

MSCI Coursework Part 1

Rivyesch Ranjan – 36330520
(Count – 2007 words)

Executive Summary

The study analysed and visualized the data “AirlinesData81” which comprised of past customers’ personal and flight details along with their survey responses. EDA and visualization techniques were used to understand the relationship between the various features and customers’ satisfaction, and discern the customers’ perception towards each of the service quality factors. Feature selection techniques were used to identify and drop features that were either redundant or had low discriminatory power over satisfaction. MDS was then performed on the remaining data and the features that heavily influence customers’ satisfaction were inflight WIFI service, onboard service, leg room, baggage handling and inflight service.

1 Introduction

In the highly competitive aviation (airline) industry, customers’ overall flight experience is the biggest differentiating factor between airlines. Besides flight safety, enhancing service quality and therefore customer satisfaction is the most crucial strategies of the airlines. Higher customer satisfaction increases customer loyalty which results in customers who are willing to purchase more and spend more on products and services.

An airline in the UK is conducting a study regarding its customer satisfaction. The aim of the study is to understand the main factors that influence customer satisfactions. Using the insights and findings from the study, the airline will then be able to target specific areas of their service that led to unhappy customers during previous travels. By improving these areas of their service, the airline would potentially be able to increase overall customer satisfaction.

The structure of the report is split into three main sections. Exploratory data analysis covers the dataset and visual exploration of the features. The next section discusses feature selection techniques used to identify important features, and redundant and multicollinear features that can be dropped. The final main section covers multidimensional scaling to visualise the multidimensional data in two-dimensional space and draw insights.

2 Exploratory Data Analysis (EDA)

2.1 Data

“AirlinesData81” consists of 10313 instances of customer’s flight information and experience. The data contains 23 features in total, some numerical and others categorical. The target variable was the customers’ satisfaction, either satisfied or neutral/dissatisfied. The other 22 features are the independent variables classified into four main groups, namely personal details, flight details, pre-boarding and onboard.

The dataset contained 20 missing values “NA”, all of which were found in Arrival delay in minutes. The features in the pre-boarding and onboard groups are supposed to be ordinal scaled responses from 1 to 5. These features were converted to ordinal factors. However, not every customer responded to every question in the survey and these instances were denoted as “0”. The instances that had missing values, invalid survey response or both were removed entirely. In total, 823 rows were dropped from the dataset which corresponds to 7.98% of the original data set. This is acceptable as less than 10% of the dataset was removed.

2.2 Visual Data Exploration

The bar graphs shown in Figure 1 displays the response of the customers to the survey questions. The 14 questions touch on various aspects of the pre-boarding and onboard customers’ flight experience. Service categories such as inflight service, onboard service, leg room service, baggage handling, inflight entertainment, seat comfort and food and drinks predominantly received ratings of 4 and 5. Other factors that were good but could be improved further are cleanliness, check-in service, online boarding, gate location and departure/arrival time convenient. The areas where the airlines fall short and needs to improve is in ease of online booking and inflight WIFI service as these received many average or poor ratings.

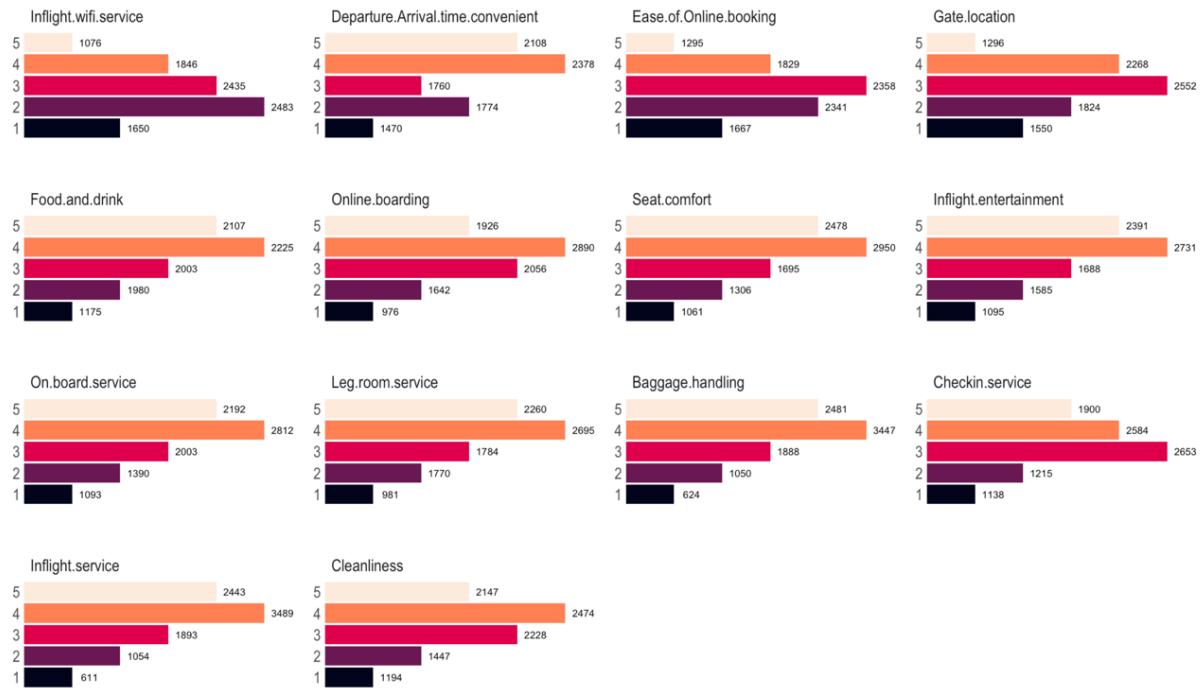


Figure 1: Survey Response

The target class reveals the airline has more neutral/dissatisfied customers than satisfied customers. The gender category is roughly equally distributed and therefore has no significant impact on satisfaction.

