



## Data Article

## Hotel booking demand datasets

Nuno Antonio<sup>a,b,\*</sup>, Ana de Almeida<sup>a,c,d</sup>, Luis Nunes<sup>a,b,d</sup><sup>a</sup> Instituto Universitário de Lisboa (ISCTE-IUL), Lisbon, Portugal<sup>b</sup> Instituto de Telecomunicações, Lisbon, Portugal<sup>c</sup> CISUC, Coimbra, Portugal<sup>d</sup> ISTAR-IUL, Lisbon, Portugal

## ARTICLE INFO

## Article history:

Received 5 October 2018

Accepted 26 November 2018

Available online 29 November 2018

## ABSTRACT

This data article describes two datasets with hotel demand data. One of the hotels (H1) is a resort hotel and the other is a city hotel (H2). Both datasets share the same structure, with 31 variables describing the 40,060 observations of H1 and 79,330 observations of H2. Each observation represents a hotel booking. Both datasets comprehend bookings due to arrive between the 1st of July of 2015 and the 31st of August 2017, including bookings that effectively arrived and bookings that were canceled. Since this is hotel real data, all data elements pertaining hotel or customer identification were deleted. Due to the scarcity of real business data for scientific and educational purposes, these datasets can have an important role for research and education in revenue management, machine learning, or data mining, as well as in other fields.

© 2018 The Authors. Published by Elsevier Inc. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

## Specifications table

Subject area	Hospitality Management
More specific subject area	Revenue Management
Type of data	Text files and R objects
How data was acquired	Extraction from hotels' Property Management System (PMS) SQL databases

\* Corresponding author.

E-mail address: [nuno\\_miguel\\_antonio@iscte-iul.pt](mailto:nuno_miguel_antonio@iscte-iul.pt) (N. Antonio).

Data format	Mixed (raw and preprocessed)
Experimental factors	Some of the variables were engineered from other variables from different database tables. The data point time for each observation was defined as the day prior to each booking's arrival
Experimental features	Data was extracted via TSQL queries executed directly in the hotels' PMS databases and R was employed to perform data analysis
Data source location	Both hotels are located in Portugal: H1 at the resort region of Algarve and H2 at the city of Lisbon
Data accessibility	Data is supplied with the paper

Value of the data

- Descriptive analytics can be employed to further understand patterns, trends, and anomalies in data;
- Used to perform research in different problems like: bookings cancellation prediction, customer segmentation, customer satiation, seasonality, among others;
- Researchers can use the datasets to benchmark bookings' prediction cancellation models against results already known (e.g. [1]);
- Machine learning researchers can use the datasets for benchmarking the performance of different algorithms for solving the same type of problem (classification, segmentation, or other);
- Educators can use the datasets for machine learning classification or segmentation problems;
- Educators can use the datasets to obtain either statistics or data mining training.

1. Data

In tourism and travel related industries, most of the research on Revenue Management demand forecasting and prediction problems employ data from the aviation industry, in the format known as the Passenger Name Record (PNR). This is a format developed by the aviation industry [2]. However, the remaining tourism and travel industries like hospitality, cruising, theme parks, etc., have different requirements and particularities that cannot be fully explored without industry's specific data. Hence, two hotel datasets with demand data are shared to help in overcoming this limitation.

The datasets now made available were collected aiming at the development of prediction models to classify a hotel booking's likelihood to be canceled. Nevertheless, due to the characteristics of the variables included in these datasets, their use goes beyond this cancellation prediction problem.

One of the most important properties in data for prediction models is not to promote leakage of future information [3]. In order to prevent this from happening, the timestamp of the target variable must occur after the input variables' timestamp. Thus, instead of directly extracting variables from the bookings database table, when available, the variables' values were extracted from the bookings change log, with a timestamp relative to the day prior to arrival date (for all the bookings created before their arrival date).



Fig. 1. Diagram of PMS database tables where variables were extracted from.

