



**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

12 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	Data	3
2.1	Exploratory Data Analysis	4
2.2	Data Preprocessing	11
2.2.1	Missing data treatment	11
2.2.2	Feature Engineering	12
2.2.3	Feature Encoding	19
3	A train-test split of data	21
4	Baseline Model	21
4.1	Fitting and performance	21
4.2	Validation	21
5	The Challenger Models	22
5.1	Random Forest	22
5.1.1	Fitting and performance	22
5.1.2	Validation	22
5.2	Logistic Regression	22
5.2.1	Fitting and performance	22
5.2.2	Validation	31
5.3	State-of-the-art: Neural Network model	33
5.3.1	Fitting and performance	33
5.3.2	Validation	34
6	Conclusion	34
7	Bibliography	34

1 Introduction

Tbd...

2 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 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 Exploratory Data Analysis

Before data wrangling and trying to fit some models, it is great to get an overview of our data set and verify whether it makes sense. We start from reviewing the few top lines of our data and summarizing it to see basic statistical analysis. For each quantitative variable we get the minimum, maximum, mean, median and the inter quartile range information. However, for each categorical variable we get the number of observations in each category.

```
## # A tibble: 6 x 23
##   satisfaction Gender `Customer Type`    Age `Type of Travel` `Class
```

```

##   <fct>      <fct>  <fct>          <dbl> <fct>      <fct>
## 1 satisfied  Female Loyal Customer  65 Personal Travel Eco
## 2 satisfied  Male   Loyal Customer  47 Personal Travel Busi~
## 3 satisfied  Female Loyal Customer  15 Personal Travel Eco
## 4 satisfied  Female Loyal Customer  60 Personal Travel Eco
## 5 satisfied  Female Loyal Customer  70 Personal Travel Eco
## 6 satisfied  Male   Loyal Customer  30 Personal Travel Eco
## # ... with 17 more variables: `Flight Distance` <dbl>, `Seat comfort` <fct>,
## # `Departure/Arrival time convenient` <fct>, `Food and drink` <fct>, `Gate
## # location` <fct>, `Inflight wifi service` <fct>, `Inflight
## # entertainment` <fct>, `Online support` <fct>, `Ease of Online
## # booking` <fct>, `On-board service` <fct>, `Leg room service` <fct>,
## # `Baggage handling` <fct>, `Checkin service` <fct>, Cleanliness <fct>,
## # `Online boarding` <fct>, `Departure Delay in Minutes` <dbl>, `Arrival Delay
## # in Minutes` <dbl>

##      satisfaction     Gender      Customer Type      Age
## satisfied    :71087  Female:65899  Loyal Customer  :106100  Min.   : 7.00
## dissatisfied:58793  Male  :63981  disloyal Customer: 23780  1st Qu.:27.00
##                                         Median :40.00
##                                         Mean   :39.43
##                                         3rd Qu.:51.00
##                                         Max.   :85.00
##
##      Type of Travel     Class      Flight Distance  Seat comfort
## Personal Travel:40187  Eco   :58309  Min.   : 50   0: 4797
## Business travel:89693 Business:62160  1st Qu.:1359   1:20949
##                           Eco Plus: 9411  Median :1925   4:28398
##                                         Mean   :1981   5:17827
##                                         3rd Qu.:2544   2:28726
##                                         Max.   :6951   3:29183
##
##      Departure/Arrival time convenient Food and drink Gate location
## 0: 6664                               0: 5945       2:24518
## 1:20828                             1:21076       3:33546
## 2:22794                             2:27146       4:30088
## 3:23184                             3:28150       1:22565
## 4:29593                             4:27216       5:19161

```

```

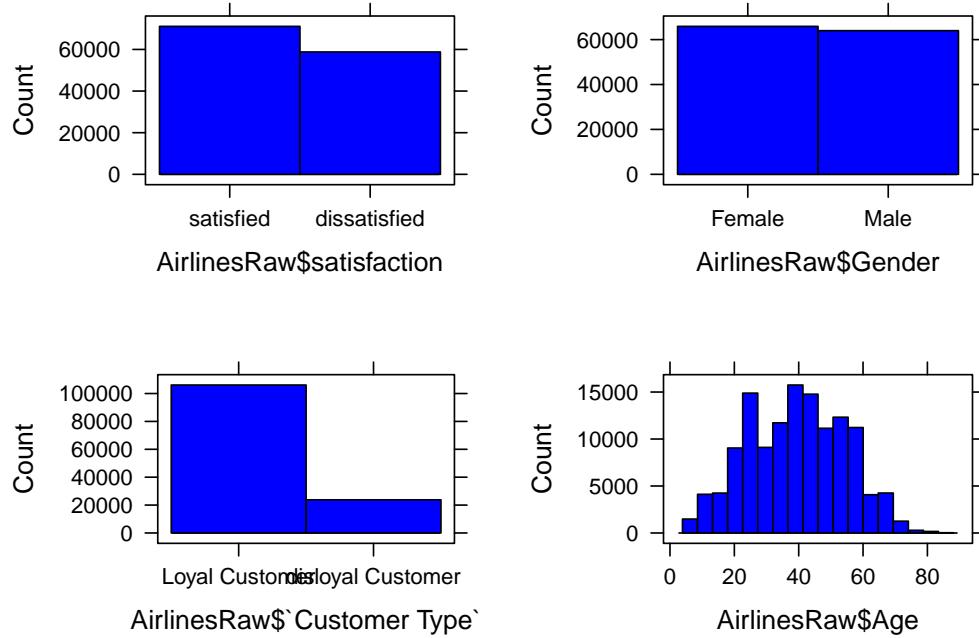
## 5:26817           5:20347      0: 2
##
## Inflight wifi service Inflight entertainment Online support
## 2:27045          4:41879      2:17260
## 0: 132           2:19183      3:21609
## 3:27602          0: 2978       4:41510
## 4:31560          3:24200      5:35563
## 5:28830          5:29831      1:13937
## 1:14711          1:11809      0: 1
##
## Ease of Online booking On-board service Leg room service Baggage handling
## 3:22418          3:27037      0: 444       3:24485
## 2:19951          4:40675      4:39698      4:48240
## 1:13436          1:13265      3:22467      1: 7975
## 5:34137          2:17174      2:21745      2:13432
## 4:39920          5:31724      5:34385      5:35748
## 0: 18            0: 5         1:11141
##
## Checkin service Cleanliness Online boarding Departure Delay in Minutes
## 5:27005          3:23984      2:18573      Min.   : 0.00
## 2:15486          4:48795      3:30780      1st Qu.: 0.00
## 4:36481          1: 7768       5:29973      Median  : 0.00
## 3:35538          2:13412       4:35181      Mean    : 14.71
## 1:15369          5:35916       1:15359      3rd Qu.: 12.00
## 0: 1              0: 5         0: 14        Max.   :1592.00
##
## Arrival Delay in Minutes
## Min.   : 0.00
## 1st Qu.: 0.00
## Median : 0.00
## Mean   : 15.09
## 3rd Qu.: 13.00
## Max.   :1584.00
## NA's   :393

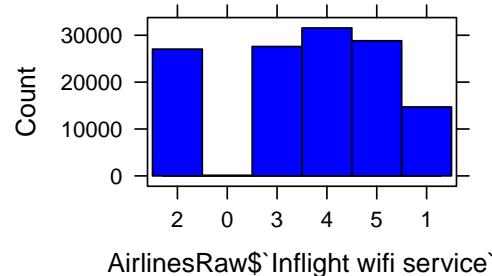
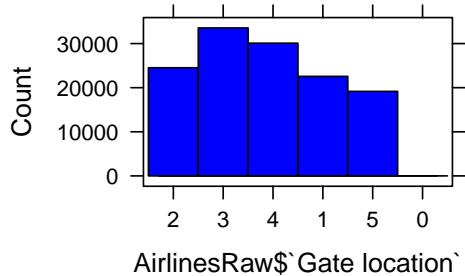
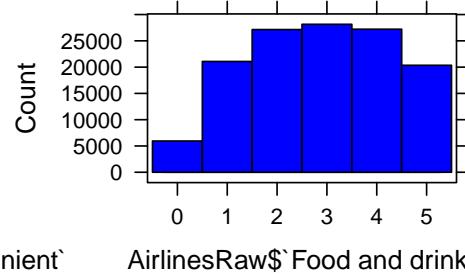
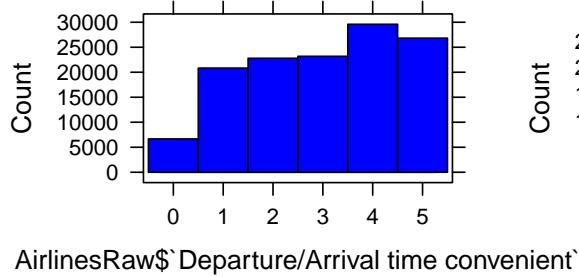
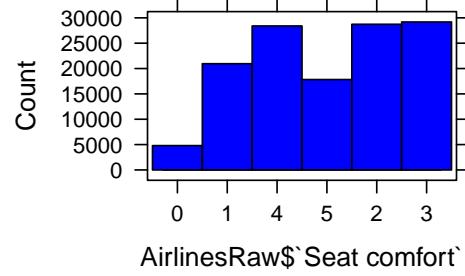
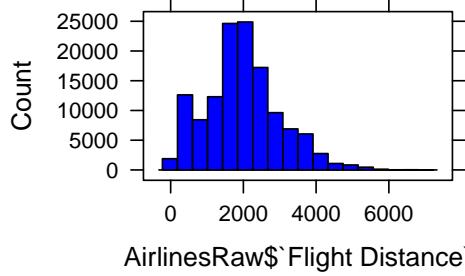
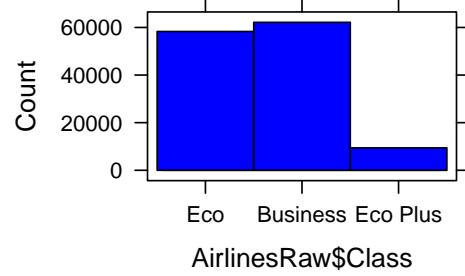
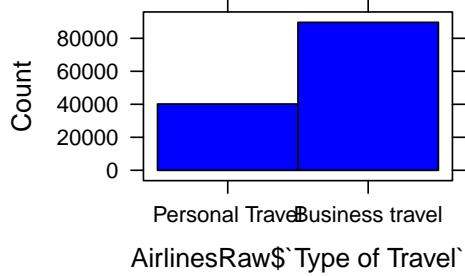
```