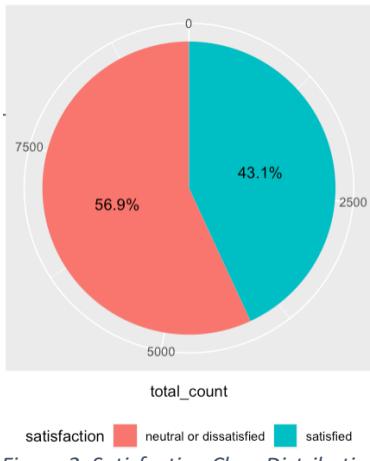


Figure 2: Satisfaction Class Distribution

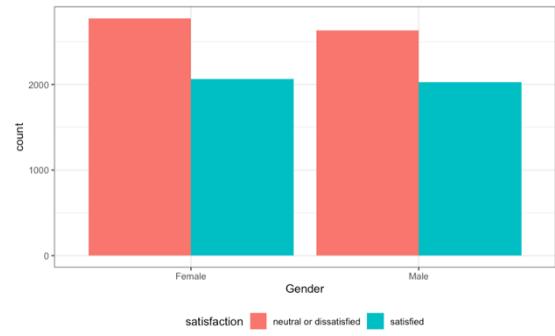


Figure 3: Satisfaction Distribution over the Population based on Gender

The following figures represent some of the features that are bad predictors of satisfaction. In the case of Figure 4 regardless of the rating there are more neutral/dissatisfied than satisfied customers. Figure 5 and 6 is similar in the sense that at high ratings the difference between satisfied and neutral/dissatisfied customers is only marginal and can be regarded as roughly the same.

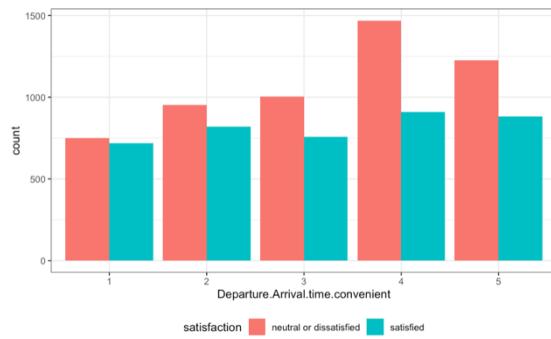


Figure 4: Satisfaction Distribution over the Population based on Departure/Arrival Time Convenience

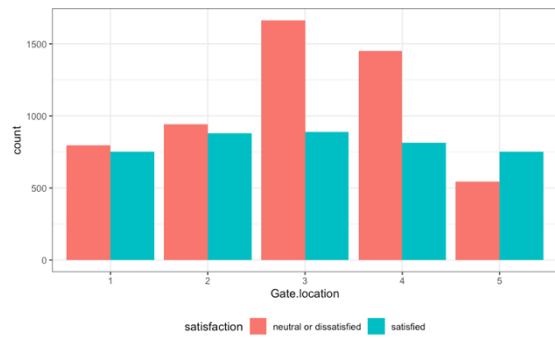


Figure 5: Satisfaction Distribution over the Population based on Gate Location

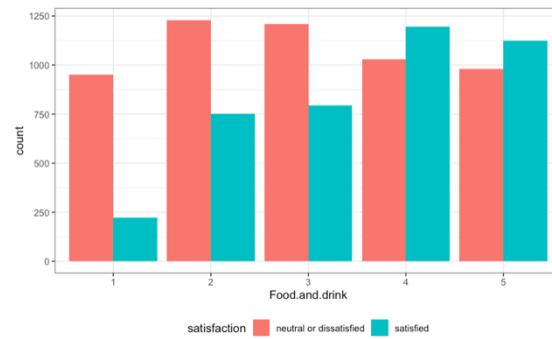


Figure 6: Satisfaction Distribution over the Population based on Food and Drinks

Customers that travel business class are generally more satisfied than those that travel using other classes. Eco Plus class proves to improve satisfaction compared to Eco class although very few customers opt for this class. Figures 8 to 10 represent some features that are useful in explaining customers' satisfaction. A high rating shows significantly more satisfied

customers. Inflight Entertainment, Seat Comfort and On-board Service exhibit similar distributions to that shown in Figure 10.

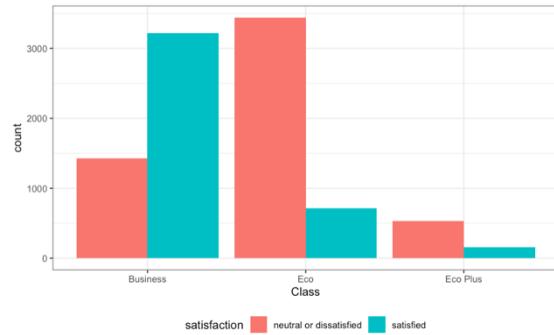


Figure 7: Satisfaction Distribution over the Population based on Class

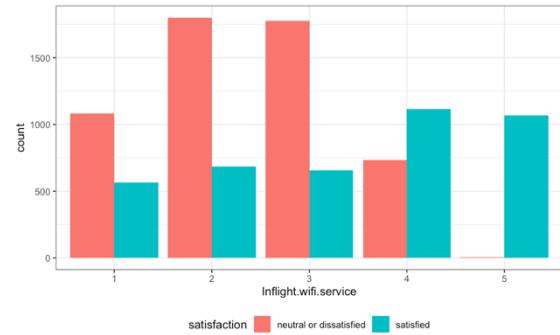


Figure 8: Satisfaction Distribution over the Population based on Inflight WIFI Service

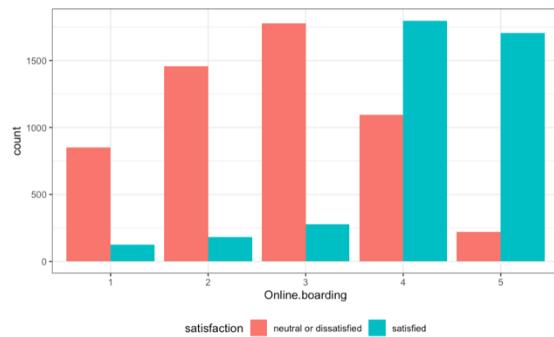


Figure 9: Satisfaction Distribution over the Population based on Online Boarding

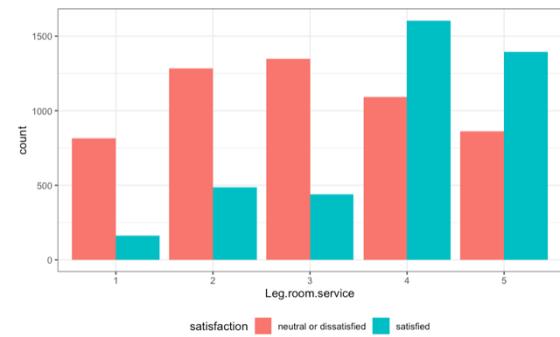


Figure 10: Satisfaction Distribution over the Population based on Leg Room Service

From Figure 11 people younger than 38 years and older than 61 years are relatively more dissatisfied/neutral while people aged between 38 to 61 tend to be more satisfied. Majority of customers are middle aged while there are very few customers past the age of 70. Figure 12 depicts that more customers feel neutral/dissatisfied for short distance flights. This is seen to be less pronounced past 1400 miles. Customers that travel long distance are generally more satisfied. Majority of the flights taken range between short to medium distance, with there being only a few flights that travel past 4000 miles.