**Table 1**  
Variables description.

Variable	Type	Description	Source/Engineering
ADR	Numeric	Average Daily Rate as defined by [5]	BO, BL and TR / Calculated by dividing the sum of all lodging transactions by the total number of staying nights
Adults	Integer	Number of adults	BO and BL
Agent	Categorical	ID of the travel agency that made the booking <sup>a</sup>	BO and BL
ArrivalDateDayOfMonth	Integer	Day of the month of the arrival date	BO and BL
ArrivalDateMonth	Categorical	Month of arrival date with 12 categories: “January” to “December”	BO and BL
ArrivalDateWeekNumber	Integer	Week number of the arrival date	BO and BL
ArrivalDateYear	Integer	Year of arrival date	BO and BL
AssignedRoomType	Categorical	Code for the type of room assigned to the booking. Sometimes the assigned room type differs from the reserved room type due to hotel operation reasons (e.g. overbooking) or by customer request. Code is presented instead of designation for anonymity reasons	BO and BL
Babies	Integer	Number of babies	BO and BL
BookingChanges	Integer	Number of changes/amendments made to the booking from the moment the booking was entered on the PMS until the moment of check-in or cancellation	BO and BL/Calculated by adding the number of unique iterations that change some of the booking attributes, namely: persons, arrival date, nights, reserved room type or meal
Children	Integer	Number of children	BO and BL/Sum of both payable and non-payable children
Company	Categorical	ID of the company/entity that made the booking or responsible for paying the booking. ID is presented instead of designation for anonymity reasons	BO and BL
Country	Categorical	Country of origin. Categories are represented in the ISO 3155–3:2013 format [6]	BO, BL and NT
CustomerType	Categorical	Type of booking, assuming one of four categories: Contract – when the booking has an allotment or other type of contract associated to it; Group – when the booking is associated to a group; Transient – when the booking is not part of a group or contract, and is not associated to other transient booking; Transient-party – when the booking is transient, but is associated to at least other transient booking	BO and BL
DaysInWaitingList	Integer	Number of days the booking was in the waiting list before it was confirmed to the customer	BO/Calculated by subtracting the date the booking was confirmed to the customer from the date the booking entered on the PMS
DepositType	Categorical	Indication on if the customer made a deposit to guarantee the booking. This variable can assume three categories: No Deposit – no deposit was made;  Non Refund – a deposit was made in the value of the total stay cost;  Refundable – a deposit was made with a value under the total cost of stay.	BO and TR/Value calculated based on the payments identified for the booking in the transaction (TR) table before the booking’s arrival or cancellation date. In case no payments were found the value is “No Deposit”. If the payment was equal or exceeded the total cost of stay, the value is set as “Non Refund”. Otherwise the value is set as “Refundable”

**Table 1** (continued)

Variable	Type	Description	Source/Engineering
<i>DistributionChannel</i>	Categorical	Booking distribution channel. The term “TA” means “Travel Agents” and “TO” means “Tour Operators”	BO, BL and DC
<i>IsCanceled</i>	Categorical	Value indicating if the booking was canceled (1) or not (0)	BO
<i>IsRepeatedGuest</i>	Categorical	Value indicating if the booking name was from a repeated guest (1) or not (0)	BO, BL and C/ Variable created by verifying if a profile was associated with the booking customer. If so, and if the customer profile creation date was prior to the creation date for the booking on the PMS database it was assumed the booking was from a repeated guest
<i>LeadTime</i>	Integer	Number of days that elapsed between the entering date of the booking into the PMS and the arrival date	BO and BL/ Subtraction of the entering date from the arrival date
<i>MarketSegment</i>	Categorical	Market segment designation. In categories, the term “TA” means “Travel Agents” and “TO” means “Tour Operators”	BO, BL and MS
<i>Meal</i>	Categorical	Type of meal booked. Categories are presented in standard hospitality meal packages: Undefined/SC – no meal package; BB – Bed & Breakfast; HB – Half board (breakfast and one other meal – usually dinner); FB – Full board (breakfast, lunch and dinner)	BO, BL and ML
<i>PreviousBookingsNotCanceled</i>	Integer	Number of previous bookings not cancelled by the customer prior to the current booking	BO and BL / In case there was no customer profile associated with the booking, the value is set to 0. Otherwise, the value is the number of bookings with the same customer profile created before the current booking and not canceled.
<i>PreviousCancellations</i>	Integer	Number of previous bookings that were cancelled by the customer prior to the current booking	BO and BL/ In case there was no customer profile associated with the booking, the value is set to 0. Otherwise, the value is the number of bookings with the same customer profile created before the current booking and canceled.
<i>RequiredCardParkingSpaces</i>	Integer	Number of car parking spaces required by the customer	BO and BL
<i>ReservationStatus</i>	Categorical	Reservation last status, assuming one of three categories: Canceled – booking was canceled by the customer; Check-Out – customer has checked in but already departed; No-Show – customer did not check-in and did inform the hotel of the reason why	BO

