

DATA MINING

Group 6 - Hotelbook

Gerard Álvarez, Tomás Calaf, Natalia Dai,
Adrià Espinoza and Mario Martín

23/10/2023

List of contents

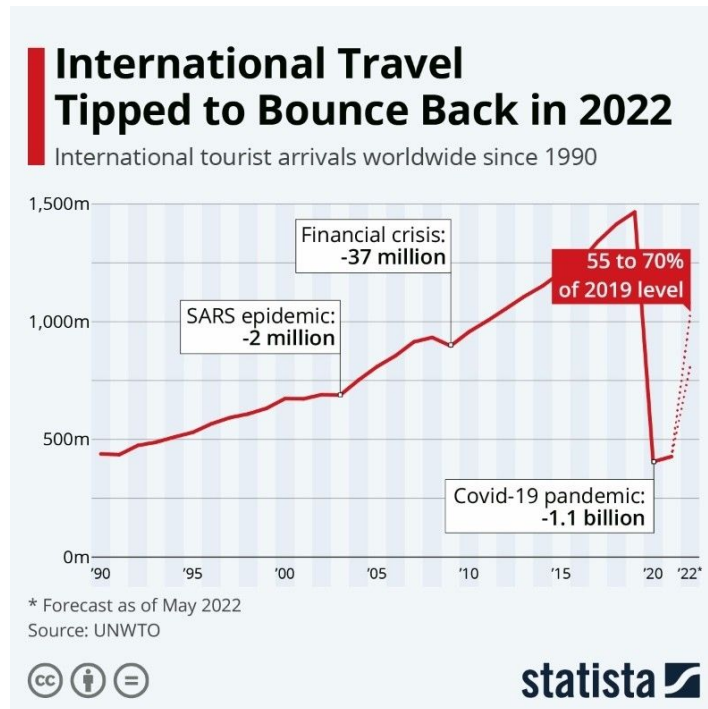
- Introduction
- Data mining process
- Descriptive analysis of variables
- Univariate descriptive analysis
- Preprocessing
- PCA
- Clustering
- Profiling
- PCA vs Clustering
- Conclusions
- Original vs Final scheduling

Introduction

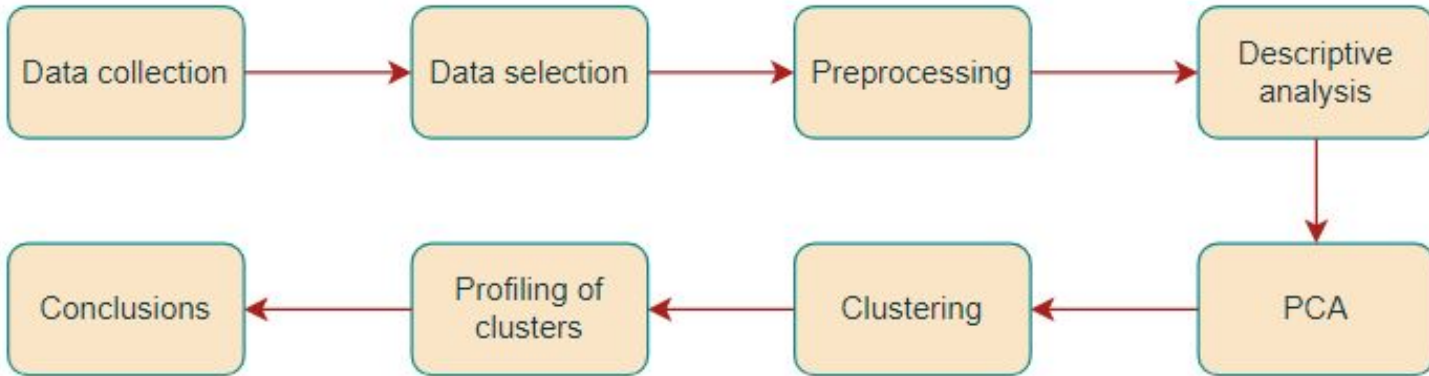
- ➔ Detect the trends among the tourists when they book an hotel
- ➔ Categorize all bookings into groups



https://www.kaggle.com/datasets/jessemostipak/hotel-booking-demand?select=hotel_booking_s.csv



