

# Responsive Split Questionnaire Survey Design for the Estimation of Tourist Expenditure

Ang Khay Wee, Carol Anne Hargreaves\*

Department of Statistics and Applied Probability, National University of Singapore, Singapore, Singapore

## Abstract

This survey design is motivated by the need to address the consequences such as non-response, high respondent burden and poor data quality, that lengthy surveys are often associated with. Responsive split design is a survey design technique that directly addresses these consequences. The objective of this study was to reduce the respondent burden and to increase the response rate by using the Split Questionnaire Design (SQD). We leveraged on the idea of issuing only relevant surveys to each respondent as this directly reduces the connotation of being probed with irrelevant questions which is regarded as a primary reason for non-response. To achieve this, we developed a responsive design by utilising the prior information that we collected as part of the questionnaire, and then created decision rules for the administration of the relevant micro-surveys. The results that we have attained are promising in terms of the trade-off between precision and responsiveness. The responsive design offered the advantage of being more responsive to the tourist and was able to detect rare expenditures. This suggests that the responsive design provides improvements in issues regarding missing and rare events. Besides meeting the objectives, we also demonstrated that a secondary advantage of responsive split questionnaires also brings about significant cost savings to the survey organisation.

## Keywords

Responsive Split Questionnaire Design, Split Questionnaire Design, Machine Learning, Tourism

Received: May 10, 2021 / Accepted: July 22, 2021 / Published online: July 28, 2021

@ 2021 The Authors. Published by American Institute of Science. This Open Access article is under the CC BY license.

<http://creativecommons.org/licenses/by/4.0/>

## 1. Introduction

Many organisations use surveys to collect important information to generate reports to understand market trends and behaviour or have an overview of their performance standing. Very often these reports could then be analysed to inform the organization on the best course of action to be taken to meet certain goals. Stimulated by increasing demands for more detailed information from respondents, stakeholders of organisation may request for more survey items to be included in the questionnaires. Hence, there is a tendency for these questionnaires to be long.

In this paper, we analyse tourism expenditure in 15 different domains of the industry. They are keen to know what the expenditure amount in industries are ranging from

*accommodations, dining to medical*, to list a few. The objective of this study is to estimate the mean expenditure ( $\widehat{y}_k$ ) by tourists on each of the expenditure categories shown in Table 1 below.

**Table 1.** List of Expenditure Categories.

|                      |                      |
|----------------------|----------------------|
| 1. Hotel             | 9. Sightseeing       |
| 2. Homestay          | 10. Attractions      |
| 3. Hostel            | 11. Entertainment    |
| 4. Service Apartment | 12. Medical Services |
| 5. Hawker            | 13. Shopping         |
| 6. Casual Dining     | 14. Businesses       |
| 7. Fine Dining       | 15. Education        |
| 8. Transport         |                      |

The tourist expenditure survey is lengthy and often associated with high respondent burden which Rolstad [13] defined by the amount of effort required to complete the

\* Corresponding author

E-mail address: carol.hargreaves@nus.edu.sg (C. A. Hargreaves), khayweeang@u.nus.edu (A. K. Wee)

survey [1]. High respondent burden has the consequences of introducing negative impacts in the form of high non-response rate and low quality of data collected [9]. This issue can be attributed to the likely more time required to complete the survey and consequently could increase boredom and fatigue levels among respondent [1]. By the new millennium, scientific and government surveys became more complex and often posed great uncertainty in design parameters and operational features. Survey populations' resistance to survey participation continued to increase. Survey cost structures were becoming even more dependent on decisions being made in the field or data collection centers, often with no evidentiary basis to measure or respond to cost fluctuations. To counteract these issues, Chun [4] introduced the responsive and adaptive design (RAD), a scientific framework driven by cost-quality tradeoff analysis and optimization that enables the most efficient production of high-quality data. Murphy [12] discussed the critical and complex design decisions associated with transitioning an interviewer-administered survey to a self-administered, postal, web/paper survey. Murphy's [12] approach embeds adaptive, responsive, and tailored (ART) design principles and data visualization during a multi-phased data collection operation to project the outcomes of each phase in preparation for subsequent phases. Ali [2] tackle the problem of splitting a long (potentially time consuming) questionnaire into two parts, where each participant only responds to a fraction of the questions, and all respondents obtain a common portion of questions. We propose a method that combines regression models to the two independent samples (questionnaires) in the survey.

## 2. Split Questionnaire Design

Chipperfield [5] used Split Questionnaire Designs to collect only the data that was needed and showed empirically and theoretically a significant reduction in respondent burden with a negligible impact on the accuracy of estimates by not collecting data from respondents who identified as contributing little to the accuracy of estimates.

We propose to employ the use of Split Questionnaire Design (SQD) to tackle the issue of high respondent burden. SQD tackles the problem by only administering different subset(s) of questionnaire to particular respondents. For instance we can split the full questionnaire into 15 subset(s) which we call micro-survey(s), and  $y_{ik}$  represents the expenditure  $k^{\text{th}}$  collected from the micro-survey on the  $k^{\text{th}}$  expenditure where  $k \in [0, 15]$ ,  $i = 1, 2, \dots, N$  represents the respondent  $i$  among the population of size  $N$ . The following equation represents how expenditure by respondent  $i$  from the full questionnaire is split into 15 micro-surveys.

$$y_i = (y_{i,1}, y_{i,2}, \dots, y_{i,15}) \quad (1)$$

From (1) SQD can be achieved by instead of administering all 15 micro-surveys, we can arbitrarily choose a subset of these 15 micro-surveys, effectively only collecting some of the  $y_{i,k}$ 's from each respondent. This can be demonstrated in Table 1 where each column signifies the expenditure from 1 to  $k$  and each row signifies 1 respondent. The shaded region indicates which expenditure  $y_k$  is collected from each respondent.

