

Taxi-cancellation-case.csv is the dataset for this case study.

Business Situation

In late 2013, the taxi company Yourcabs.com in Bangalore, India was facing a problem with the drivers using their platform—not all drivers were showing up for their scheduled calls. Drivers would cancel their acceptance of a call, and, if the cancellation did not occur with adequate notice, the customer would be delayed or even left high and dry.

Bangalore is a key tech center in India, and technology was transforming the taxi industry. Yourcabs.com featured an online booking system (though customers could phone in as well), and presented itself as a taxi booking portal. The Uber ride sharing service started its Bangalore operations in mid-2014.

Yourcabs.com had collected data on its bookings from 2011 to 2013, and posted a contest on Kaggle, in coordination with the Indian School of Business, to see what it could learn about the problem of cab cancellations.

The data presented for this case are a randomly selected subset of the original data, with 10,000 rows, one row for each booking. There are 17 input variables, including user (customer) ID, vehicle model, whether the booking was made online or via a mobile app, type of travel, type of booking package, geographic information, and the date and time of the scheduled trip. The target variable of interest is the binary indicator of whether a ride was canceled. The overall cancellation rate is between 7% and 8%.

Assignment

1. How can a predictive model based on these data be used by Yourcabs.com?
2. How can a profiling model (identifying predictors that distinguish canceled/uncanceled trips) be used by Yourcabs.com?
3. Explore, prepare, and transform the data to facilitate predictive modeling. Here are some hints:

- In exploratory modeling, it is useful to move fairly soon to at least an initial model without solving *all* data preparation issues. One example is the GPS information—other geographic information is available so you could defer the challenge of how to interpret/use the GPS information.

⁵Copyright © Datastats, LLC and Galit Shmueli 2017; used with permission.

- How will you deal with missing data, such as cases where NULL is indicated?
 - Think about what useful information might be held within the date and time fields (the booking timestamp and the trip timestamp). The data file contains a worksheet with some hints on how to extract features from the date/time field.
 - Think also about the categorical variables, and how to deal with them. Should we turn them all into dummies? Use only some?
4. Fit several predictive models of your choice. Do they provide information on how the predictor variables relate to cancellations?
 5. Report the predictive performance of your model in terms of error rates (the confusion matrix). How well does the model perform? Can the model be used in practice?
 6. Examine the predictive performance of your model in terms of ranking (lift). How well does the model perform? Can the model be used in practice?

21.6 SEGMENTING CONSUMERS OF BATH SOAP⁶

BathSoap.csv is the dataset for this case study.

Business Situation

CRISA is an Asian market research agency that specializes in tracking consumer purchase behavior in consumer goods (both durable and nondurable). In one major research project, CRISA tracks numerous consumer product categories (e.g., “detergents”), and, within each category, perhaps dozens of brands. To track purchase behavior, CRISA constituted household panels in over 100 cities and towns in India, covering most of the Indian urban market. The households were carefully selected using stratified sampling to ensure a representative sample; a subset of 600 records is analyzed here. The strata were defined on the basis of socioeconomic status and the market (a collection of cities).

CRISA has both transaction data (each row is a transaction) and household data (each row is a household), and for the household data it maintains the following information:

- Demographics of the households (updated annually)
- Possession of durable goods (car, washing machine, etc., updated annually; an “affluence index” is computed from this information)
- Purchase data of product categories and brands (updated monthly)

⁶Copyright © Cytel, Inc. and Resampling Stats, Inc. 2017; used with permission.