

# Travel the World: Analyzing and Predicting Booking Behavior using E-mail Travel Receipts

Nemanja Djuric<sup>†</sup>, Mihajlo Grbovic<sup>†</sup>, Vladan Radosavljevic<sup>†</sup>,  
Jaikit Savla<sup>‡</sup>, Varun Bhagwan<sup>‡</sup>, Doug Sharp<sup>‡</sup>

<sup>†</sup>Yahoo Labs, Sunnyvale, CA, USA    <sup>‡</sup>Yahoo Inc., Sunnyvale, CA, USA  
{nemanja, mihajlo, vladan, jaikit, vbhagwan, dsharp}@yahoo-inc.com

## ABSTRACT

Tourism industry has grown tremendously in the previous several decades. Despite its global impact, there still remain a number of open questions related to better understanding of tourists and their habits. In this work we analyze the largest data set of travel receipts considered thus far, and focus on exploring and modeling booking behavior of on-line customers. We extract useful, actionable insights into the booking behavior, and tackle the task of predicting the booking time. The presented results can be directly used to improve booking experience of customers and optimize targeting campaigns of travel operators.

## 1. INTRODUCTION

Tourism and hospitality sectors have witnessed enormous growth in the past century. Originating as a pastime of the wealthy, only very recently technological advances and global economic progress have brought traveling for leisure closer to nearly all layers of the society [5]. The rise from modest beginnings has been astounding, and today tourism accounts for nearly 10% of the world's gross domestic product, employing more than 277 million people worldwide [7]. The upward trend is projected to continue in the future at unabated rates, further highlighting the impact and the importance of the thriving industry.

Given such deep economic footprint, it is highly beneficial to understand customers' behavior in order to satisfy the growing demand. A number of challenging problems has brought attention of researchers, who proposed diverse approaches towards painting a clearer picture of a modern tourist. More specifically, in [3] authors focus on understanding how tourists make their travel plans, while authors of [2] tackle the task of profiling of hotel patrons and learning actionable behavioral patterns using rule mining. In the absence of travel data, some studies attempt to recreate travel information using online photos and other resources [8].

In addition to modeling tourist behavior when they are already on the road, understanding decision process that led to booking of the trip is equally important, as it allows tourist companies to actively influence decisions of a potential customer at its source through better ad targeting [4]. For example, effect of online reviews on booking decisions is explored in [6], while [1] analyzes impact of information

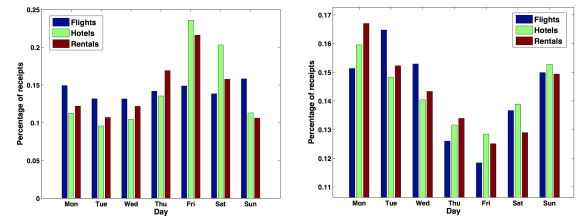


Figure 1: Distribution of a) check-in; and b) booking days

overload. However, such studies mostly rely on surveys comprising few hundred examples, limiting their generality. We address this issue, and explore booking behavior of millions of online customers by mining their travel receipts.

## 2. ANALYSIS OF BOOKING BEHAVIOR

We collected more than 25 million travel receipts, received by a subset of Yahoo Mail users from January through April, 2015. We considered flight, hotel, and car rental receipts, for which we extracted booking timestamp, travel (for flights) and check-in and check-out dates (for hotels and rentals). The data was anonymized (users were assigned random IDs).

First, let us consider weekly distribution of check-in and booking days. In Figure 1a we can see that the largest number of hotel check-ins and rental pick-ups happens on Fridays and Saturdays, while flights had largely uniform distribution. On the other hand, as seen in Figure 1b, the largest number of bookings for hotels and rentals is made on Mondays, and on Tuesdays for flights. The least number of bookings in all cases is made on Fridays. Another interesting finding was a huge increase in hotel check-ins on Valentine's day (not shown), more than 2 times higher number than rates one week before or after. Curiously, we did not see the corresponding increase in flights and rentals, implying that most of these hotel stays were local.

### 2.1 Analysis of individual channels

We take a closer look at the booking channels independently. In Figure 2a we analyze airline booking times segregated by companies. Around 22% of flights are booked within one week ahead of the trip. However, when we take a look at each company individually, this fluctuates from 10% for RyanAir to 35% for US Airways. The results suggests that customers plan further ahead for lower-cost companies which often have special deals and promotions in such cases.