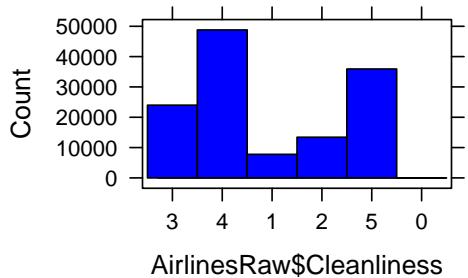
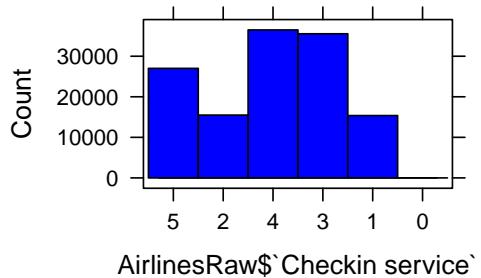
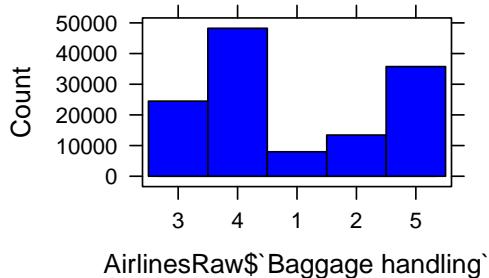
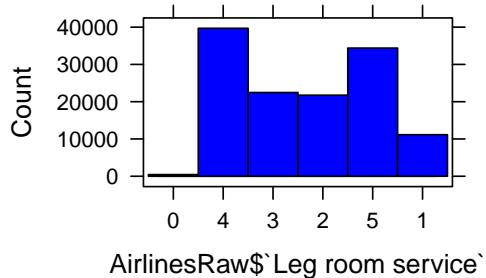
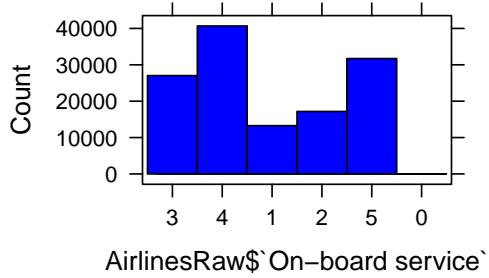
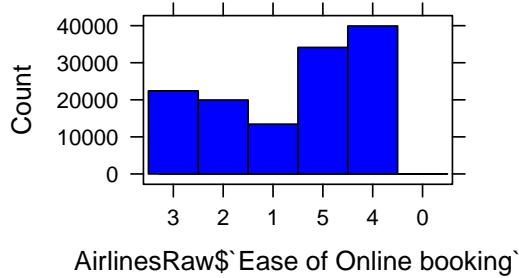
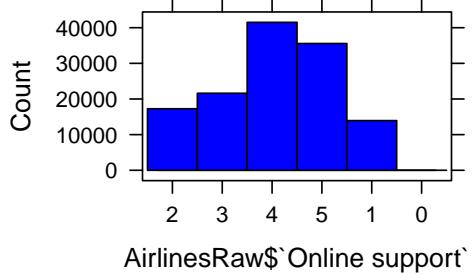
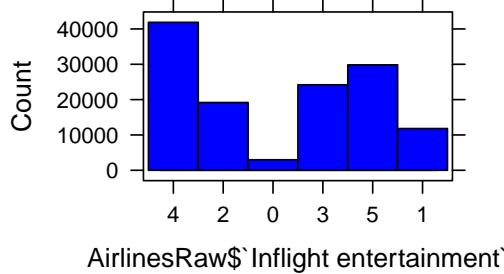
In this project we will predict the satisfaction level of a customer in function of the other variables based on past data. All results seems to be reasonable, but we notice that there are missing values in some columns. However, we

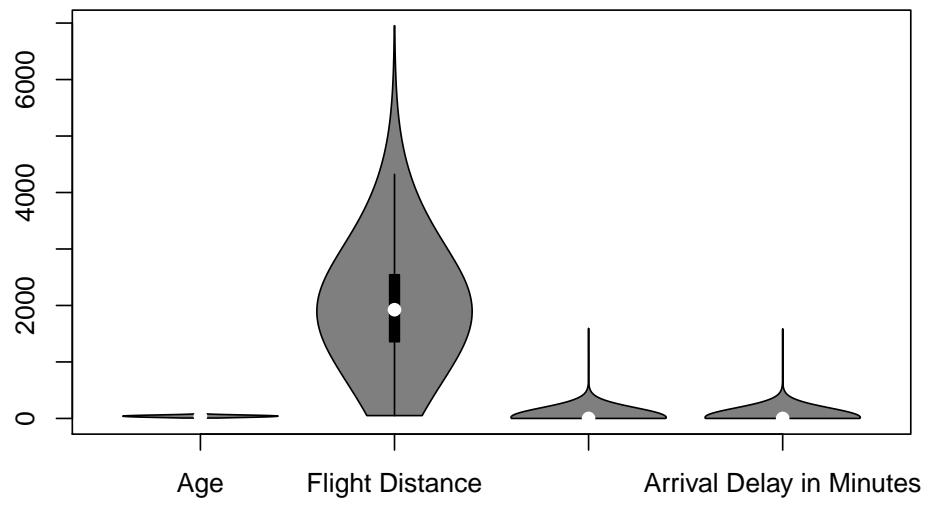
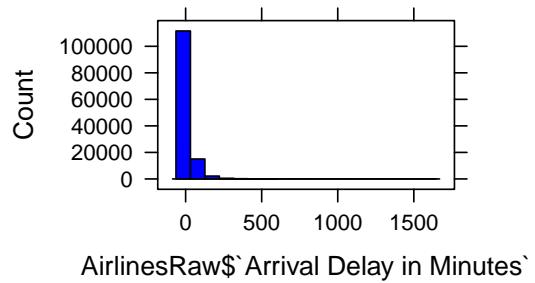
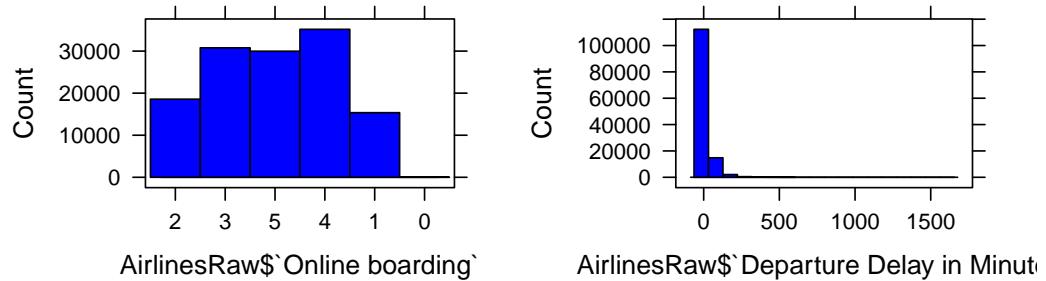
will take care about it later.\