Data mining process

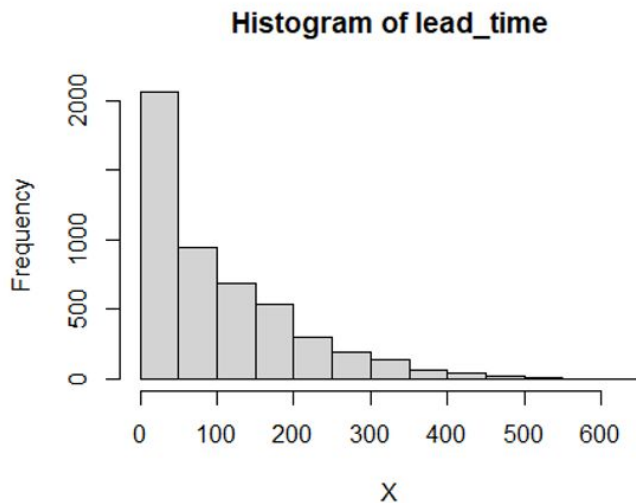


Descriptive analysis of variables

- → hotel
 - is_canceled
 - lead_time
 - arrival_date_month
 - arrival_date_week_number
 - arrival_date_day_of_month
 - stays_in_weekend_nights;
stays_in_week_nights
 - adults; children; babies
 - meal
 - country
 - days_in_waiting_list
 - customer_type
 - adr
 - reservation_status

Univariate descriptive analysis

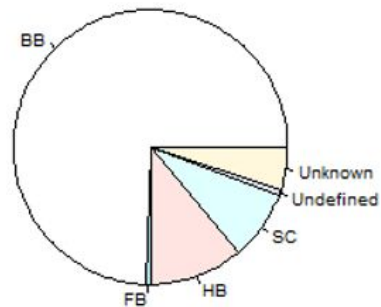
➔ Numerical



Variable	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	sd	vc
lead_time	0.00	18.00	72.63	101.30	155.00	629.00	101.47	1.00

➔ Categorical

Pie of meal

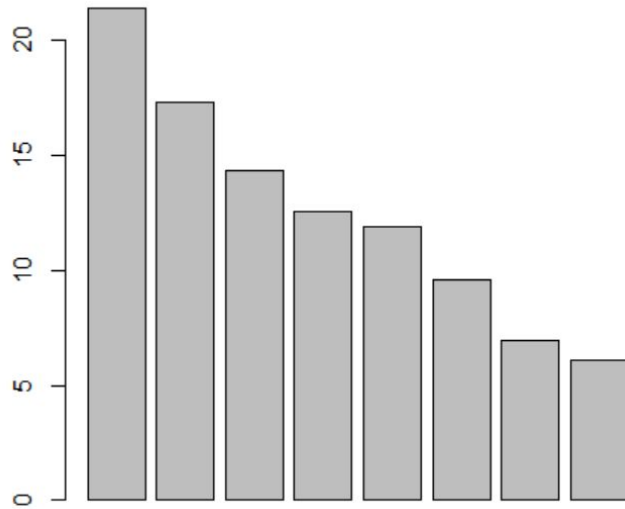


Preprocessing

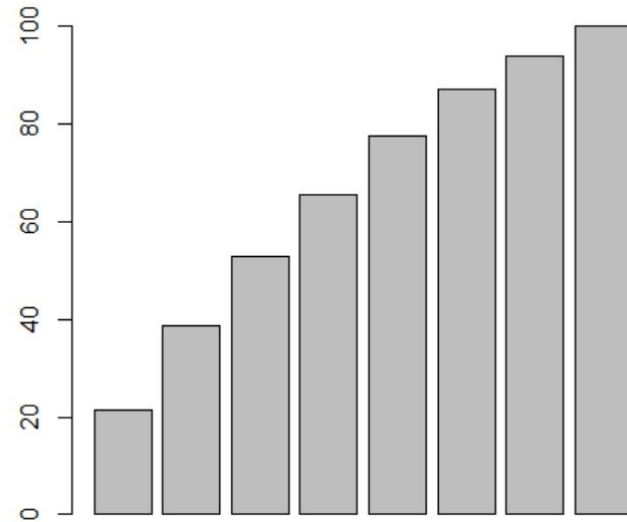
- - Determining working matrix, filtering
 - Creation of the missing values
 - Identification and treatment of missing data
 - Identification and treatment of outliers/errors