**Table 1** (continued)

Variable	Type	Description	Source/Engineering
<i>ReservationStatusDate</i>	Date	Date at which the last status was set. This variable can be used in conjunction with the <i>ReservationStatus</i> to understand when was the booking canceled or when did the customer checked-out of the hotel	BO
<i>ReservedRoomType</i>	Categorical	Code of room type reserved. Code is presented instead of designation for anonymity reasons	BO and BL
<i>StaysInWeekendNights</i>	Integer	Number of weekend nights (Saturday or Sunday) the guest stayed or booked to stay at the hotel	BO and BL/ Calculated by counting the number of weekend nights from the total number of nights
<i>StaysInWeekNights</i>	Integer	Number of week nights (Monday to Friday) the guest stayed or booked to stay at the hotel	BO and BL/Calculated by counting the number of week nights from the total number of nights
<i>TotalOfSpecialRequests</i>	Integer	Number of special requests made by the customer (e.g. twin bed or high floor)	BO and BL/Sum of all special requests

<sup>a</sup> ID is presented instead of designation for anonymity reasons.

Not all variables in these datasets come from the bookings or change log database tables. Some come from other tables, and some are engineered from different variables from different tables. A diagram presenting the PMS database tables from where variables were extracted is presented in Fig. 1. A detailed description of each variable is offered in the following section.

## 2. Experimental design, materials and methods

Data was obtained directly from the hotels' PMS databases' servers by executing a TSQL query on SQL Server Studio Manager, the integrated environment tool for managing Microsoft SQL databases [4]. This query first collected the value or ID (in the case of foreign keys) of each variable in the BO table. The BL table was then checked for any alteration with respect to the day prior to the arrival. If an alteration was found, the value used was the one present in the BL table. For all the variables holding values in related tables (like meals, distribution channels, nationalities or market segments), their related values were retrieved. A detailed description of the extracted variables, their origin, and the engineering procedures employed in its creation is shown in Table 1.

The PMS assured no missing data exists in its database tables. However, in some categorical variables like Agent or Company, "NULL" is presented as one of the categories. This should not be considered a missing value, but rather as "not applicable". For example, if a booking "Agent" is defined as "NULL" it means that the booking did not come from a travel agent.

Summary statistics for both hotels datasets are presented in Tables 2–7. These statistics were obtained using the 'skimr' R package [7].

A word of caution is due for those not so familiar with hotel operations. In hotel industry it is quite common for customers to change their booking's attributes, like the number of persons, staying duration, or room type preferences, either at the time of their check-in or during their stay. It is also common for hotels not to know the correct nationality of the customer until the moment of check-in. Therefore, even though the capture of data took considered a timespan prior to arrival date, it is understandable that the distribution of some variables differ between non canceled and canceled bookings. Consequently, the use of these datasets may require this difference in distribution to be taken into account. This difference can be seen in the table plots of Fig. 2 and Fig. 3. Table plots are a powerful visualization method and were produced with the tabplot R package [8] that allow for the exploration and analysis of large multivariate datasets. In table plots each column represents a variable and each row a bin with a pre-defined number of observations. In these two figures, each bin

**Table 2**

H1 dataset summary statistics – Date variables.

Variable	Min	Max	Median	Unique
<i>ReservationStatusDate</i>	2014-11-18	2017-09-14	2016-07-31	913

**Table 3**

H1 dataset summary statistics – Categorical variables.

