

## GROUP (UP TO 3 PEOPLE) WEB ANALYTICS ASSIGNMENT #3 (30 POINTS)

DUE DATE: Wednesday, 6/16/2021

- This assignment will analyze the data (**HotelClickStream.xls**) and interpret the results. This dataset includes clickstream data of online transactions for online hotel booking. Appendix includes the detailed description for the variables.
- **Please follow the instructions very carefully to do this assignment! Please copy/summarize your key results for each question to a word file along with your answers to produce the final report for submission.**

1. Please first create the following 2 additional variables into your data

- 1) **REF\_D** (create a dummy variable indicating whether the transaction was referenced from other website, if not, the final booking website was directly accessed. If no information provided for the variable *REF\_DOMAIN\_NAME*, *REF\_D* = 0; otherwise *REF\_D* = 1)
- 2) **LOG\_PRICE** (take the log transformation of the variable *PROD\_TOTPRICE* using the LOG function in excel)

- a) (5 points) Please provide a summary table showing the top 10 *domain names* (*DOMAIN\_NAME*) that generated the most volume of transactions the report should look like the following Table (Hint: one way to do this is to use the COUNTIF function in excel). Please summarize briefly your observations from the results.

Rank	Domain Names	# of Transactions
------	--------------	-------------------

- b) (5 points) Please provide a summary table showing the top 10 *reference domain names* (*REF\_DOMAIN\_NAME*) that generated the most volume of transactions the report should look like the following Table. Please summarize briefly your observations from the results.

Rank	Reference Domain Names	# of Transactions
------	------------------------	-------------------

2. (6 points) Please use the Binary Outcome (Logistic/Logit) regression technique to answer the question on “*what are the factors that influence people’s decision on whether to book directly on a hotel website or from other third party website?*” Please use *DIRECT\_D* as your Dependent Variable (DV); and *REF\_D*, *LOG\_PRICE*, *TRANS\_FREQ*, *DURATION*, *HOUSEHOLD\_SIZE*, *CHILDREN\_D*, and *CONNECTIONSPEED\_D* as your Independent Variables (IV). Please report and interpret your regression results, which should include the interpretation of the regression coefficients.
3. a) (6 points) Please use the Count Data (Poisson) regression model to answer the question on “*what are the factors that influence people’s booking frequencies?*” Please use *TRANS\_FREQ* as your DV; and *REF\_D*, *LOG\_PRICE*, *PAGES\_VIEWED*, *HOUSEHOLD\_SIZE*, *CHILDREN\_D*, and *CONNECTIONSPEED\_D* as your IVs. Please report and interpret your regression results, which should include the interpretation of the regression coefficients.
- b) (6 points) Please repeat the analysis in question a) using the Negative Binomial Regression model. Please report and interpret your regression results and coefficients.
- c) (2 points) Please summarize your observations by comparing the results from a) and b).

## Appendix

- This is a sample data selected from a large online clickstream dataset. The data was collected by tracking over 100,000 unique household online shopping behavior. This small sample data includes transactions for booking hotels online.

### Variable Descriptions

Variable	Description and Measure
<i>ID</i>	Unique transaction ID
<i>DOMAIN_ID</i>	Unique ID for the web domain
<i>MACHINE_ID</i>	Unique ID for the computer (household) on which the transaction was made
<i>SITE_SESSION_ID</i>	Unique ID for the session in which the transaction was made
<i>TRANS_FREQ</i>	Total number of transactions for the household.
<i>DOMAIN_NAME</i>	The website (domain) name where the transaction was made
<i>DIRECT_D</i>	A dummy variable indicating whether the transaction is incurred directly from a hotel website (1) or other third_party travel website (0).
<i>PROD_NAME</i>	The product (e.g., hotel or packages) purchased by the household
<i>PROD_TOTPRICE</i>	Total price paid for this transaction
<i>REF_DOMAIN_NAME</i>	The referring website (domain) name through which the final purchase website was reached
<i>DURATION</i>	Total time spent at a site (mins)
<i>PAGES_VIEWED</i>	Total pages viewed at a site
<i>HOUSEHOLD_SIZE</i>	Total number of people in the household
<i>CHILDREN_D</i>	A dummy variable indicating whether the household has any children.
<i>CONNECTIONSPEED_D</i>	A dummy variable indicating whether the household has high speed internet connection