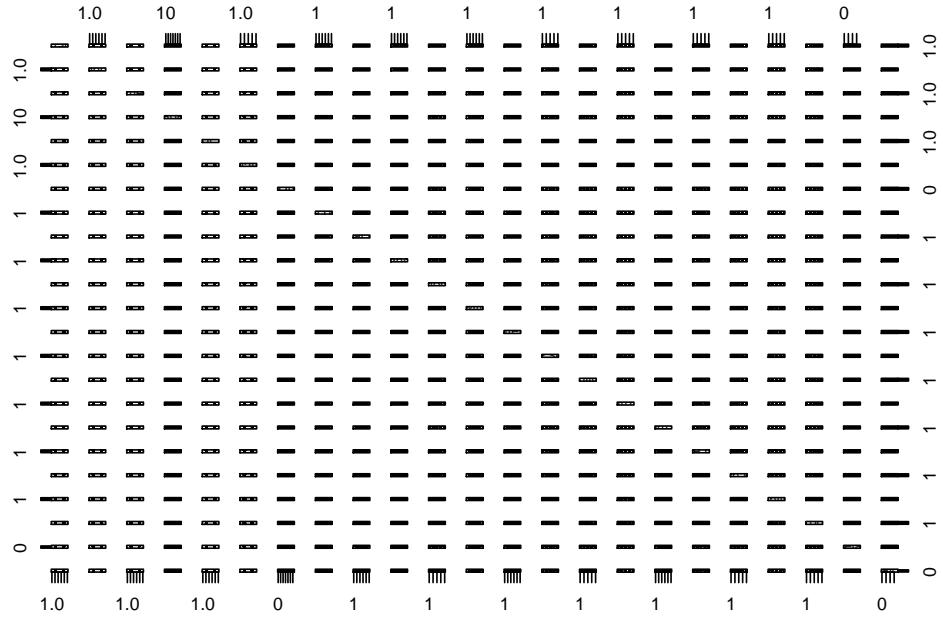
Additionally, we plot some histograms and vioplots to view overall data distribution and scatterplots for relationships between variables.











2.2 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.

During data preprocessing step we will gain insight about data statistics and information it conveys. First, we'll deal with NAs and outliers. Then we will refactor the features with variable binning and select a subset of features which are likely to display predictive power for our problem. In the end we will encode the variables and prepare a train-test split for model development.

2.2.1 Missing data treatment

Table 2: NA breakdown per feature. NAs span a small portion of data

Feature	NA_count	pct_of_data
ScheduleNote	6664	5.13%
FoodNote	5945	4.58%
SeatNote	4797	3.69%
EntertainmentNote	2978	2.29%
LegRoomNote	444	0.34%
ArrivalDelay	393	0.3%
WifiNote	132	0.1%
eBookingNote	18	0.01%
eBoardingNote	14	0.01%
CleanNote	5	0%
ServiceNote	5	0%
GateNote	2	0%
CheckInNote	1	0%
eSupportNote	1	0%

After examining the data it seems we don't have any critical issue related to missing values. *NAs* are present in 12 variables, but they constitute a minuscule portion of a very large dataset (see fig. 2). We considered employing an imputation strategy based on median, but given that NAs constitute roughly 0.08 of all observations, even if we drop them we would still have 119255 observations left to work with. Based on that we decided not to introduce imputed values to the dataset, but rather work with pure data.

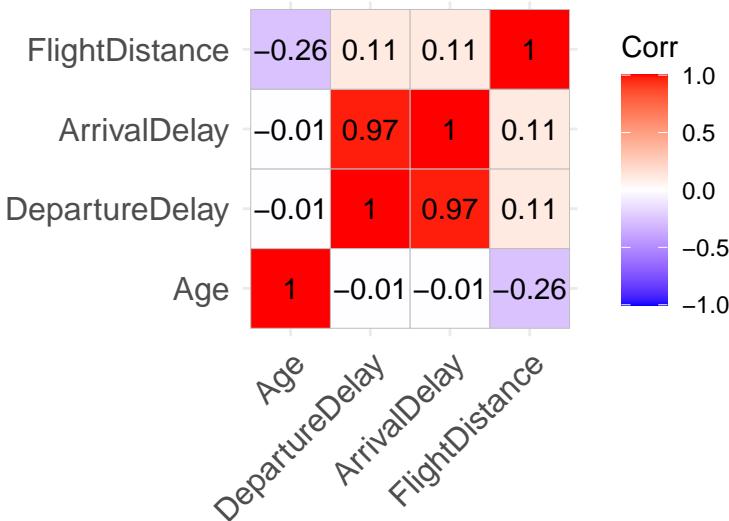
2.2.2 Feature Engineering

The more is not always the better. Feature engineering is a pre-modeling stage which serves identifying features which are significant and filtering out the ones that are not. *Filtering feature selection methods* allow one to discard redundant features in a model independent way. By reducing the number of variables or discretizing them we simplify the model and increase its interpretability. It is also a step which tackles multicollinearity (high linear codependency of explanatory variables) which kills stability and predictive

power of some models.

Across the following sections we are going to introduce a few additional modifications of variables to the dataset and see if it makes sense to keep them. Since we have a big share of 0s in column `DepartureDelay`, we wanted to check if it makes sense to include a binary variable `IsDelayed` as a predictor. We are also going to try discretizing `Age`, `DepartureDelay` and `FlightDistance` continuous variables into bins based on Weight of Evidence metric and evaluate their predictive power.

2.2.2.1 Continuous Features We have four continuous variables in our dataset: `Age`, `DepartureDelay`, `ArrivalDelay` and `FlightDistance`. We start the analysis by analyzing their codependence structure. We note a high linear relationship between `ArrivalDelay` and `DepartureDelay`, visible both in the correlation matrix (fig. ??) and on figure 1. We can safely drop `ArrivalDelay`, since it doesn't introduce new information and additionally contaminates the dataset with NAs.