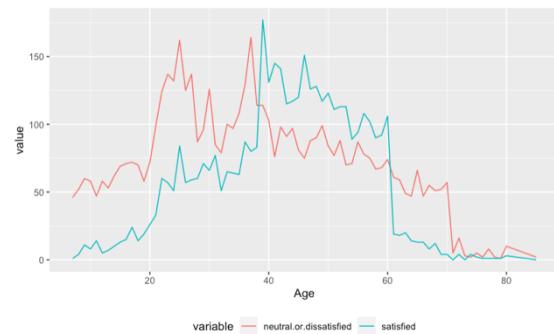


Figure 11: Satisfaction Distribution over the Population based on Age

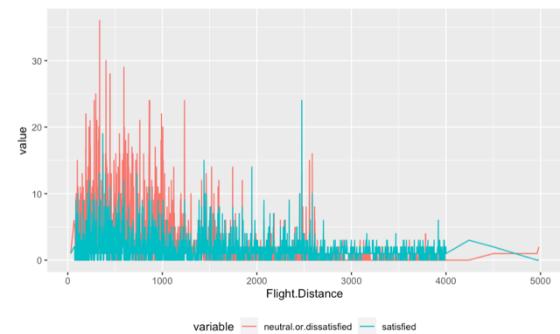


Figure 12: Satisfaction Distribution over the Population based on Flight Distance

Figure 13 shows that almost all customers travelling for business travel opted for business class. Customers travelling for personal reasons mainly chose either eco or eco plus. Figure 14 depicts that personal travel customers make shorter distance trips. Contrastingly, flights taken for business purposes tend to be longer distance travel.

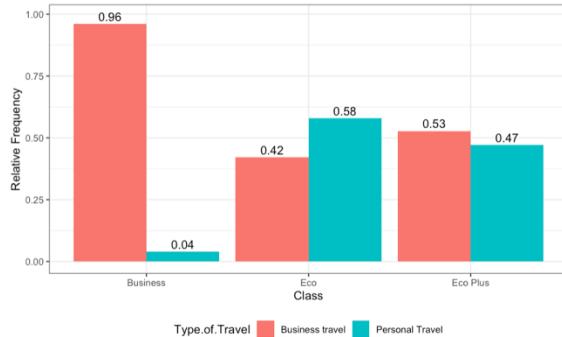


Figure 13: Conditional Probability of Class given Type of Travel

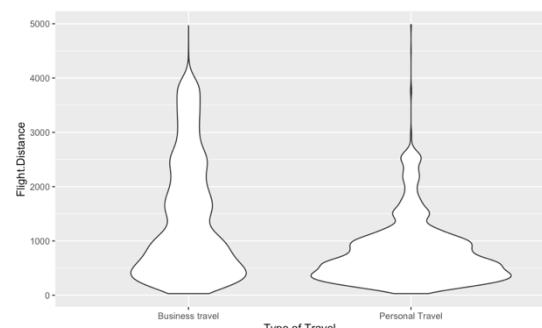


Figure 14: Distribution of Flight Distance for Different Types of Travel

3 Feature Selection

3.1 Information Value (IV)

A higher IV values indicate stronger relationship between the feature and the target variable, and hence a stronger predictive power. Based on the IV results in Figure 15, the variables that offer the most discriminatory power in order is online boarding, inflight Wi-Fi service, class, type of travel and inflight entertainment. The other variables that also have relatively high IV is seat comfort, leg room service, on board service and cleanliness.

The variables that offer the least in terms of discriminatory power are gender, departure arrival time convenient, departure delay in minutes, arrival delay in minutes and gate location. These variables had an IV lower than 0.2 and could potentially be dropped.

	Variable	IV
12	Online.boarding	2.0876096507
7	Inflight.wifi.service	1.7353250714
5	Class	1.2243093501
4	Type.of.Travel	1.1103434418
14	Inflight.entertainment	0.9989876700
13	Seat.comfort	0.7729616705
16	Leg.room.service	0.6125061232
15	On.board.service	0.5794272166
20	Cleanliness	0.5286069724
17	Baggage.handling	0.4479620931
19	Inflight.service	0.4267773205
6	Flight.Distance	0.3953947092
9	Ease.of.Online.booking	0.3845272768
18	Checkin.service	0.3124117310
3	Age	0.2845647828
2	Customer.Type	0.2416255449
11	Food.and.drink	0.2381024379
10	Gate.location	0.1187230914
22	Arrival.Delay.in.Minutes	0.0427471527
21	Departure.Delay.in.Minutes	0.0270911254
8	Departure.Arrival.time.convenient	0.0222156670
1	Gender	0.0003020206

Figure 15: IV Results

3.2 Multicollinearity

From the correlation heatmap in Figure 16, Ease of Online Boarding and Inflight WIFI Service, Cleanliness and Food and Drinks, and Arrival Delay in Minutes and Departure Delay in Minutes are correlated. The latter has the highest correlation of 96% and is visualised in Figure 17. Referring to the correlation coefficients found using the `assoc()` function, one variable from each correlated pair could potentially be dropped. The aforementioned pairs all had a correlation of 60% or higher.

