        1:13436
##  3:28150    1:22565    4:31560    3:24200                5:35563        5:34137
##  4:27216    5:19161    5:28830    5:29831                1:13937        4:39920
##  5:20347    0: 2        1:14711    1:11809                0: 1           0: 18
##
##  ServiceNote LegRoomNote BaggageNote CheckInNote CleanNote eBoardingNote
##  3:27037      0: 444        3:24485      5:27005      3:23984      2:18573
##  4:40675      4:39698        4:48240      2:15486      4:48795      3:30780
##  1:13265      3:22467        1: 7975      4:36481      1: 7768      5:29973
##  2:17174      2:21745        2:13432      3:35538      2:13412      4:35181
##  5:31724      5:34385        5:35748      1:15369      5:35916      1:15359
##  0: 5         1:11141                0: 1         0: 5         0: 14
##
##  DepartureDelay  ArrivalDelay  IsSatisfied IsFemale IsLoyal
##  Min.      : 0.00  Min.      : 0.00  0:58793     0:63981     0: 23780
##  1st Qu.: 0.00  1st Qu.: 0.00  1:71087     1:65899     1:106100
##  Median : 0.00  Median : 0.00
##  Mean    : 14.71  Mean    : 15.09
##  3rd Qu.: 12.00  3rd Qu.: 13.00
##  Max.    :1592.00  Max.    :1584.00
##
##  NA's      :393
##  IsPersonalTravel
```

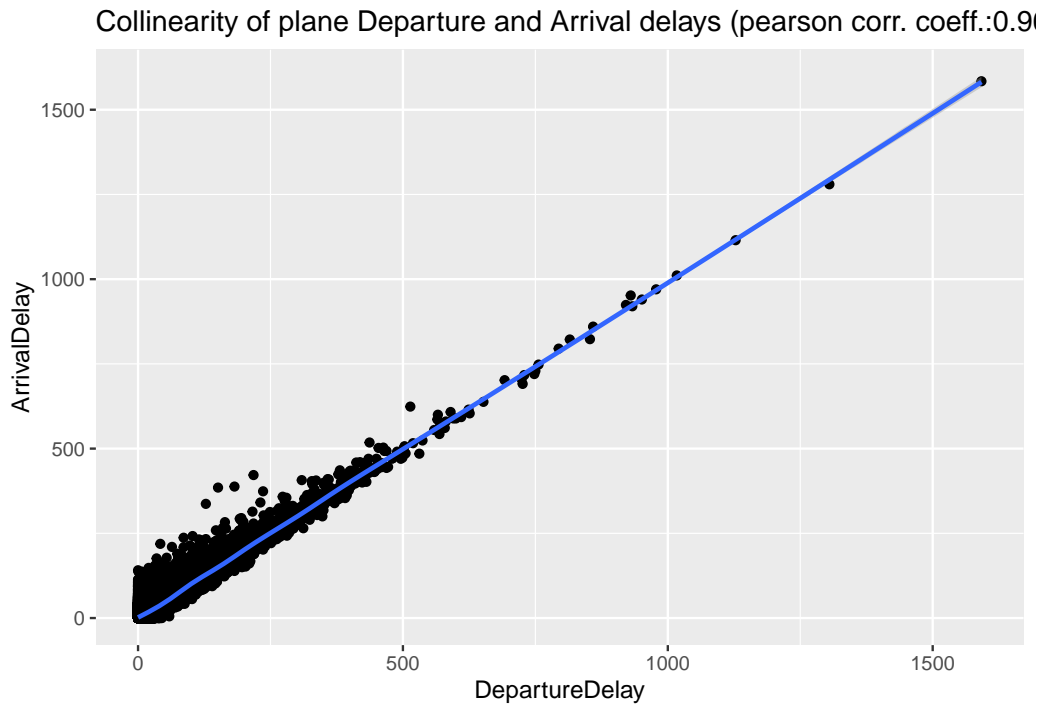
```
## 0:89693
## 1:40187
##
##
##
##
##
##
```

2.1.4 NA treatment

In the dataset we only have *NAs* for the `ArrivalDelay` feature. As there is an option to assign 0 to the delay, and there is non-zero data about `DepartureDelay` we believe these *NAs* are a genuine data loss. We stand before three possible choices:

- impute the missing values
- drop them from the dataset
- discard the feature from the dataset

Technically, we could do all of them. Looking however at the scatterplot of `ArrivalDelay` vs `DepartureDelay` we see these variables are extremely highly correlated, especially in their positive tails. This seems in line with intuition as the airplane departure delay should translate to delay in it's arrival roughly linearly.



Therefore we could *technically* easily impute the values by regressing `ArrivalDelay` on `DepartureDelay` - however given the high correlation of those variables (0.9653 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.

2.1.5 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.

2.2 Categorical Variables

2.2.1 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.

2.3 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.

2.3.1 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.

2.4 Data Binning

2.4.1 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.

2.4.2 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