Variable	Unique	Top counts
<i>Agent</i>	186	240: 13 095, NULL: 8 209, 250: 2 869, 241: 1 721
<i>ArrivalDateMonth</i>	12	Aug: 4 894, Jul: 4 573, Apr: 3 609, May: 3 559
<i>AssignedRoomType</i>	11	A: 17 046, D: 10 339, E: 5 638, C: 2 214
<i>Company</i>	236	NULL: 36 952, 223: 784, 281: 138, 154: 133
<i>Country</i>	125	PRT: 17 630, GBR: 6 814, ESP: 3 957, IRL: 2 166
<i>CustomerType</i>	4	Tra.: 30 209, Tra.-Party: 7 791, Con.: 1 776, Gro.: 284
<i>DepositType</i>	3	No Dep.: 38 199, Non-Refund.: 1 719, Ref.: 142
<i>DistributionChannel</i>	4	TA/TO: 28 295, Dir.: 7 865, Cor.: 3 269, Und.: 1
<i>IsCanceled</i>	2	0: 28 938, 1: 11 122
<i>IsRepeatedGuest</i>	2	0: 38 282, 1: 1 778
<i>MarketSegment</i>	6	Onl.: 17 729, Off.: 7472, Dir.: 6 513, Gro.: 5 836
<i>Meal</i>	5	BB: 30 005, HB: 8 046, Und.: 1 169, FB: 754
<i>ReservationStatus</i>	3	C.Out: 28 938, Can.: 10 831, No-Show: 291
<i>ReservedRoomType</i>	10	A: 23 399, D: 7 433, E: 4 892, G: 1610

**Table 4**

H1 dataset summary statistics – Integer and numeric variables.

Variable	Mean	SD	P0	P25	Median	P75	P100
<i>ADR</i>	94.95	61.44	-6.38	50	75	125	508
<i>Adults</i>	1.87	0.7	0	2	2	2	55
<i>ArrivalDateOfMonth</i>	15.82	8.88	1	8	16	24	31
<i>ArrivalDateWeekNumber</i>	27.14	14.01	1	16	28	38	53
<i>ArrivalDateYear</i>	2016.12	0.72	2015	2016	2016	2017	2017
<i>Babies</i>	0.014	0.12	0	0	0	0	2
<i>BookingChanges</i>	0.29	0.73	0	0	0	0	17
<i>Children</i>	0.13	0.45	0	0	0	0	10
<i>DaysInWaitingList</i>	0.53	7.43	0	0	0	0	185
<i>LeadTime</i>	92.68	97.29	0	10	57	155	737
<i>PreviousBookingsNotCanceled</i>	0.15	1	0	0	0	0	30
<i>PreviousCancellations</i>	0.1	1.34	0	0	0	0	26
<i>RequiredCarParkingSpaces</i>	0.14	0.35	0	0	0	0	8
<i>StaysInWeekendNights</i>	1.19	1.15	0	0	1	2	19
<i>StaysInWeekNights</i>	3.13	2.46	0	1	3	5	50
<i>TotalOfSpecialRequests</i>	0.62	0.81	0	0	0	1	5

contains 100 observations. The bars in each variable show the mean value for numeric variables or the frequency of each level for categorical variables. Analyzing these figures it is possible to verify that, for both of the hotels, the distribution of variables like *Adults*, *Children*, *StaysInWeekendNights*, *StaysInWeekNights*, *Meal*, *Country* and *AssignedRoomType* is clearly different between non-canceled and canceled bookings.

**Table 5**

H2 dataset summary statistics – Date variables.

Variable	Min	Max	Median	Unique
<i>ReservationStatusDate</i>	2014-10-17	2017-09-07	2016-08-10	864

**Table 6**

H2 dataset summary statistics – Categorical variables.

Variable	Unique	Top counts
<i>Agent</i>	224	9: 31 955, NULL: 8 131, 1: 7 137, 14: 3 640
<i>ArrivalDateMonth</i>	12	Aug: 8 983, May: 8 232, Jul: 8 088, Jun: 7 894
<i>AssignedRoomType</i>	9	A: 57 007, D: 14 983, E: 2 168, F: 2 018
<i>Company</i>	208	NULL: 75 641, 40: 924, 67: 267, 45: 250
<i>Country</i>	166	PRT: 30 960, FRA: 8 804, DEU: 6 084, GBR: 5315
<i>CustomerType</i>	4	Tra.:59 404, Tra.-P.: 17 333, Con.: 2 300, Gro.:293
<i>DepositType</i>	3	No Dep.: 66 442, Non-Refund.: 12 868, Ref.: 20
<i>DistributionChannel</i>	5	TA/TO: 68 945, Dir.: 6 780, Cor.: 3 408, GDS: 193
<i>IsCanceled</i>	2	0: 46 228, 1: 33 102
<i>IsRepeatedGuest</i>	2	0: 77 298, 1: 2 032
<i>MarketSegment</i>	8	Onl.: 38 748, Off: 16 747, Gro.: 13 975, Dir.: 6 093
<i>Meal</i>	4	BB: 62 305, SC: 10 564, HB: 6 417, FB: 44
<i>ReservationStatus</i>	3	C.Out: 46 228, Can.: 32 186, No-Show: 916
<i>ReservedRoomType</i>	8	A: 62 595, D: 11768, F: 1 791, E: 1 553

**Table 7**

H2 dataset summary statistics – Integer and numeric variables.