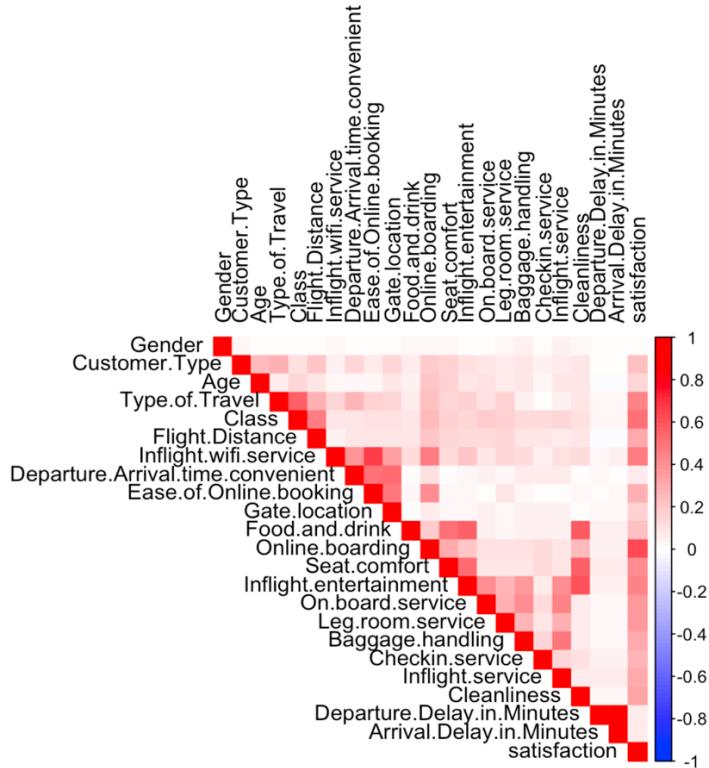


Figure 16: Correlation Heatmap

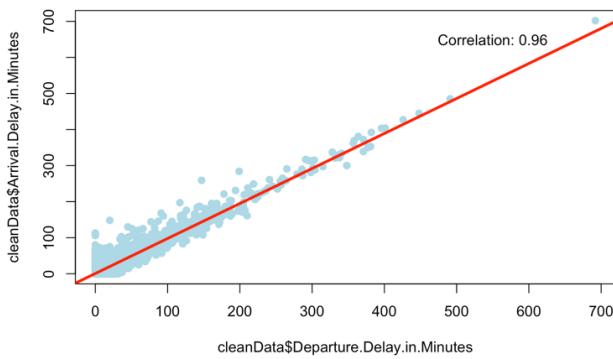


Figure 17: Correlation of Arrival Delay and Departure Delay

Generalised variance inflation factor (GVIF) was also used to identify multicollinearity and substantiate the results found from the correlation heatmap. Dummy variables were created due to the presence of categorical features. GVIF was performed using a binomial logistic regression model. A high GVIF score infers that the independent variable has high correlation

with other independent variables. Arrival Delay in Minutes, Departure Delay in Minutes and Inflight WIFI Service had a relatively high GVIF score as seen in Figure 18.

cleanData.Age	cleanData.Flight.Distance
1.178986	1.355370
cleanData.Departure.Delay.in.Minutes	cleanData.Arrival.Delay.in.Minutes
13.564936	13.616862
GenderMale	Customer.TypeLoyal.Customer
1.046430	1.834342
Type.of.TravelPersonal.Travel	ClassEco
2.591003	1.866224
ClassEco.Plus	Inflight.wifi.service.L
1.254982	7.045629
Inflight.wifi.service.Q	Inflight.wifi.service.C
4.717366	4.614777
Inflight.wifi.service.4	Departure.Arrival.time.convenient.L
2.162503	2.628719
Departure.Arrival.time.convenient.Q	Departure.Arrival.time.convenient.C
2.438997	2.492397
Departure.Arrival.time.convenient.4	Ease.of.Online.booking.L
2.349862	2.942564
Ease.of.Online.booking.Q	Ease.of.Online.booking.C
2.595982	2.783291
Ease.of.Online.booking.4	Gate.location.L
2.443321	2.360329
Gate.location.Q	Gate.location.C
2.356557	2.125942
Gate.location.4	Food.and.drink.L
1.934436	1.820090
Food.and.drink.Q	Food.and.drink.C
1.760712	1.871274
Food.and.drink.4	Online.boarding.L
1.870208	1.789492
Online.boarding.Q	Online.boarding.C
1.364842	1.840143
Online.boarding.4	Seat.comfort.L
1.486487	2.149764
Seat.comfort.Q	Seat.comfort.C
1.876090	2.103872
Seat.comfort.4	Inflight.entertainment.L
1.726498	3.375473
Inflight.entertainment.Q	Inflight.entertainment.C
2.989353	3.687170
Inflight.entertainment.4	On.board.service.L
3.446937	1.911951
On.board.service.Q	On.board.service.C
1.550335	1.905068
On.board.service.4	Leg.room.service.L
1.757058	1.410678
Leg.room.service.Q	Leg.room.service.C
1.303540	1.485563
Leg.room.service.4	Baggage.handling.L
1.240710	2.295813
Baggage.handling.Q	Baggage.handling.C
1.949676	2.139844
Baggage.handling.4	Checkin.service.L
1.627738	1.309420
Checkin.service.Q	Checkin.service.C
1.103777	1.126621
Checkin.service.4	Inflight.service.L
1.145648	2.334371
Inflight.service.Q	Inflight.service.C
2.030903	2.422257
Inflight.service.4	Cleanliness.L
1.982723	2.250795
Cleanliness.Q	Cleanliness.C
1.839840	2.064026
Cleanliness.4	
2.032478	

Figure 18: VIF Results

3.3 Regularisation

Only Lasso and Elastic Net regularization are valid options to identify features that could be eliminated, since Ridge does not eliminate any features. The data set was first split into training set (80%) and test set (20%). Cross validation was done to identify the combination of the best regularization parameter (α) between 0 and 1 and the optimal penalty parameter (λ) that minimizes the mean squared error (MSE) in the training set. Figure 19 shows that an Elastic Net binomial regularization model with $\alpha=0.6$ should be fit. Figure 20 below shows the minimum penalty parameter (λ_{\min}) and the penalty parameter with 1 standard error (λ_{1SE}). The selected penalty is λ_{1SE} since the MSE is lowest when tested on the test data using this penalty.