PCA

- 80% of information → 6 Dimensions
- PC2 & PC3 → 45° angle

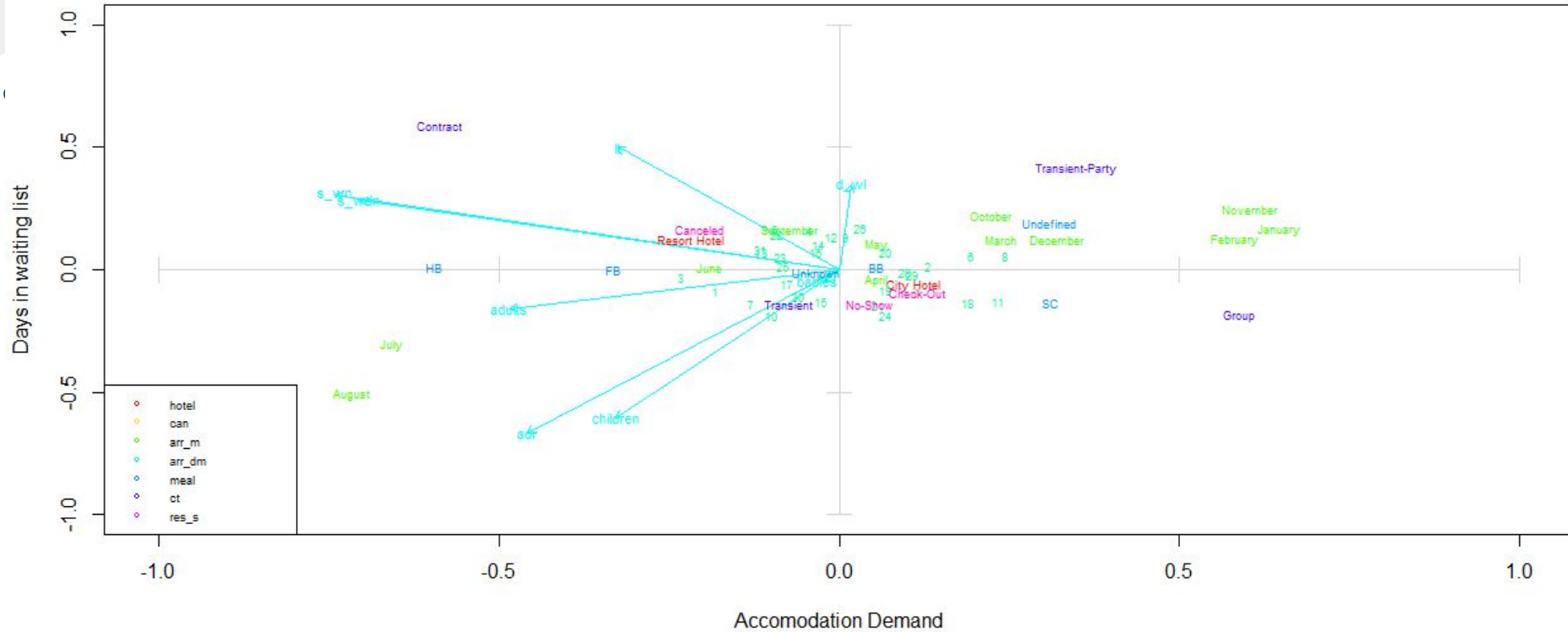