Variable	Mean	SD	P0	P25	Median	P75	P100
<i>ADR</i>	105.3	43.6	0	79.2	99.9	126	5400
<i>Adults</i>	1.85	0.51	0	2	2	2	4
<i>ArrivalDateOfMonth</i>	15.79	8.73	1	8	16	23	31
<i>ArrivalDateWeekNumber</i>	27.18	13.4	1	17	27	38	53
<i>ArrivalDateYear</i>	2016.17	0.7	2015	2016	2016	2017	2017
<i>Babies</i>	0.0049	0.084	0	0	0	0	10
<i>BookingChanges</i>	0.19	0.61	0	0	0	0	21
<i>Children</i>	0.091	0.37	0	0	0	0	3
<i>DaysInWaitingList</i>	3.23	20.87	0	0	0	0	391
<i>LeadTime</i>	109.74	110.95	0	23	74	163	629
<i>PreviousBookingsNotCanceled</i>	0.13	1.69	0	0	0	0	72
<i>PreviousCancellations</i>	0.08	0.42	0	0	0	0	32
<i>RequiredCarParkingSpaces</i>	0.024	0.15	0	0	0	0	3
<i>StaysInWeekendNights</i>	0.8	0.89	0	0	1	2	16
<i>StaysInWeekNights</i>	2.18	1.46	0	1	2	3	41
<i>TotalOfSpecialRequests</i>	0.55	0.78	0	0	0	1	5





## Acknowledgements

The authors would like to thank the hotels' administration for allowing their data to be shared publicly.

## Transparency document. Supplementary material

Transparency document associated with this article can be found in the online version at <https://doi.org/10.1016/j.dib.2018.11.126>.

## Appendix A. Supporting information

Supplementary data associated with this article can be found in the online version at <https://doi.org/10.1016/j.dib.2018.11.126>.

## References

- [1] N. Antonio, A. Almeida, L. Nunes, Predicting hotel bookings cancellation with a machine learning classification model, in: Proceedings of the 16th IEEE International Conference Machine Learning Application, IEEE, Cancun, Mexico, pp. 1049–1054, doi:10.1109/ICMLA.2017.00-11, 2017.
- [2] International Civil Aviation Organization, Guidelines on Passenger Name Record (PNR) data, (2010). ([https://www.iata.org/iata/passenger-data-toolkit/assets/doc\\_library/04-pnr/New%20Doc%209944%201st%20Edition%20PNR.pdf](https://www.iata.org/iata/passenger-data-toolkit/assets/doc_library/04-pnr/New%20Doc%209944%201st%20Edition%20PNR.pdf)) (accessed 17 February 2016).
- [3] D. Abbott, *Applied Predictive Analytics: Principles and Techniques for the Professional Data Analyst*, Wiley, Indianapolis, IN, USA, 2014.
- [4] Microsoft, SQL Server Management Studio (SSMS), (2017). (<https://docs.microsoft.com/en-us/sql/ssms/sql-server-management-studio-ssms>) (accessed 24 March 2018).
- [5] American Hotel & Lodging Association, Uniform System of Accounts for the Lodging Industry, 11th Revised edition, Educational Institute, New York, 2014.
- [6] International Standards Organization, ISO country codes 3166-3:2013, (<https://www.iso.org/obp/ui/#iso:std:iso:3166:-3:ed-2:v1:en,fr>) (accessed 24 March 2018), 2013.
- [7] A. McNamara, E.A. de la Rubia, H. Zhu, S. Ellis, M. Quinn, skimr: Compact and flexible summaries of data. R package version 1.0.1 (<https://CRAN.R-project.org/package=skimr>), 2018.
- [8] M. Tennekes, E. de Jonge, tabplot: Tableplot, a visualization of large datasets (<https://CRAN.R-project.org/package=tabplot>), 2017.