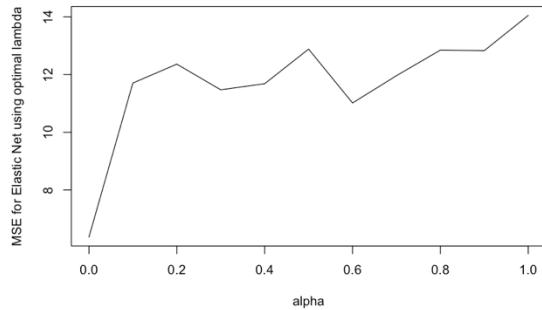


Figure 19: Effect of α on MSE of Regularisation Model

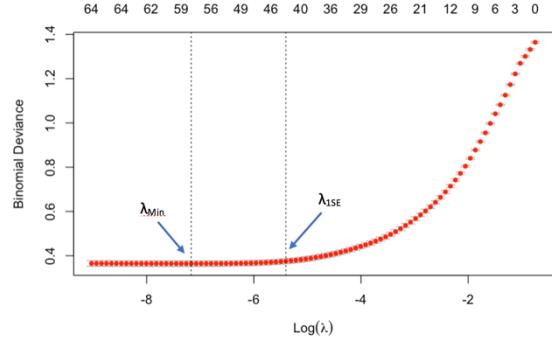


Figure 20: Cross validation result of Effect of λ on Deviance

The coefficients show that Age, Departure Delay in Minutes, Gender and Food and Drinks can be shrunk down to zero and hence infers can be removed.

	s1	
(Intercept)	-1.430735e+00	Seat.comfort.Q
cleanData.Age	.	Seat.comfort.C
cleanData.Flight.Distance	1.393191e-05	Seat.comfort.4
cleanData.Departure.Delay.in.Minutes	.	Inflight.entertainment.L
cleanData.Arrival.Delay.in.Minutes	-3.693295e-03	Inflight.entertainment.Q
GenderMale	.	Inflight.entertainment.C
Customer.TypeLoyal.Customer	2.634696e+00	Inflight.entertainment.4
Type.of.TravelPersonal.Travel	-3.018750e+00	On.board.service.L
ClassEco	-8.742660e-01	On.board.service.Q
ClassEco.Plus	-8.769470e-01	On.board.service.C
Inflight.wifi.service.L	3.072282e+00	On.board.service.4
Inflight.wifi.service.Q	1.993979e+00	Leg.room.service.L
Inflight.wifi.service.C	1.562650e-01	Leg.room.service.Q
Inflight.wifi.service.4	.	Leg.room.service.C
Departure.Arrival.time.convenient.L	-8.310276e-01	Leg.room.service.4
Departure.Arrival.time.convenient.Q	-2.065457e-01	Baggage.handling.L
Departure.Arrival.time.convenient.C	6.448128e-02	Baggage.handling.Q
Departure.Arrival.time.convenient.4	2.291564e-02	Baggage.handling.C
Ease.of.Online.booking.L	6.264063e-01	Baggage.handling.4
Ease.of.Online.booking.Q	.	Checkin.service.L
Ease.of.Online.booking.C	.	Checkin.service.Q
Ease.of.Online.booking.4	.	Checkin.service.C
Gate.location.L	-4.725227e-01	Checkin.service.4
Gate.location.Q	-1.255652e-01	Inflight.service.L
Gate.location.C	9.861708e-02	Inflight.service.Q
Gate.location.4	.	Inflight.service.C
Food.and.drink.L	.	Inflight.service.4
Food.and.drink.Q	.	Cleanliness.L
Food.and.drink.C	.	Cleanliness.Q
Food.and.drink.4	.	Cleanliness.C
Online.boarding.L	1.970021e+00	Cleanliness.4
Online.boarding.Q	1.135313e+00	.
Online.boarding.C	.	.
Online.boarding.4	-5.297474e-01	.
Seat.comfort.L	.	.

Figure 21: Elastic Net Coefficients

The table shown below provides a summary of the results from the techniques above and more importantly the features that were dropped before performing MDS. Considering the limitations of each techniques, some amount of judgement was used in deciding variables to drop and retain. Departure delay in minutes and gender were dropped since majority of the techniques indicate they are not useful. Since departure delay in minutes was dropped, arrival delay in minutes was retained. Gate location was dropped as this is out of the control of the airlines. Based on the visual exploration customer type was removed as it was similar to class and type of travel. Whilst Elastic Net recommended food and drink to be dropped, this is regarded as a factor that influence customers' satisfaction based on logical reasoning and hence is retained.

Independent Variables	Visual Data Exploration	IV	Correlation Heatmap	VIF	Regularisation	MDS
Gender	Bad Predictor	Low	Low	Low	Zero	Dropped
Customer Type	Bad Predictor	Moderate	Low	Low	Non-zero	Dropped
Age	Average Predictor	Moderate	Low	Low	Zero	Retained
Type of Travel	Good Predictor	High	Low	Low	Non-zero	Retained
Class	Good Predictor	High	Low	Low	Non-zero	Retained
Flight Distance	Good Predictor	Moderate	Low	Low	Non-zero	Retained
Inflight WIFI Service	Good Predictor	High	High	High	Non-zero	Retained
Departure/Arrival time convenient	Bad Predictor	Low	Low	Low	Non-zero	Dropped
Ease of Online booking	Average Predictor	Moderate	High	Low	Non-zero	Retained
Gate location	Bad Predictor	Low	Low	Low	Non-zero	Dropped
Food and drink	Bad Predictor	Moderate	High	Low	Zero	Retained
Online boarding	Good Predictor	High	Low	Low	Non-zero	Retained
Seat comfort	Good Predictor	High	Low	Low	Non-zero	Retained
Inflight entertainment	Good Predictor	High	Low	Low	Non-zero	Retained
On-board service	Good Predictor	High	Low	Low	Non-zero	Retained
Leg room service	Good Predictor	High	Low	Low	Non-zero	Retained
Baggage handling	Average Predictor	Moderate	Low	Low	Non-zero	Retained
Checkin service	Average Predictor	Moderate	Low	Low	Non-zero	Retained
Inflight service	Average Predictor	Moderate	Low	Low	Non-zero	Retained
Cleanliness	Good Predictor	High	High	Low	Non-zero	Retained
Departure Delay in Minutes	Bad Predictor	Low	High	High	Zero	Dropped
Arrival Delay in Minutes	Bad Predictor	Low	High	High	Non-zero	Retained

Figure 22: Summary of EDA and Feature Selection

4 Multidimensional Scaling (MDS)

The “gower” coefficient was used to calculate the dissimilarity matrix since the data contains a mix of numerical and categorical features. To reduce the computational time and memory required, repeated random subsampling is done to identify the optimal number of dimensions (k). The most frequently occurring k was selected as the optimal k which in this case is 2 as seen in Figure 23.