PCx information

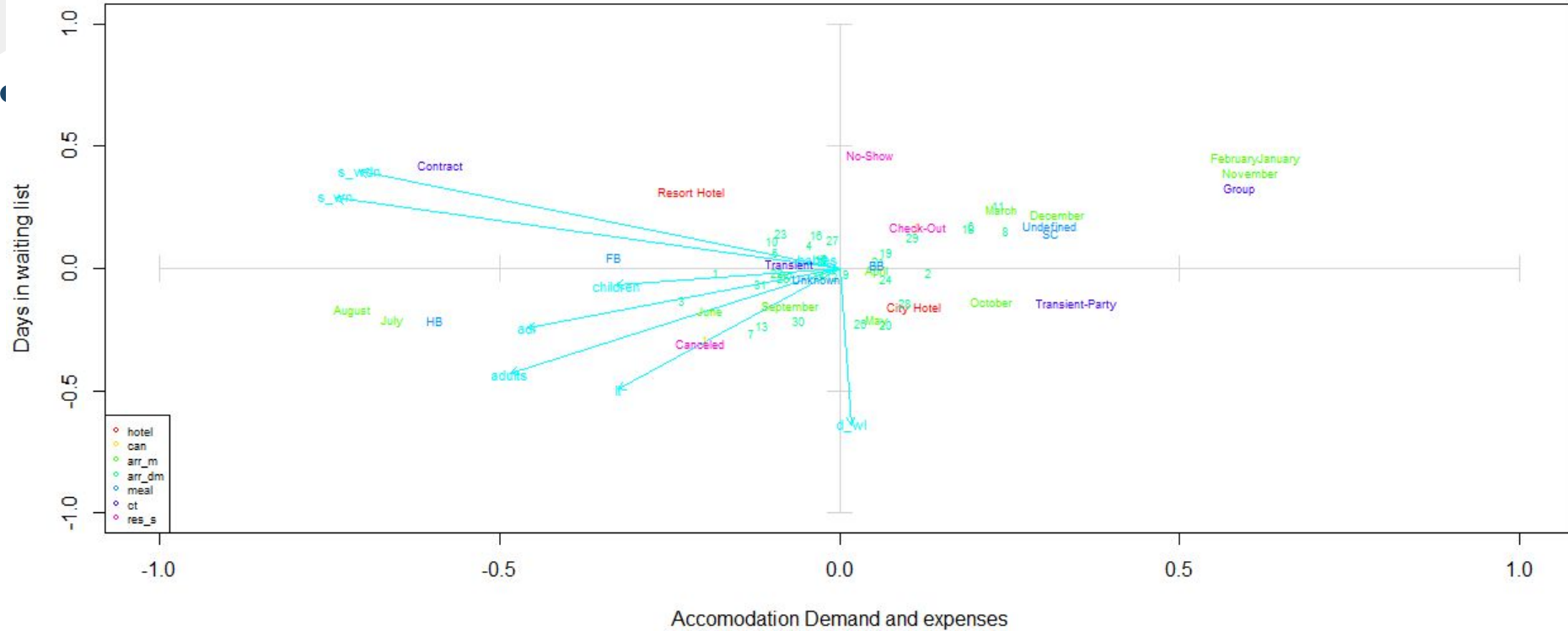


Total information

PCA - PC1 & PC2



PCA - PC1 & PC3

