

EDA and Data Preprocessing Overview

Feature Engineering:

- Drop the duplicate rows .
- Drop the rows with column " adult=0 "
- Converting the datatype of **children** and **agent** from [float](#) to [integer](#) .
- Remove column **company**.
- Replace null values with 0 in **agent** feature.

Analyses:

In our model we want to predict the probability of the customer to cancel the booking or not.

I want to find out which factors are most important in canclation by find out the correlated value in according to "is_canceled" to set up a [LogisticRegression](#), so I found out the following features :

total_of_special_requests	-23.465777
required_car_parking_spaces	-19.549782
booking_changes	-14.438099
is_repeated_guest	-8.479342
company	-8.299480
previous_bookings_not_canceled	-5.735772
agent	-4.652945
babies	-3.249109
arrival_date_day_of_month	-0.613008
stays_in_weekend_nights	-0.179108
children	0.503625
arrival_date_week_number	0.814807
arrival_date_year	1.665986
stays_in_week_nights	2.476463
adr	4.755660
days_in_waiting_list	5.418582
adults	6.001721
previous_cancellations	11.013281
lead_time	29.312336
is_canceled	100.000000