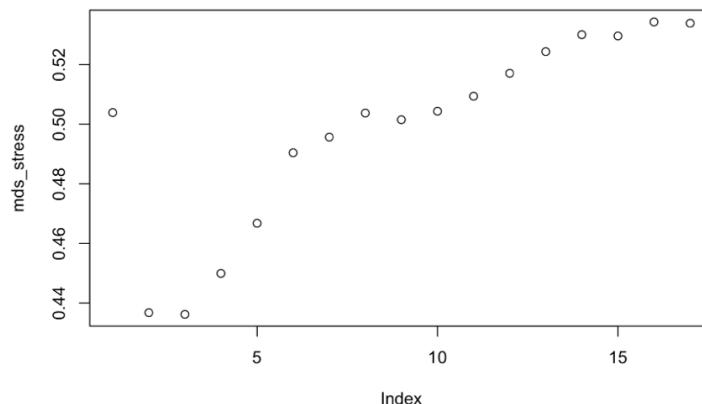


Figure 23: Stress Plot to Identify "k"

The MDS is plotted using k=2 for the entire dataset in Figure 24. The dimensional reduction done by MDS found two new dimensions D1 and D2 that are combinations of the original features that was able to find three distinct clusters. Filling in the colour of each point based on the satisfaction illustrates that majority of the satisfied customers correspond to the lowest of the three clusters when D1 is positive and D2 is close to zero or slightly negative.

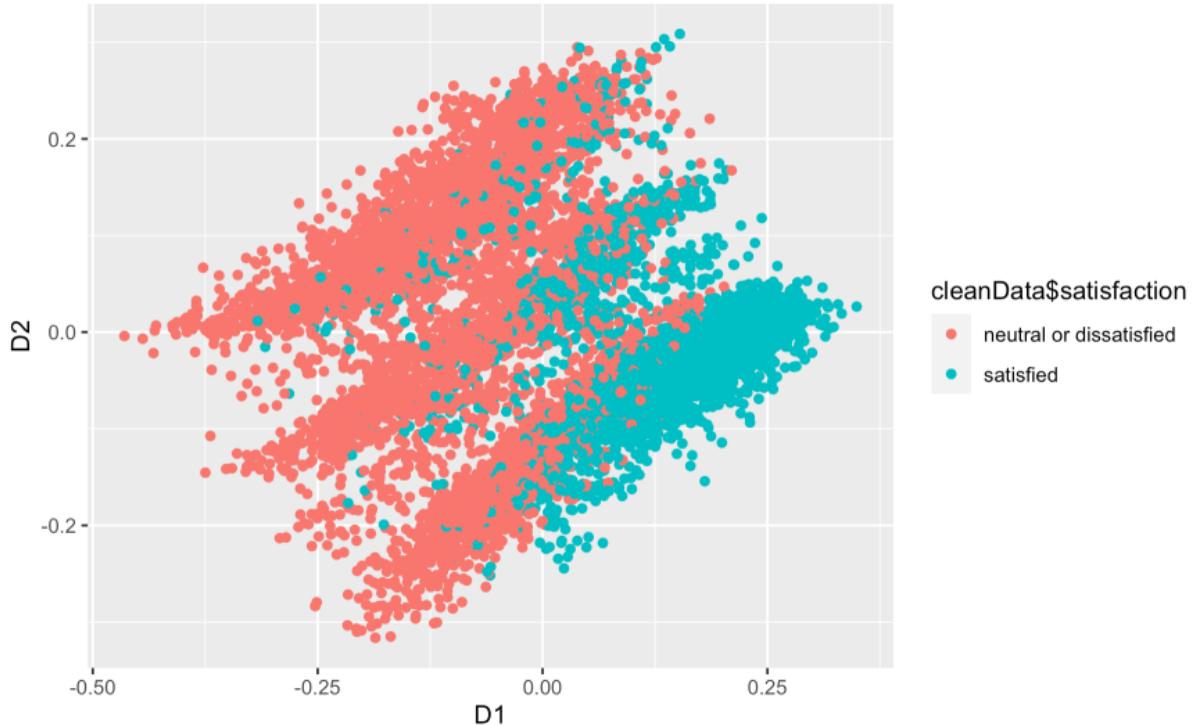


Figure 24: MDS with Satisfaction Colour Scheme Overlaid

Stepwise regression using the AICc criterion was used to identify the features that best explain D1 and D2. The coefficients of the model with the lowest AIC value were examined for each new dimension. This is shown in Figures 25 and 26.

```
Call: glm(formula = D1 ~ Inflight.entertainment.L + ClassEco + Online.boarding.L +
  Type.of.TravelPersonal.Travel + Checkin.service.L + Leg.room.service.L +
  Seat.comfort.L + On.board.service.L + Inflight.wifi.service.L +
  Food.and.drink.L + Inflight.service.L + Cleanliness.L + ClassEco.Plus +
  Baggage.handling.L + Ease.of.Online.booking.L + cleanData.Flight.Distance +
  Inflight.wifi.service.Q + Cleanliness.C + Inflight.service.Q +
  Inflight.wifi.service.C + cleanliness.Q + Food.and.drink.C +
  Baggage.handling.Q + Food.and.drink.Q + Gate.location.Q +
  Inflight.service.C + Seat.comfort.C + Ease.of.Online.booking.C +
  Departure.Arrival.time.convenient.Q + Departure.Arrival.time.convenient.4 +
  Customer.TypeLoyal.Customer + Baggage.handling.4 + On.board.service.Q +
  Checkin.service.Q + Cleanliness.Q + Leg.room.service.Q +
  On.board.service.C + Cleanliness.4 + GenderMale + Departure.Arrival.time.convenient.C +
  Gate.location.4 + Online.boarding.Q + Online.boarding.C +
  Leg.room.service.C + Inflight.service.4 + Baggage.handling.C +
  Departure.Arrival.time.convenient.L + Gate.location.L + Inflight.wifi.service.4 +
  On.board.service.4 + Inflight.entertainment.4 + Food.and.drink.4,
  data = mds_d1)
```

