

One in a Thousand: Navigating the Noise of Hotel Reviews

Amit Halbreich	20817393	amithalbreich	Amit.Halbreich@mail.huji.ac.il
Omer Mushlion	208271197	omer_mushlion	Omer.Mushlion@mail.huji.ac.il
Omri Marom	319133666	omrimar	Omri.Marom@mail.huji.ac.il

Problem description

In this work we plan to explore the *indicativeness* of online hotel reviews across various aspects (for example: breakfast, service quality, cleanliness) to determine the overall usefulness of these platforms.

The project consists of two main parts: first, we plan to explore and process the reviews of each hotel separately to encode them into a compact and meaningful representation, highlighting contradicting reviews and redundant or conflicting information; Second, we plan to use this summarized representation of hotel reviews to compare between different hotels and try to identify meaningful patterns and insights.

Data

We implemented a JS script to scrape user reviews from booking.com and packaged it into a Google Chrome extension for easy use. When running the script on the reviews page of a hotel on booking.com, it automatically sifts through all reviews pages and scrapes the reviews into a .csv file. The script's code and instructions for running it will be provided within the final project submission. For each hotel, we get a .csv file with the following scheme:

Review Title	Negative Reviews	Positive Reviews	Rating	Stay Date	Review Date	Room Type	# of Nights	Traveler Type	Overall Average Rating
text	text	text	int {1, 10}	date mm + yy	date dd+mm +yy	categorical	int	categorical	float [1, 10]

For example, for the Radisson Blu Tromsø Hotel, Norway we have a .csv file with 1,935 rows for 1,935 English reviews, and this is one of them:

Excellent hotel in	sauna hours could've	great location and	10	Nov 22	08/11/2022	Standard Room	3	Couple	8.3
--------------------	----------------------	--------------------	----	--------	------------	---------------	---	--------	-----