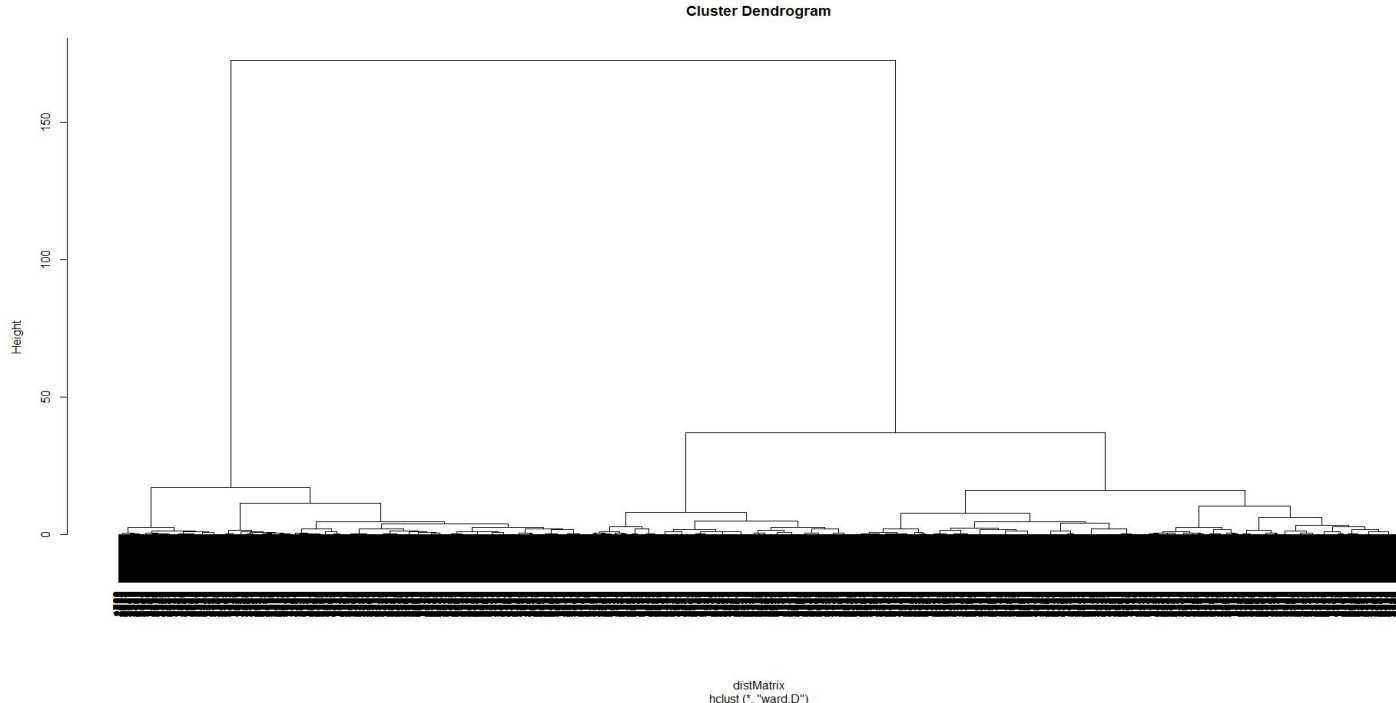


PCA - Conclusions

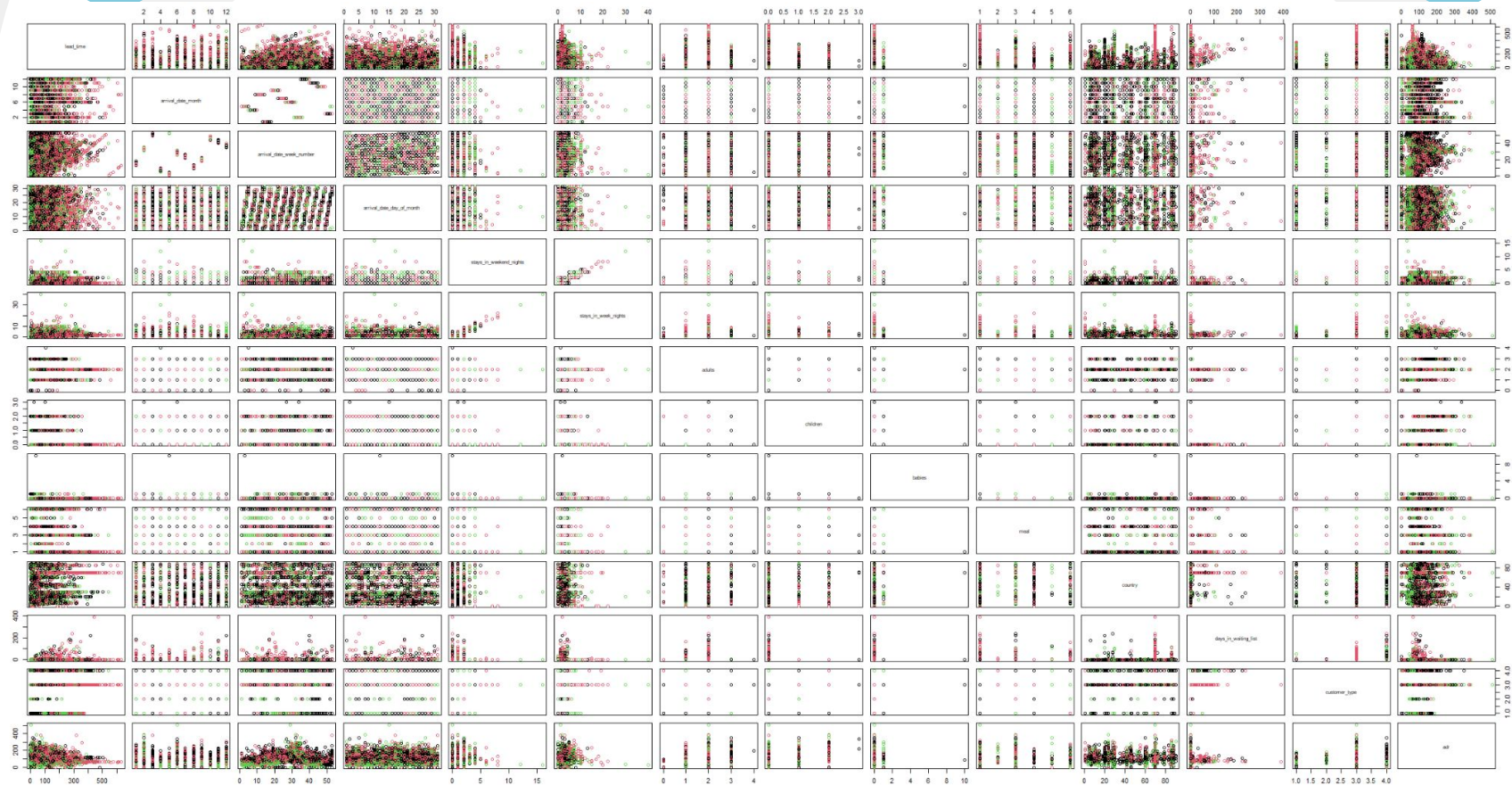
- - Summer months
 - The number of weekday nights and weekend nights and type of customer
 - Contract > transient > transient-party > group