Figure 25: Original Features that Explain D1

```
Call: glm(formula = D2 ~ Type.of.TravelPersonal.Travel + Inflight.entertainment.L +
  ClassEco + Cleanliness.L + Food.and.drink.L + ClassEco.Plus +
  On.board.service.L + Seat.comfort.L + Inflight.service.L +
  CleanData.Flight.Distance + Checkin.service.L + Baggage.handling.L +
  Online.boarding.L + Inflight.entertainment.C + Inflight.entertainment.Q +
  Inflight.wifi.service.Q + Ease.of.Online.booking.L + Inflight.wifi.service.L +
  Seat.comfort.Q + Seat.comfort.C + Food.and.drink.C + Departure.Arrival.time.convenient.Q +
  Leg.room.service.L + Cleanliness.Q + Cleanliness.C + Baggage.handling.Q +
  cleanliness.Q + Ease.of.Online.booking.Q + Baggage.handling.C +
  Customer.TypeLoyal.Customer + Food.and.drink.Q + Checkin.service.Q +
  Checkin.service.C + Ease.of.Online.booking.Q + Gate.location.Q +
  Inflight.service.Q + Inflight.service.C + Leg.room.service.Q +
  Departure.Arrival.time.convenient.C + Ease.of.Online.booking.C +
  GenderMale + On.board.service.C + Seat.comfort.4 + Inflight.wifi.service.4 +
  Inflight.wifi.service.C + Online.boarding.C + Ease.of.Online.booking.4 +
  Online.boarding.4 + Leg.room.service.4 + Inflight.entertainment.4 +
  On.board.service.Q + Leg.room.service.C + Departure.Arrival.time.convenient.L,
  data = mds_d2)
```

Figure 26: Original Features that Explain D2

The association of D1 and D2 with the original features were also examined. The findings from both approaches show that D1 is highly correlated to Inflight Entertainment, Seat Comfort, Class, Online Boarding, Cleanliness, On-board Service, Food and Drink, and Type of Travel. D2 is highly correlated to type of travel and class. It can be inferred that D1 is related to the customers' in-flight experience and D2 is related to the class type.

The following plots shown below are all the same MDS plot shown in Figure 24 with the colour scheme changed to represent the various original features in the dataset. Since they are all the same plot, the regions that represent Satisfaction is the same.

Many customers that travel for business is satisfied whilst relatively few customers that travel for personal reasons are satisfied. Customers that travel for business are more often than not satisfied. Eco plus class does a slightly better job at keeping customers satisfied than Eco class.

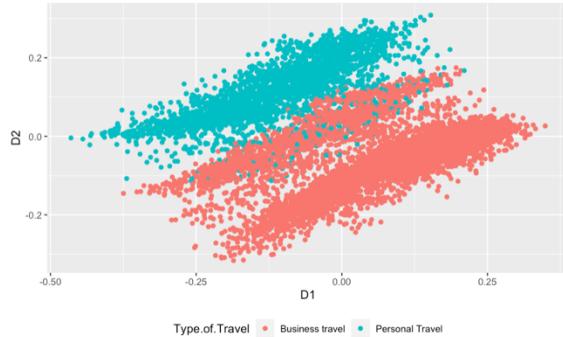


Figure 27: MDS - Type of Travel

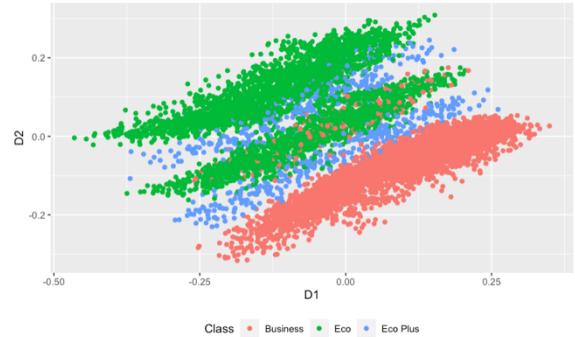


Figure 28: MDS - Class

Customers that travel further distances are found to be more satisfied with their overall experience compared to those that travel shorter distances.

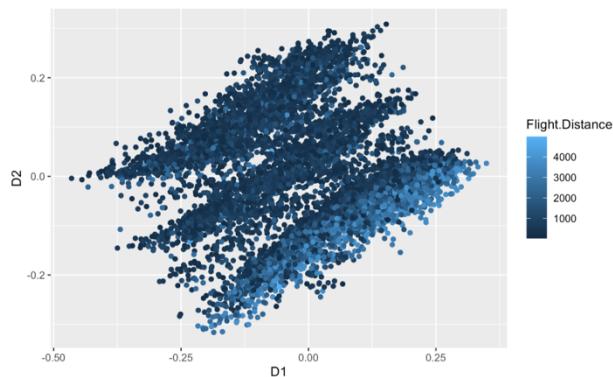


Figure 29: MDS - Flight Distance

Figures 30 to 35 all show similar trends where majority of the high ratings for each service category correspond to the regions that represent satisfied customers. Meanwhile, when the ratings are either average or poor customers are neutral/dissatisfied. This suggest that these features are the areas that most affect customer satisfaction. Therefore, the airlines should focus on improving these areas.