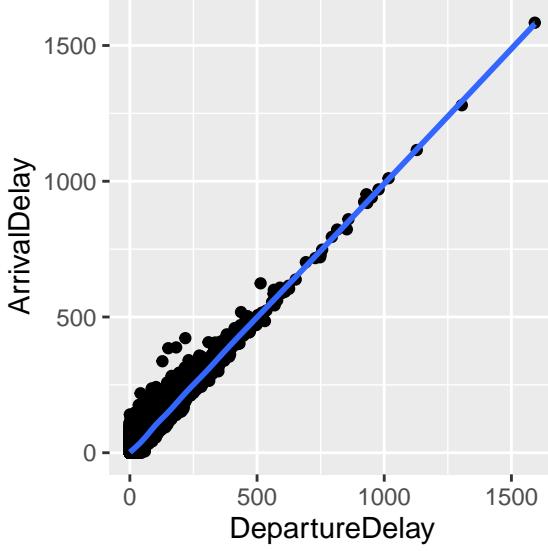
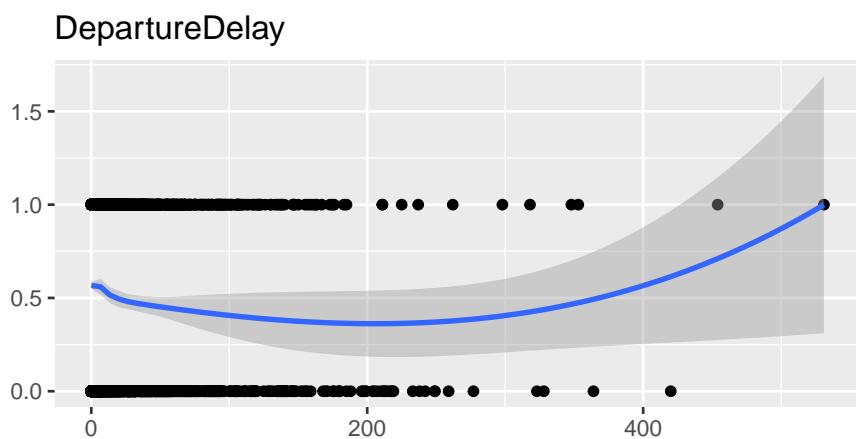
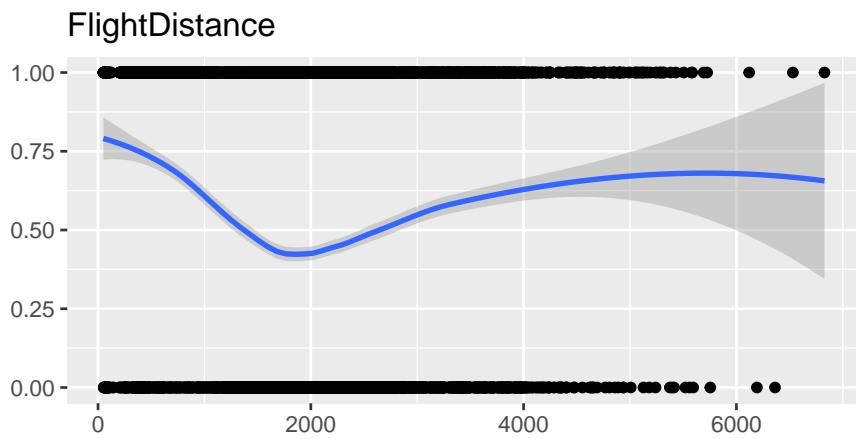
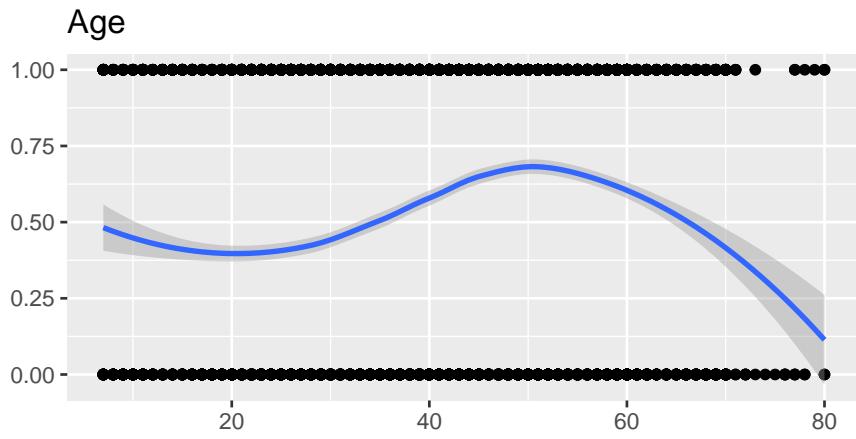


Figure 1: Strong linear relationship between departure delay and arrival delay allows to drop one of them from the dataset.

Next, we're going to examine the loess estimator of satisfaction as a function of the remaining continuous variables to see if any of them looks flat enough to raise suspicion regarding its utility. Flatness of loess implies that average satisfaction does not change in the explanatory variable, hence the explanatory variable doesn't convey much information. In our case, this is not visible on the figure ??, so we cannot discard any feature based on that. Note, the plot has been generated on a randomized data sub-sample for computational complexity reduction. We made sure however to take a sample large enough, so that the standard errors are tamed, and to ensure the relationship shape is stable regardless of the random seed chosen.



Now we are going to look at possible binnings of our features. We'll use `woeBinning::woe.binning` function which chooses the binning to maximize the information value of the feature. If the optimized binning will yield $IV < 0.1$, we will discard the variable. Otherwise we'll analyze the bins to ensure they are not over-optimized to an unreasonable degree. Based on loess plot shapes/regimes the expectation is to have not more than four bins for `Age` and `FlightDistance`, and a maximum of two bins for `DepartureDelay`.

WOE Table for Age

	Final.Bin	Total.Count	Total.Distr.	0.Rate	WOE	IV
1	<= 28	30896	25.9%	58.6%	-51.4	0.068
2	<= 40	29312	24.6%	49.0%	-12.6	0.004
3	<= 59	47672	40.0%	33.9%	50.1	0.096
4	<= Inf	11375	9.5%	53.2%	-29.3	0.008
6	Total	119255	100.0%	45.9%	NA	0.177

WOE Table for FlightDistance

	Final.Bin	Total.Count	Total.Distr.	0.Rate	WOE	IV
1	<= 1359	29839	25.0%	33.3%	53.0	0.067
2	<= 3053	71534	60.0%	53.4%	-30.2	0.055
3	<= Inf	17882	15.0%	36.7%	38.0	0.021
5	Total	119255	100.0%	45.9%	NA	0.143

WOE Table for DepartureDelay

	Final.Bin	Total.Count	Total.Distr.	0.Rate	WOE	IV
1	<= 19	95920	80.4%	44.0%	7.6	0.005
2	<= Inf	23335	19.6%	53.6%	-31.1	0.019
4	Total	119255	100.0%	45.9%	NA	0.024

From the WOE tables above we see that `DepartureDelay` is a variable of low predictive power, hence we won't use it in modelling. For the other variables, as the data binning chosen by the algorithm seems reasonable given

the ex-ante expectations, we're going to keep them.

2.2.2.2 Ordinal & Categorical Features The main challenge of the data preparation in this dataset is the proper treatment of passenger notes. 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 their 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 same 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 satisfaction 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.

We aim to overcome this problem, by binning notes into wider classes, depending on how well they explain and affect the overall satisfaction. Since we have a lot of 5-leveled `Note` factors, there's a strong suspicion that in such large set there must exist some adjacent levels such that overall satisfaction is invariant to displacements in that group of levels. In other words, we could collapse notes of 1, 2&3 to one group if they carried similar information. Hence we will again let `woe.binning` automatically select bins and then verify the result.

Since there are 12 `Note` variables, we will only display one of the WOE tables as an example. However for all it has been verified that **adjacent** levels have been binned, so the binning is plausible, and the *IV* of the newly binned features are above 0.1.

	Final.Bin	Total.Cou	Total.Distr	1.Count	0.Count	1.Distr.	0.Distr.	0.Rate	WOE	IV
1	11713	9.8%		1561	10152	2.4%	18.6%	86.7%	-	0.329
									203.8	
3 + 2	38010	31.9%		11063	26947	17.1%	49.3%	70.9%	-	0.340
									105.6	
5 + 4	69532	58.3%		51945	17587	80.4%	32.2%	25.3%	91.7	0.443
Total	119255	100.0%		64569	54686	100.0%	100.0%	45.9%	NA	1.111

Next, we will check the information value for other categorical variables. They are binary, so binning has not been applied to them in the earlier step. The table 7 presents features and their IV. We see that there are features with IV smaller than 0.1 - we are going to drop those from the dataset.

Lastly, a brief look at the spearman rank correlation matrix (fig. 2) shows that there are no highly correlated “note” features among ordinal variables in the dataset.