Clustering - Process

- First, we executed the script to obtain the dendrogram.
- We cut the dendrogram and obtained $nc = 3$.
- Overall class graphic for all the variables.

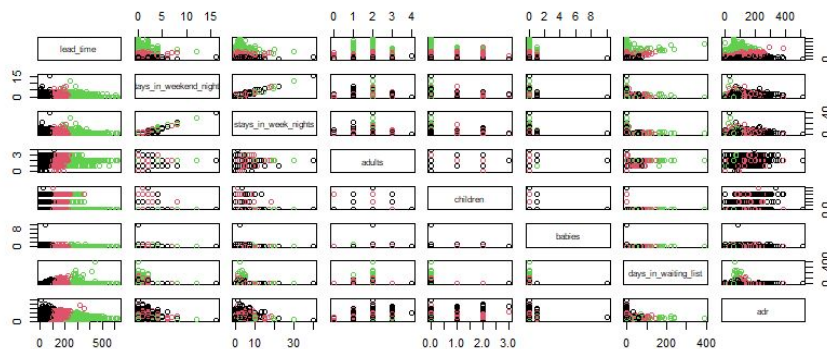


Clustering - Interpretation

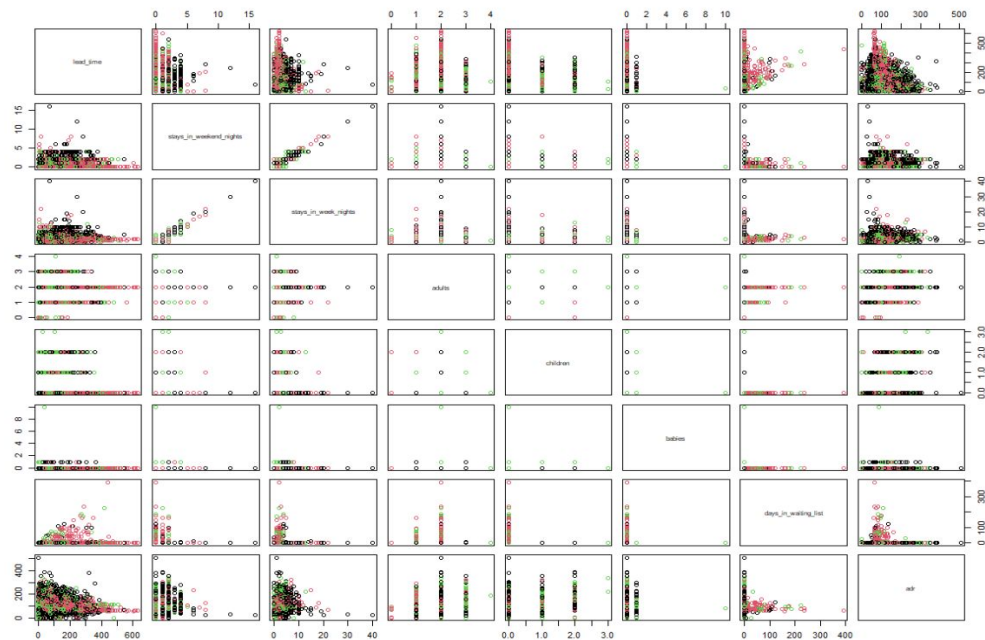


Pairplots using Gower distance

Clustering - Interpretation



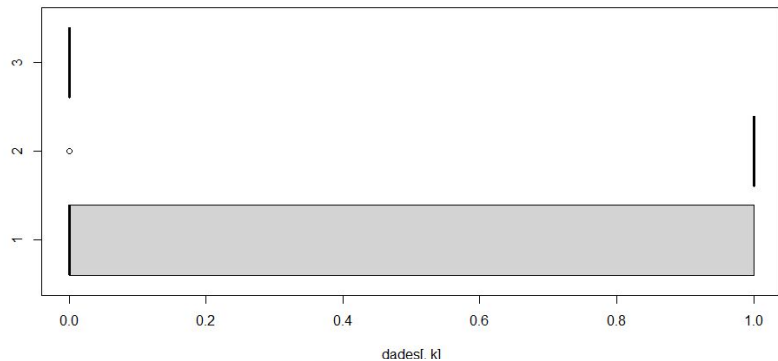
Pairplots using Euclidean distance



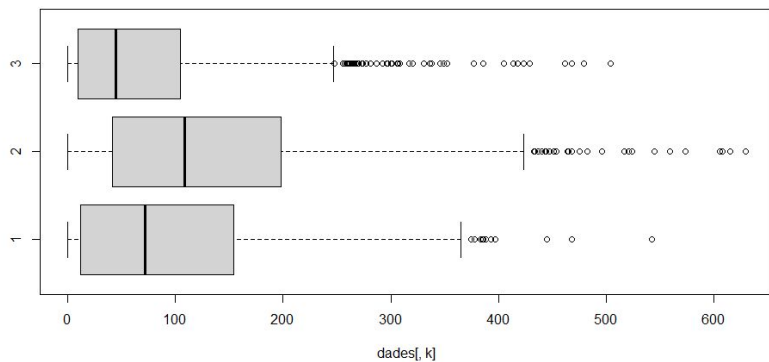
Pairplots using Gower distance
(only numerical values)