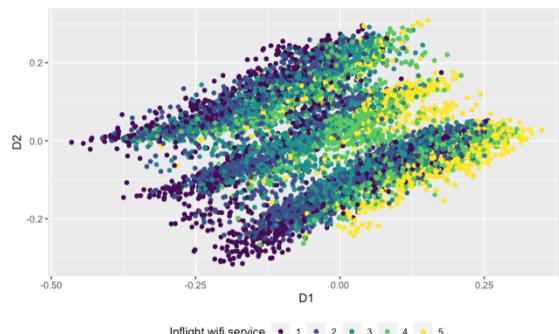


Figure 30: MDS - Inflight WiFi Service

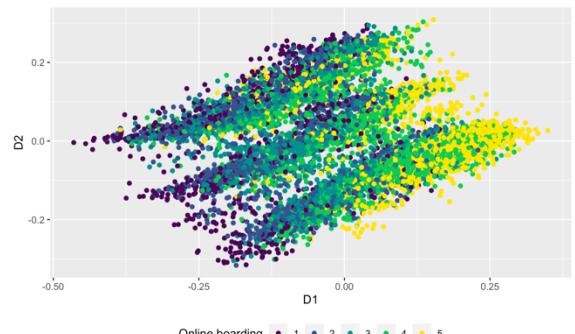


Figure 31: MDS - Online Boarding

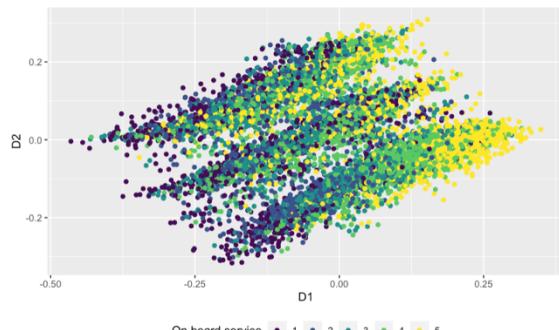


Figure 32: MDS – On-board Service

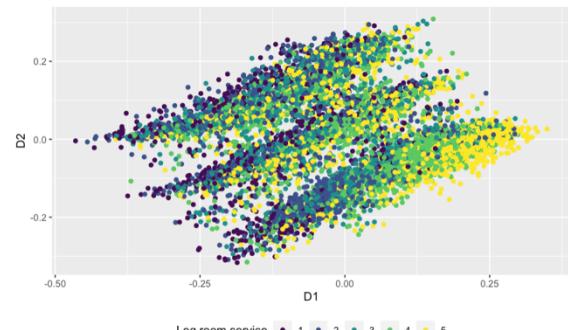


Figure 33: MDS - Leg Room

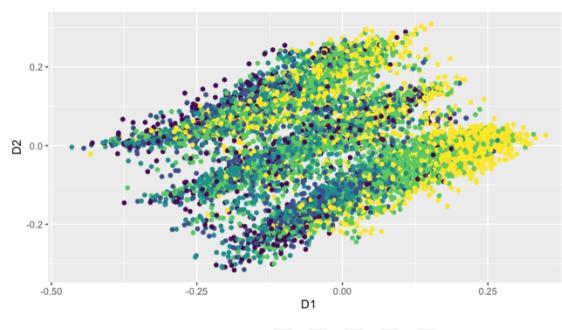


Figure 34: MDS - Baggage Handling

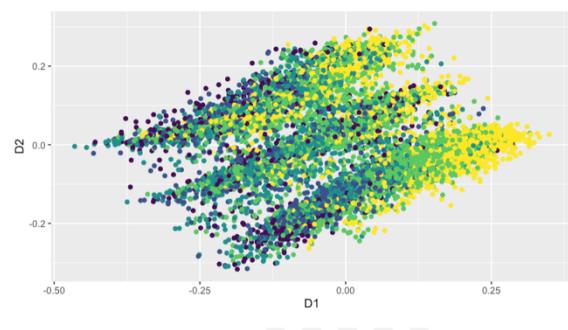


Figure 35: MDS - Inflight Service

Figures 36 to 39 illustrate that satisfied customers found the following areas of service to be good. In relation to neutral/dissatisfied customer, some rated these factors highly while others did not. This highlights that the following features were not the ultimate deciding factors in the satisfaction of customers.

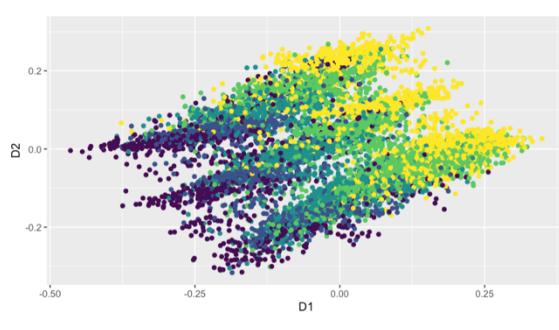


Figure 36: MDS - Seat Comfort

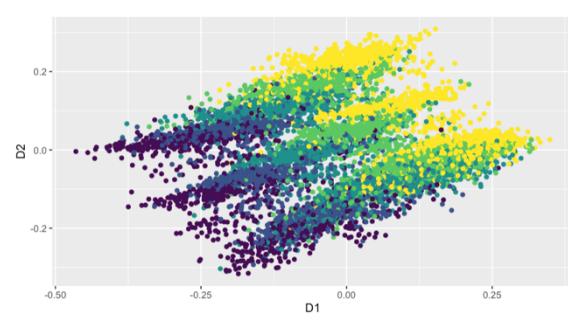


Figure 37: MDS - Cleanliness

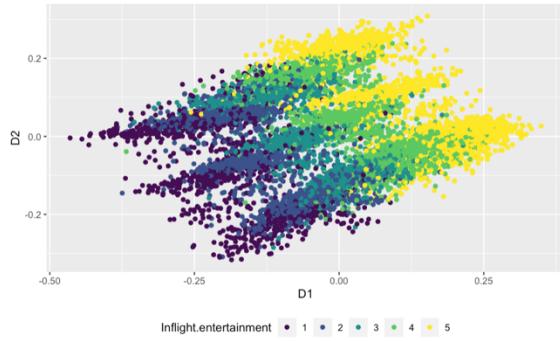


Figure 38: MDS - Inflight Entertainment

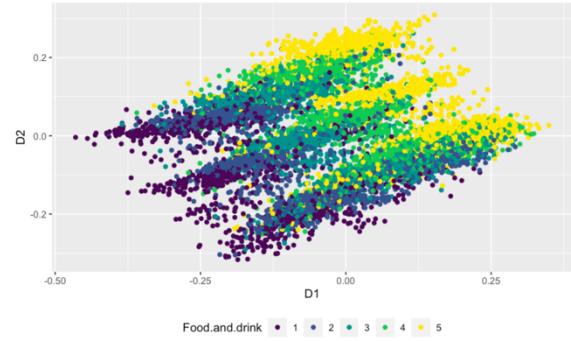


Figure 39: MDS - Food and Drink

5 Conclusion

The study identified inflight WIFI service, on-board service, leg room, baggage handling and inflight service as the features that ultimately decide customers' satisfaction. The airlines should target these factors and improve them further as higher ratings in these categories suggest an increase in satisfaction. The analysis also advocates that the airlines should focus on improving their services for non-business class customers and short distance flights. Based on the survey conducted, the airlines should improve the ease of online booking and inflight WIFI service as customers feel these areas of the airlines service are lackluster.