Table 7: Information Value for all variables.

varName	IV
IsPersonalTravel	0.0532949
FlightDistance	0.1433763
Age	0.1769584
IsFemale	0.1893377
FoodNote	0.2330999
WifiNote	0.2920904
CheckInNote	0.3607163
Class	0.4001969
IsLoyal	0.4537176
CleanNote	0.4555277
BaggageNote	0.4718849
LegRoomNote	0.5714645
eBoardingNote	0.5948280
ServiceNote	0.6151573
eSupportNote	0.9223110
eBookingNote	1.1114520
SeatNote	1.4504018

varName	IV
EntertainmentNote	2.2947658

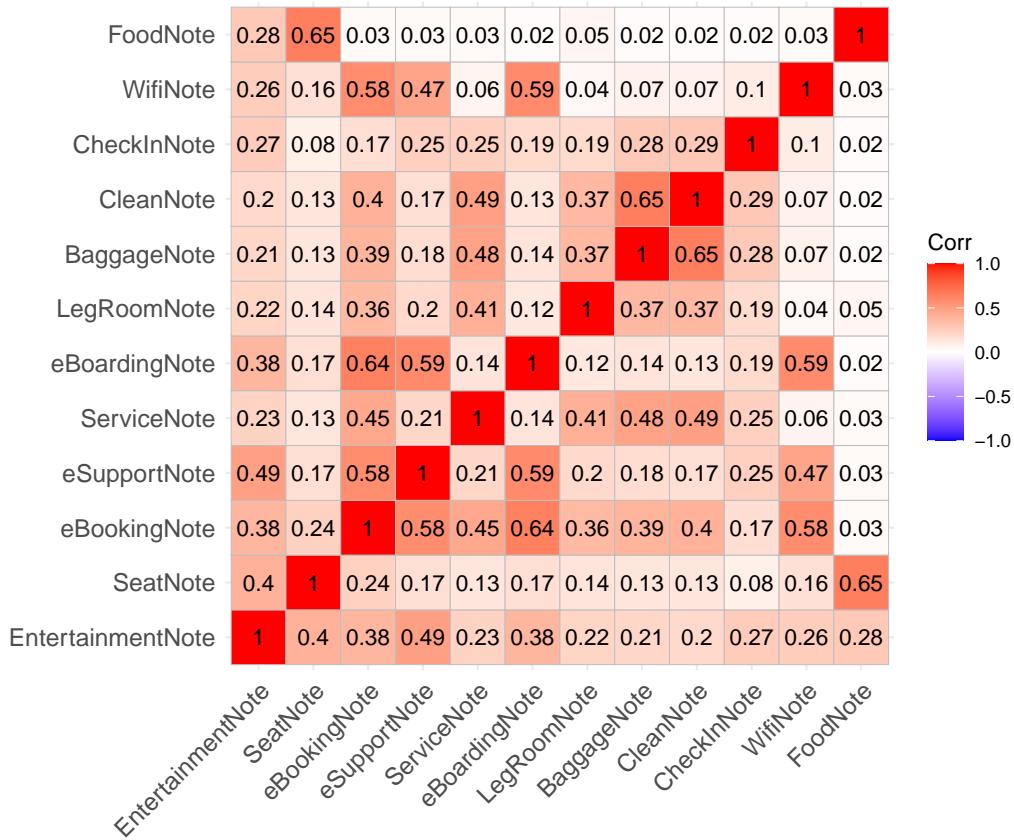


Figure 2: Spearman correlation shows no significant colinear relationships in ordinal variables

2.2.3 Feature Encoding

Some machine learning algorithms require numerical data, so we considered **ordinal encoding** and **dummy encoding** to transform our data.

We ran the following thought experiment to determine which encoding to employ. Say we use ordinal encoding and assign numbers to each factor

level. 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.

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 $4 \cdot 14 - 14 = 42$ 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:

```

## $ CheckInNote.L      <fct> 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0,...  

## $ CheckInNote.H     <fct> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,...  

## $ WifiNote.H        <fct> 0, 1, 0, 0, 1, 1, 1, 0, 1, 0, 1, 1, 1, 0, 1, 0,...  

## $ WifiNote.M        <fct> 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0,...  

## $ FoodNote.M        <fct> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,...  

## $ FoodNote.H        <fct> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,...  

## $ Age.30s            <fct> 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0,...  

## $ Age.40s50s         <fct> 1, 1, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1,...  

## $ Age.60plus          <fct> 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0,...  

## $ FlightDistance.M   <fct> 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0,...  

## $ FlightDistance.H   <fct> 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1,...  

## $ Class.Business      <fct> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0,...  

## $ Class.EcoPlus       <fct> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1,...  

## $ IsSatisfied         <fct> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,...  

## $ IsFemale             <fct> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,...  

## $ IsLoyal              <fct> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,...
```

3 A train-test split of data

4 Baseline Model

4.1 Fitting and performance

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.

4.2 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.

5 The Challenger Models

5.1 Random Forest

5.1.1 Fitting and performance

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.

5.1.2 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.