For hotels, we can see that there is a significant difference between short-term (less than 7 days) and long-term stay (7 days and more) in terms of number of days the rooms

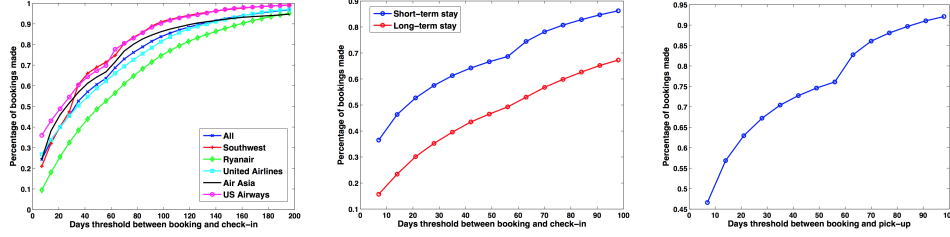


Figure 2: Cumulative distribution of receipts with respect to days booked in advance: a) flights; b) hotels; c) rentals

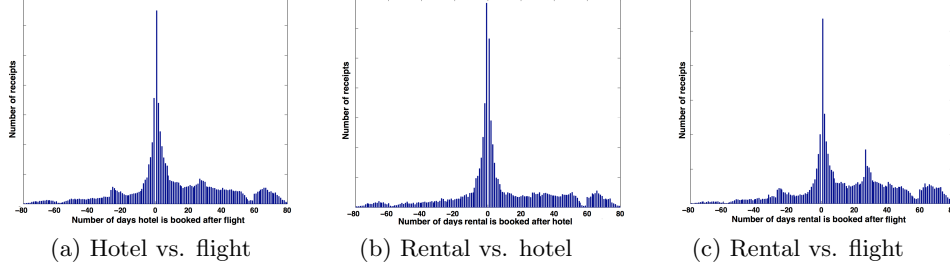


Figure 3: Number of days between coupled bookings from different channels

are booked in advance. For short-term stay nearly 40% of all bookings are made within one week in advance, while for long-term stay the percentage drops to 15%. In other words, customers are planning more in advance for longer, more expensive stays. We also observe jumps in the number of bookings 60 days in advance, which could be explained by black-out windows imposed by some loyalty programs.

Lastly, we see that more than 45% of rentals are booked within 7 days from car pick-up date. This grows to 70% one month prior to pick-up, indicating that rentals are planned much less ahead than either flights or hotels.

## 2.2 Analysis of cross-channel correlation

Next we analyzed time difference between bookings from different channels made for the same trip, where the bookings were deemed part of the same trip if the check-in times are within 24 hours. We only considered trips where the difference between booking dates and check-in date is longer than 30 days, in order to avoid the effect of late purchases which would bias the bookings to be temporally close. The results are given in Figure 3, where  $x$ -axis shows number of days between bookings from two different channels.

By considering the three subfigures, we can see that the majority of two-channel bookings are made within 10 days from each other. Moreover, we can conclude that when users are booking more than one channel for a single trip (e.g., booking a flight in addition to hotel, or rental in addition to flight), they are mostly purchasing them in the order “flight  $\rightarrow$  hotel  $\rightarrow$  rental”, going from higher demand toward smaller demand channels. Interestingly, this conclusion is also confirmed by Figure 2, where 50% of flights are booked around 30 days in advance, 50% of short-term hotel stays around 20 days in advance, while half of car rentals are booked 10 days in advance. These insights are of great importance to booking agencies (e.g., they provide guidance on how and when to target users that already booked a flight).

## 3. PREDICTING BOOKING BEHAVIOR

In this section we turn our attention to behavior prediction task. In particular, we predict when a user will

Table 1: Predicting rental time given flight booking

Method	Rel. accuracy improv.
Most frequent	0%
Age-gender bucket	14%
Logistic regression	23%

book a rental given they already booked a flight. We cast the problem as a classification task and predict labels “<5 days”, “6-10”, “11-20”, “21-40”, and “>41 or never”, while we used age, gender, and past bookings in each channel as features. We trained logistic regression (we split the data into equally-sized training and test sets), and compared to baselines that predict most frequent label and most frequent label in age-gender groups (age buckets were set to “<26”, “26-40”, “>40”). The results are given in Table 1, where we see that the proposed approach outperformed the baselines.

## 4. REFERENCES

- [1] K. Matzler and M. Waiguny. Consequences of customer confusion in online hotel booking. *Information and communication technologies in tourism 2005*, pages 306–317, 2005.
- [2] H. Min, H. Min, and A. Emam. A data mining approach to developing the profiles of hotel customers. *International Journal of Contemporary Hospitality Management*, 14(6):274–285, 2002.
- [3] M. Oppermann. A model of travel itineraries. *Journal of Travel Research*, 33(4):57–61, 1995.
- [4] Å. Rudström and P. Fagerberg. Socially enhanced travel booking: a case study. *Information Technology & Tourism*, 6(3):211–221, 2003.
- [5] E. Sezgin and M. Yolal. *Golden Age of Mass Tourism: Its History and Development*. InTech, 2012.
- [6] B. A. Sparks and V. Browning. The impact of online reviews on hotel booking intentions and perception of trust. *Tourism Management*, 32(6):1310–1323, 2011.
- [7] World Travel & Tourism Council. Economic impact of travel & tourism. Technical report, 2015.

- [8] Y.-T. Zheng, Z.-J. Zha, and T.-S. Chua. Mining travel patterns from geotagged photos. *ACM Trans. on Intelligent Systems and Technology*, 3(3):56, 2012.