Profiling graphs

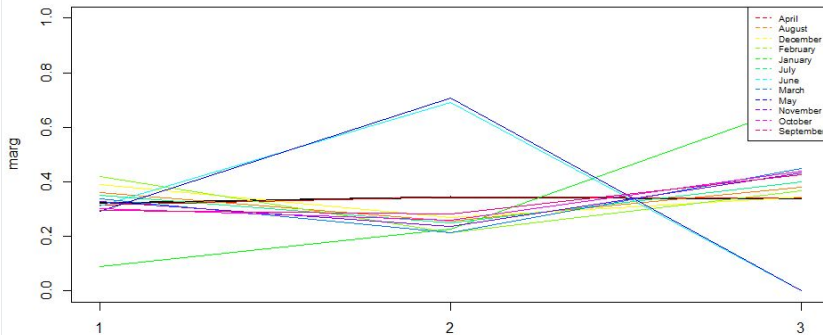
Boxplot of is_canceled vs Class



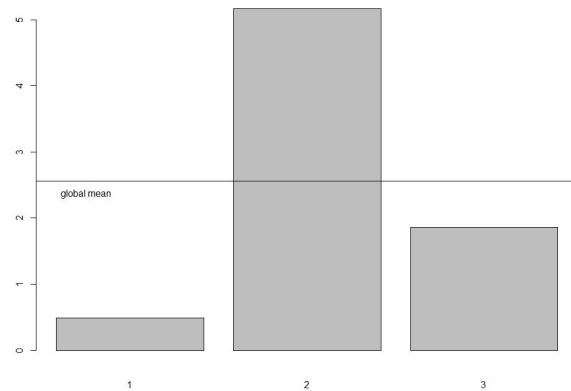
Boxplot of lead_time vs Class



Prop. of pos & neg by arrival_date_month



Means of days_in_waiting_list by Class



Final Class Profiling

- Class 1 and 2 have no noticeable separation.
- Class 2:
 - ◆ high “days_in_waiting_list” and summer values “arrival_date_of_month”
 - ◆ medium high “lead_time” and likely to be cancelled(“is_cancelled”)
- Class 2, with the majority of summer reservations, is the busiest time of the year.
- No further conclusions can be seen with class profiling.

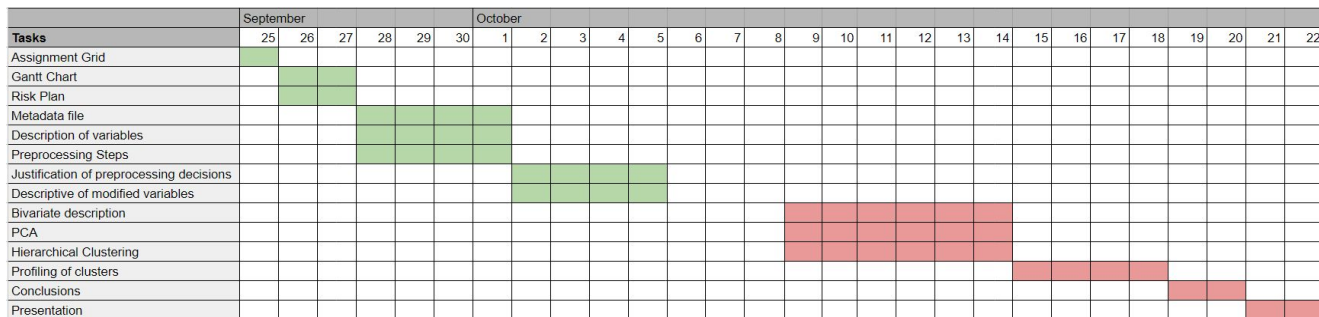
Conclusions: PCA vs Clustering

- PCA: Travel patterns, customer types and seasonal reservations.
- Clustering: Prices fluctuations over the time.
- No clear results, no good data
- PCA better than clustering

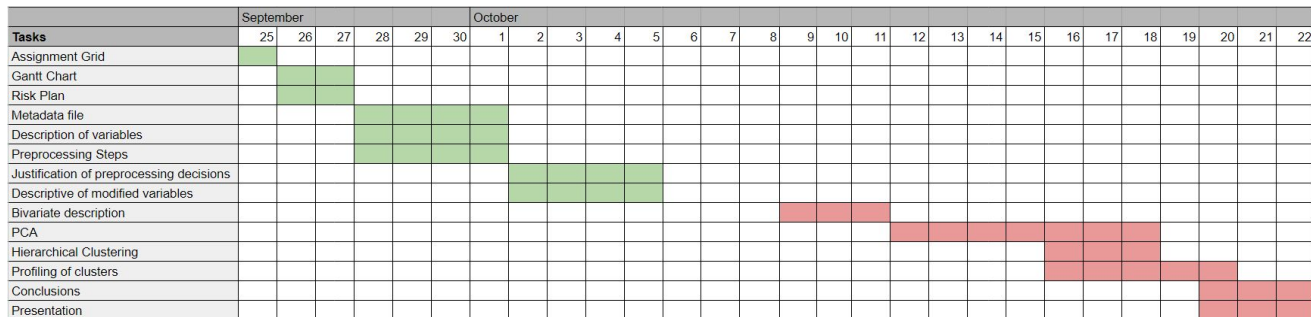
PCA



Original vs Final scheduling



Initial Gantt diagram



Final Gantt diagram

Thank you

Time to answer your
questions