5.2 Logistic Regression

5.2.1 Fitting and performance

The next model that we will study is a logistic regression. We name it for future reference *logistic1*.

```
##  
## Call:  
## glm(formula = frm, family = "binomial", data = Airlines_binned_train)  
##  
## Deviance Residuals:  
##      Min        1Q    Median        3Q       Max  
## -4.1699  -0.3623   0.0313   0.3099   3.6301  
##
```

```

## Coefficients:
##                               Estimate Std. Error z value Pr(>|z|)
## (Intercept)           -7.21962   0.08599 -83.961 < 2e-16 ***
## ClassBusiness          1.34302   0.02605  51.550 < 2e-16 ***
## ClassEcoPlus           0.03494   0.04385   0.797 0.425568
## IsFemale1              0.97658   0.02355  41.466 < 2e-16 ***
## IsLoyal1                1.91910   0.03742  51.283 < 2e-16 ***
## Age30s                 -0.02214   0.03196 -0.693 0.488560
## Age40s50s               0.21030   0.03105   6.774 1.25e-11 ***
## Age60plus              -0.37541   0.04313 -8.704 < 2e-16 ***
## FlightDistanceM        -0.19019   0.02919 -6.516 7.21e-11 ***
## FlightDistanceH         -0.13764   0.03939 -3.494 0.000475 ***
## EntertainmentNoteM      1.82628   0.02847  64.152 < 2e-16 ***
## EntertainmentNoteH      2.99102   0.04336  68.979 < 2e-16 ***
## SeatNoteM                0.85488   0.03330  25.670 < 2e-16 ***
## SeatNoteH                5.13222   0.11098  46.246 < 2e-16 ***
## eBookingNoteH            1.64792   0.06294  26.180 < 2e-16 ***
## eBookingNoteM            0.85178   0.06052  14.073 < 2e-16 ***
## eSupportNoteM            0.22928   0.03314   6.918 4.59e-12 ***
## eSupportNoteH             0.47526   0.03731  12.739 < 2e-16 ***
## ServiceNoteH              0.72232   0.04724  15.290 < 2e-16 ***
## ServiceNoteM              0.22391   0.04594   4.874 1.09e-06 ***
## eBoardingNoteH            0.19938   0.03785  5.267 1.38e-07 ***
## eBoardingNoteM            0.32782   0.03447   9.509 < 2e-16 ***
## LegRoomNoteM              0.03922   0.04776   0.821 0.411536
## LegRoomNoteH              0.75533   0.04772  15.827 < 2e-16 ***
## BaggageNoteM              0.10721   0.03359   3.192 0.001412 **
## BaggageNoteH              0.46100   0.03863  11.932 < 2e-16 ***
## CleanNoteM                0.06595   0.03497   1.886 0.059276 .
## CleanNoteH                0.48344   0.03958  12.214 < 2e-16 ***
## CheckInNoteM              0.33644   0.02874  11.706 < 2e-16 ***
## CheckInNoteH              0.90613   0.03728  24.307 < 2e-16 ***
## WifiNoteH                 -0.13901   0.05062 -2.746 0.006025 **
## WifiNoteM                  0.12726   0.04902   2.596 0.009435 **
## FoodNoteM                 -0.48163   0.03440 -14.000 < 2e-16 ***
## FoodNoteH                 -0.51322   0.04304 -11.924 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

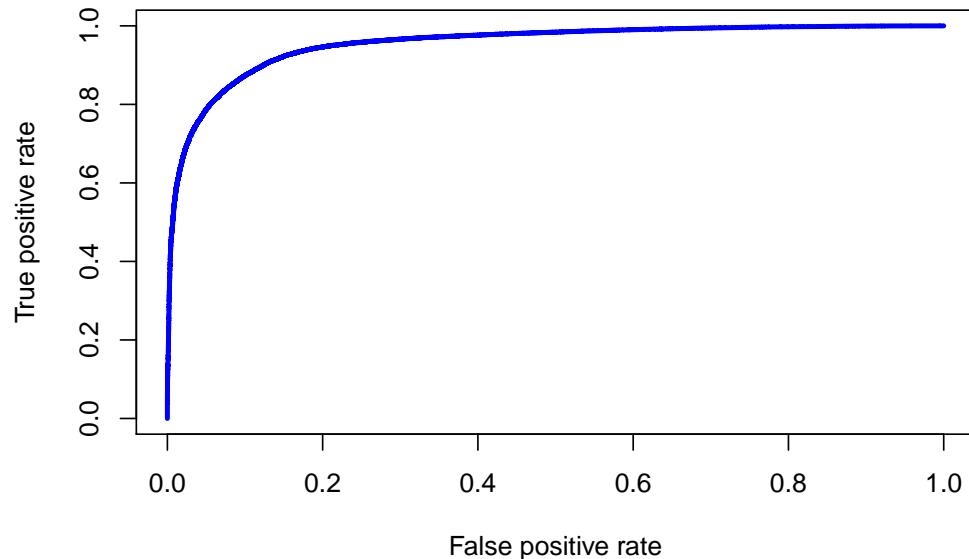
```

```

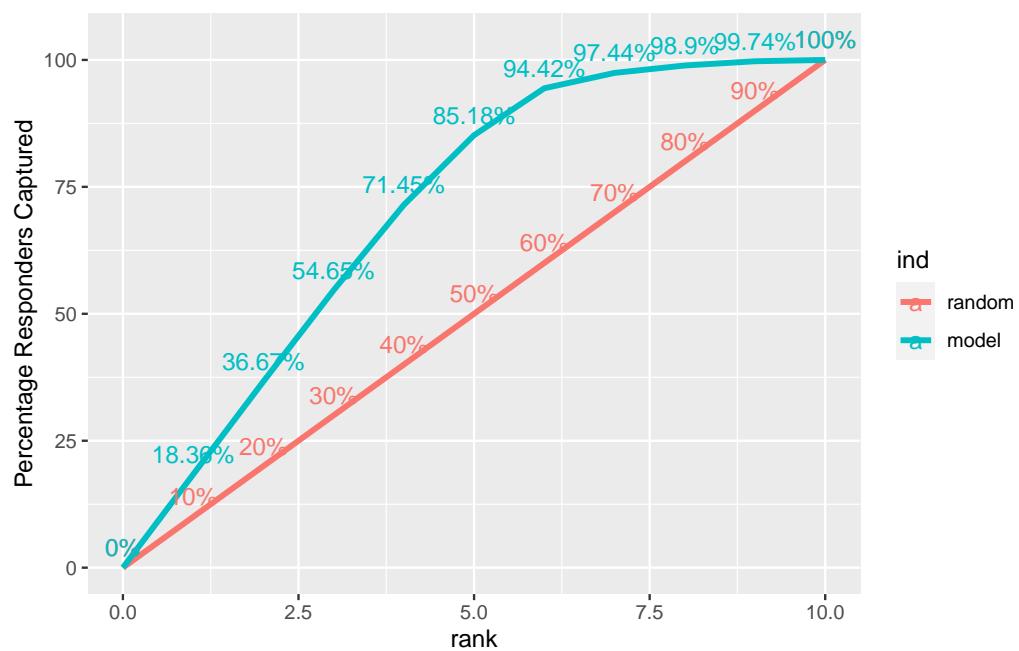
## 
## (Dispersion parameter for binomial family taken to be 1)
## 
##     Null deviance: 131599  on 95403  degrees of freedom
## Residual deviance:  51904  on 95370  degrees of freedom
## AIC: 51972
## 
## Number of Fisher Scoring iterations: 7

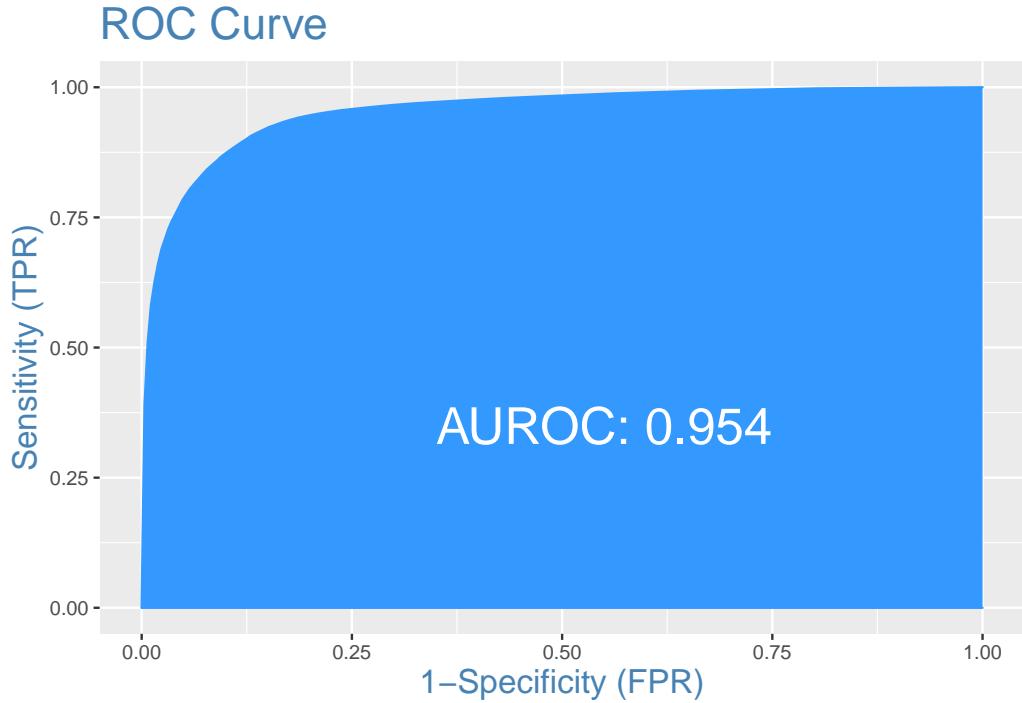
```

We notice that according to this model • business class passengers have a much higher satisfaction level than eco plus class passengers, who in their turn have a higher satisfaction level than eco class ones • satisfaction level is higher for females and loyal customers • passengers at the age of 40s and 50s have a higher satisfaction level than the others • satisfaction level is lower for long and medium flight distances than for the short ones • the higher entertainment notes, seat notes, eBooking notes, eSupport notes, eBoarding, leg room notes, clean notes and baggage notes, the higher the satisfaction level • satisfaction level is lower for low and medium service notes than for the high ones • satisfaction level is higher for high check-in notes and lower for the low ones • satisfaction level is higher for medium wifi notes than for the high ones • satisfaction level is lower for high food notes than for the medium ones



KS Plot





The AUC on the training data is 0.95429 and on the testing data 0.95410. The difference is very small. This model is a good contender. However, we will try to add some interactions between variables.

```

## 
## Call:
## glm(formula = frm2, family = "binomial", data = Airlines_binned_train)
## 
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max 
## -4.1233  -0.3476   0.0299   0.2984   3.7349 
## 
## Coefficients:
##                               Estimate Std. Error z value Pr(>|z|)    
## (Intercept)              -7.031166  0.101184 -69.489 < 2e-16 ***
## ClassBusiness             1.466122  0.072754  20.152 < 2e-16 ***
## ClassEcoPlus            -0.897868  0.277094  -3.240 0.001194 ** 
## IsFemale1                 1.029808  0.075041  13.723 < 2e-16 ***
## IsLoyal1                  0.817303  0.070453  11.601 < 2e-16 ***

```

## Age30s	0.237101	0.047282	5.015	5.31e-07	***
## Age40s50s	0.554496	0.044481	12.466	< 2e-16	***
## Age60plus	-0.341481	0.066539	-5.132	2.87e-07	***
## FlightDistanceM	0.034027	0.031187	1.091	0.275251	
## FlightDistanceH	0.030433	0.040761	0.747	0.455291	
## EntertainmentNoteM	1.739130	0.029019	59.930	< 2e-16	***
## EntertainmentNoteH	2.794782	0.044255	63.152	< 2e-16	***
## SeatNoteM	0.892443	0.034044	26.214	< 2e-16	***
## SeatNoteH	5.133180	0.110901	46.286	< 2e-16	***
## eBookingNoteH	1.533933	0.063705	24.079	< 2e-16	***
## eBookingNoteM	0.836671	0.060976	13.721	< 2e-16	***
## eSupportNoteM	0.183299	0.033562	5.462	4.72e-08	***
## eSupportNoteH	0.433249	0.038099	11.372	< 2e-16	***
## ServiceNoteH	0.687019	0.048365	14.205	< 2e-16	***
## ServiceNoteM	0.200990	0.046923	4.283	1.84e-05	***
## eBoardingNoteH	0.271834	0.037990	7.155	8.34e-13	***
## eBoardingNoteM	0.323874	0.034892	9.282	< 2e-16	***
## LegRoomNoteM	0.002897	0.048535	0.060	0.952397	
## LegRoomNoteH	0.666962	0.048623	13.717	< 2e-16	***
## BaggageNoteM	0.123945	0.034418	3.601	0.000317	***
## BaggageNoteH	0.528297	0.039999	13.208	< 2e-16	***
## CleanNoteM	0.118070	0.035968	3.283	0.001028	**
## CleanNoteH	0.589432	0.041125	14.333	< 2e-16	***
## CheckInNoteM	0.359854	0.029073	12.378	< 2e-16	***
## CheckInNoteH	0.976777	0.038301	25.503	< 2e-16	***
## WifiNoteH	-0.049311	0.050554	-0.975	0.329355	
## WifiNoteM	0.139034	0.048974	2.839	0.004527	**
## FoodNoteM	-0.489478	0.034879	-14.033	< 2e-16	***
## FoodNoteH	-0.508576	0.043207	-11.771	< 2e-16	***
## IsFemale1:IsLoyal1	1.218629	0.071359	17.078	< 2e-16	***
## ClassBusiness:IsFemale1	-1.606214	0.050757	-31.645	< 2e-16	***
## ClassEcoPlus:IsFemale1	-0.375600	0.091745	-4.094	4.24e-05	***
## ClassBusiness:IsLoyal1	0.811235	0.073118	11.095	< 2e-16	***
## ClassEcoPlus:IsLoyal1	1.219838	0.271800	4.488	7.19e-06	***
## IsFemale1:Age30s	-0.363369	0.065223	-5.571	2.53e-08	***
## IsFemale1:Age40s50s	-0.437266	0.061794	-7.076	1.48e-12	***
## IsFemale1:Age60plus	0.114841	0.088785	1.293	0.195847	
## ---					

```

## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 131599  on 95403  degrees of freedom
## Residual deviance: 50155  on 95362  degrees of freedom
## AIC: 50239
##
## Number of Fisher Scoring iterations: 7

```

We notice that all of the interactions are statistically significant and decide to keep them. Moreover, we see that FlightDistance is now not significant, so we leave it out of the model.

```

##
## Call:
## glm(formula = frm3, family = "binomial", data = Airlines_binned_train)
##
## Deviance Residuals:
##      Min        1Q     Median        3Q       Max
## -4.1196   -0.3477    0.0299    0.2982    3.7351
##
## Coefficients:
##                               Estimate Std. Error z value Pr(>|z|)
## (Intercept)                 -6.996915  0.096219 -72.719 < 2e-16 ***
## ClassBusiness                1.460452  0.072497  20.145 < 2e-16 ***
## ClassEcoPlus                -0.895669  0.277068 -3.233 0.001226 **
## IsFemale1                   1.026007  0.074930 13.693 < 2e-16 ***
## IsLoyal1                     0.813975  0.070292 11.580 < 2e-16 ***
## Age30s                       0.234788  0.047216  4.973 6.60e-07 ***
## Age40s50s                    0.550496  0.044222 12.448 < 2e-16 ***
## Age60plus                    -0.346570  0.066341 -5.224 1.75e-07 ***
## EntertainmentNoteM           1.738686  0.029015 59.923 < 2e-16 ***
## EntertainmentNoteH            2.793664  0.044233 63.159 < 2e-16 ***
## SeatNoteM                     0.892204  0.034039 26.211 < 2e-16 ***
## SeatNoteH                     5.133044  0.110899 46.286 < 2e-16 ***
## eBookingNoteH                1.533019  0.063699 24.067 < 2e-16 ***
## eBookingNoteM                0.836308  0.060977 13.715 < 2e-16 ***
## eSupportNoteM                0.183350  0.033550  5.465 4.63e-08 ***

```

```

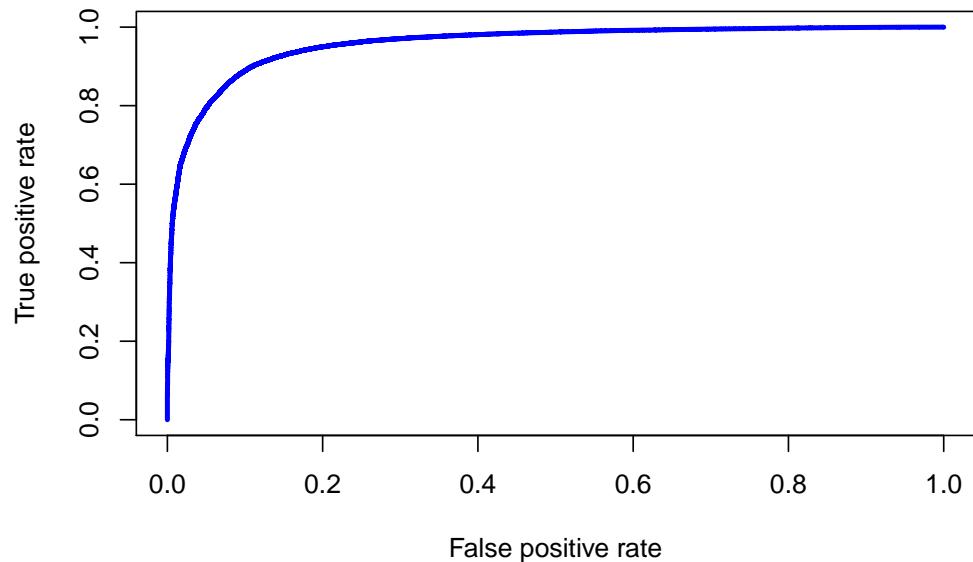
## eSupportNoteH          0.432847  0.038075 11.368 < 2e-16 ***
## ServiceNoteH          0.686734  0.048355 14.202 < 2e-16 ***
## ServiceNoteM          0.201015  0.046921  4.284 1.83e-05 ***
## eBoardingNoteH         0.272988  0.037976  7.188 6.55e-13 ***
## eBoardingNoteM         0.324184  0.034892  9.291 < 2e-16 ***
## LegRoomNoteM          0.002799  0.048534  0.058 0.954013
## LegRoomNoteH          0.666710  0.048624 13.712 < 2e-16 ***
## BaggageNoteM          0.124644  0.034413  3.622 0.000292 ***
## BaggageNoteH          0.529257  0.039988 13.235 < 2e-16 ***
## CleanNoteM             0.118476  0.035964  3.294 0.000987 ***
## CleanNoteH             0.590200  0.041116 14.355 < 2e-16 ***
## CheckInNoteM           0.360304  0.029071 12.394 < 2e-16 ***
## CheckInNoteH           0.977256  0.038299 25.517 < 2e-16 ***
## WifiNoteH              -0.048253 0.050546 -0.955 0.339765
## WifiNoteM              0.140380  0.048962  2.867 0.004142 **
## FoodNoteM              -0.489124 0.034875 -14.025 < 2e-16 ***
## FoodNoteH              -0.508233 0.043205 -11.763 < 2e-16 ***
## IsFemale1:IsLoyal1     1.218020  0.071343 17.073 < 2e-16 ***
## ClassBusiness:IsFemale1 -1.595037 0.049728 -32.075 < 2e-16 ***
## ClassEcoPlus:IsFemale1 -0.376573 0.091734 -4.105 4.04e-05 ***
## ClassBusiness:IsLoyal1   0.811640  0.072942 11.127 < 2e-16 ***
## ClassEcoPlus:IsLoyal1   1.218229  0.271791  4.482 7.39e-06 ***
## IsFemale1:Age30s        -0.365510 0.065177 -5.608 2.05e-08 ***
## IsFemale1:Age40s50s     -0.449299 0.060823 -7.387 1.50e-13 ***
## IsFemale1:Age60plus      0.100356  0.087826  1.143 0.253180
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 131599  on 95403  degrees of freedom
## Residual deviance: 50156  on 95364  degrees of freedom
## AIC: 50236
##
## Number of Fisher Scoring iterations: 7

```

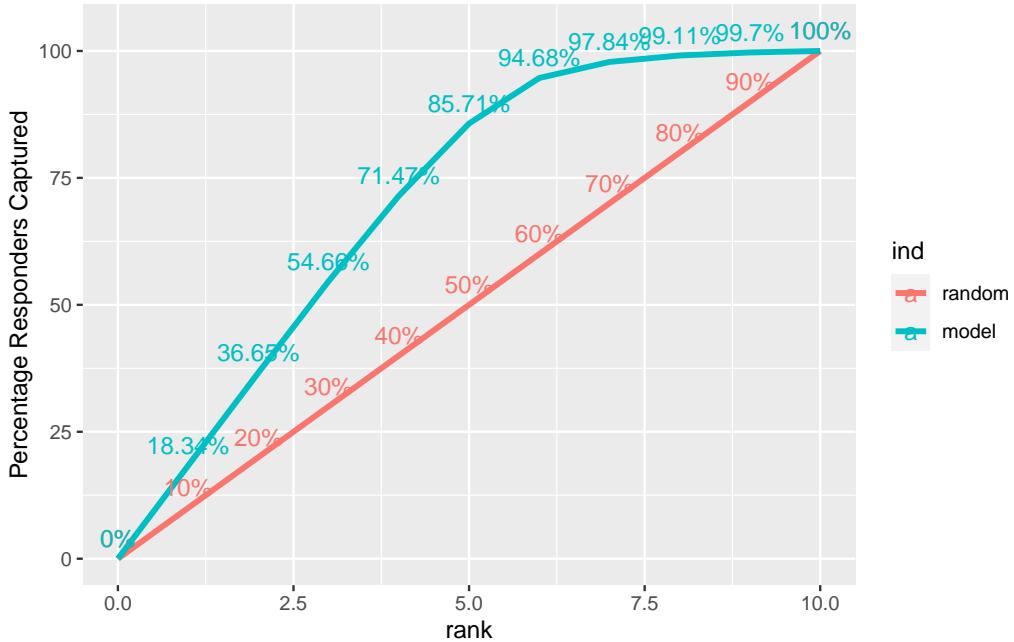
We notice that according to this model • business class passengers have a much higher satisfaction level than eco plus class passengers, who in their

turn have a higher satisfaction level than eco class ones • satisfaction level is higher for females and loyal customers • passengers at the age of 40s and 50s have a higher satisfaction level than the others • satisfaction level is lower for long and medium flight distances than for the short ones • the higher entertainment notes, seat notes, eBooking notes, eSupport notes, leg room notes, clean notes and baggage notes, the higher the satisfaction level • satisfaction level is lower for low and medium service notes that for the high ones • satisfaction level is higher for high check-in notes and lower for the low ones • satisfaction level is higher for medium wifi notes than for the high ones • satisfaction level is lower for high food notes than for the medium ones • females travelling business class are less satisfied than the ones from eco plus class. Analogically, for loyal customers.

Now we will check the performance of the model using few criteria.



KS Plot



The third model has an AUC of 0.9564 on the training data and 0.9574 on the testing data. The difference is very small. The KS is 0.79096. Moreover, we can see (on the ROC curve graph) that our classifier is not far from being perfect. These findings are great and we will proceed now to calculate the optimal cutoff.

```
## [1] 0.4899762
```

The optimal cut-off (that gives the minimum mis-classification error) is 0.48998. Above this cutoff level we assume satisfaction.

5.2.2 Validation

To validate the model we will use the cross validation method and opt for the Monte Carlo Cross Validation.

We use the Monte Carlo Cross Validation with a test data-set that spans 30% of our observations and 70% in the training data-set. We will draw 200 times a training data-set of 0.7 and study the AUC on the testing data-set.

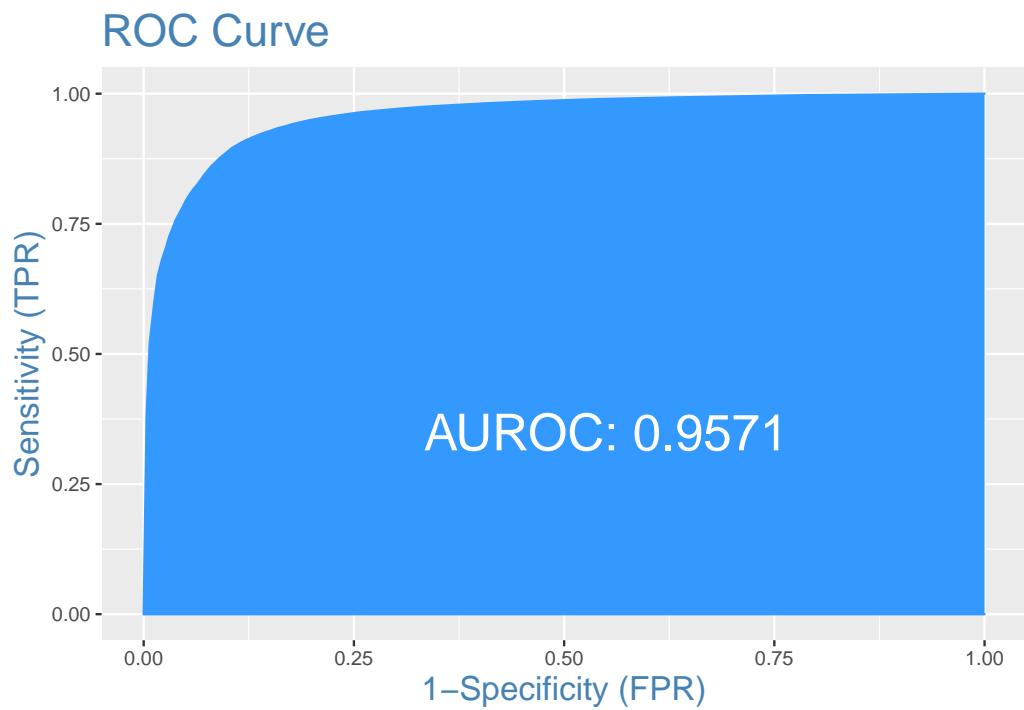
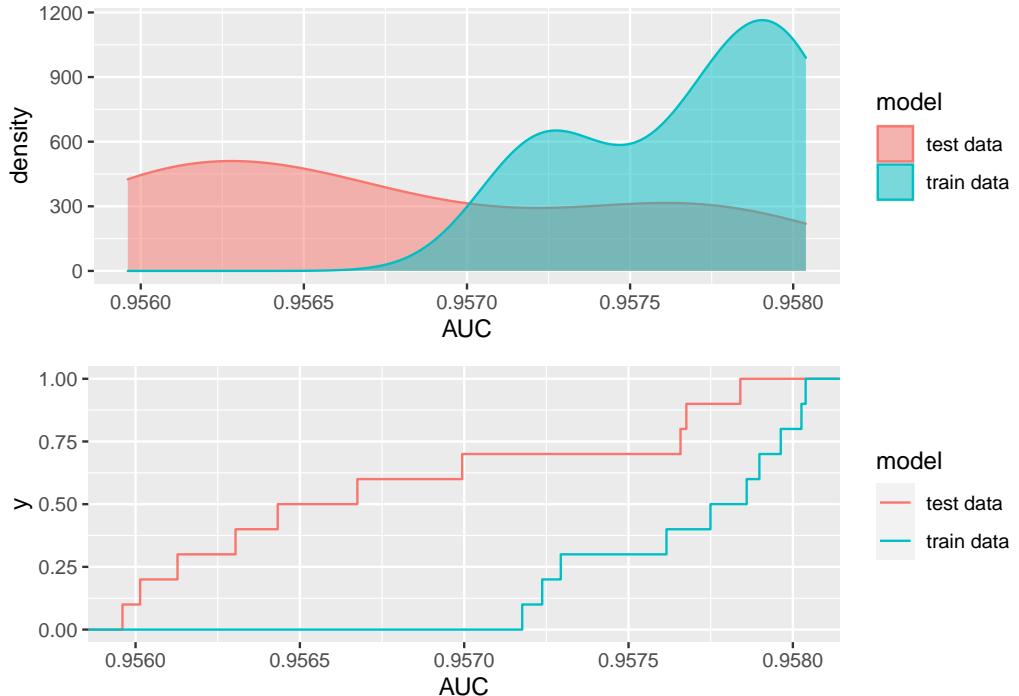


Figure 3: The ROC (receiver operating curve) for our model



```
## [1] 0.9567682
## [1] 0.0007291288
```

We notice that the model is a little over-fit, but even for the test data-set the performance is acceptable. All variables retained have at least one of the categories with a $p - value$ that is smaller than 0.01. We could consider to leave out the categories that have a higher $p - value$ to make the model more robust. We choose, however, not to do that because all coefficients are somehow logical.\ The median of the observed values for the AUC of the test data is 0.9571835, the average is 0.9571316 with a standard deviation of 0.0008522686.

5.3 State-of-the-art: Neural Network model

5.3.1 Fitting and performance

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.

5.3.2 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 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.

7 Bibliography