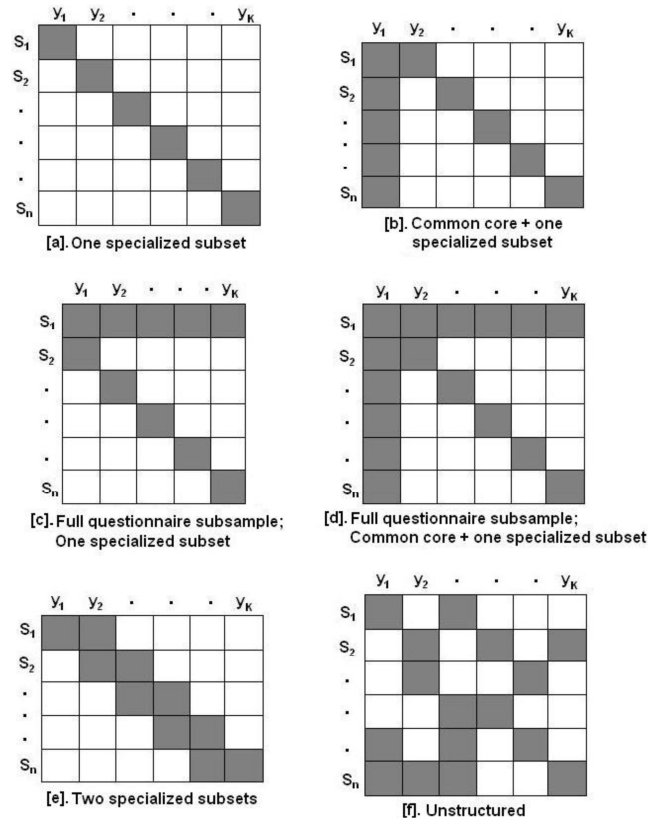


Figure 1. Illustration of the Administration of SQD.

A key concept required in SQD is a *decision rule* to inform on the ways in which we distribute the micro-surveys. With a lack of prior information in SQD, a sensible decision rule that survey designers may employ is using random issuance of each micro-survey with probability of 0.5. This means that each respondent will have a probability of 0.5 of receiving each micro-survey about the  $k^{\text{th}}$  expenditure.

## 3. Responsive Split Questionnaire Design

We further improve on the concept of SQD by utilising the heterogeneous property of spenders, and then we create a *decision rule* on which the micro-survey design is to be administered. The assumption of the heterogeneous property of spenders is a valid assumption as tourists generally do not

share the same expenditure behaviours. We can expect different tourists to have different tendencies to spend on different expenditure categories. Hence, we improve over the SQD by only administering micro-surveys that are relevant to the tourist. The phenomenon of capturing the relevancy of the micro-survey for a particular tourist is what we term, Responsive Split Questionnaire Design (RSQD). Referring to the RSQD in Figure 1, each shaded region signifies the expenditure information most relevant to be collected from a particular respondent. This allows us to achieve targeted and tailored administration of the survey [8, 10, 15].

The procedure of RSQD involves the following components:

1. Designation of a set of Core variables to be used to predict the likely expenditure with a certain probability  $p_{ik}$ .
2. A two phase design where the first phase is dedicated to collecting initial on-boarding Core questions while, the second phase is the actual administration of the relevant micro-surveys.

Gonzalez [9] proposed that the demographic and auxiliary information can be designated as variables for the CORE questions. Where demographic information refers to a set of attributes of the biography information of the respondents whereas the auxiliary information refers to any additional information e.g. any information about the trip. The derivation of a decision rule is obtained using any classification model, for example, a Logistic Regression or Decision Tree model.

In the RSQD, each respondent has a unique set of probabilities assigned to each of the micro-surveys on the  $k^{th}$  expenditure. This  $p_{ik}$  signifies the chance that the respondent  $i$  would incur the  $k^{th}$  expenditure category and hence the probability to be administered the micro-survey.

## 4. Survey Estimates

The RSQD introduces, a reduced sample size for each of the reported  $k$  expenditures. For a full questionnaire, we would have  $N$  reports of expenditure information for each of the  $k$  expenditures, however under the SQD design, we only have a sub sample of the entire population reporting their responses. The consequence of this phenomenon is an increase in the variance of estimation. Equation (2) shows the theoretical sampling variance of the estimated mean of  $\widehat{y}_k$  which is the estimated mean of expenditure of category  $k$ . Where  $\frac{n}{N}$  is known as the sampling fraction of sampling  $n$  respondents out of population of size  $N$  and  $S^2$  is the population variance for  $y_k$  [7].

$$Var(\widehat{y}_k) = (1 - \frac{n}{N}) \frac{S^2}{n} \tag{2}$$

Under SQD, we expect only a subsample of the  $N$

respondents to receive the micro-survey on expenditure  $k$ . For instance, given a new reduced sampling fraction, such as a reduction by half we can compute the increase in variance of our estimate. We let  $Var(\widehat{y}_{SQD})$  represent the variance of our estimate under SQD and  $Var(\widehat{y}_{SPD})$  to represent the variance of our estimate if the full questionnaire were to be administered. The ratio in equation (3) allows us to know the magnitude of increased in variance. Theoretically, a 50% reduction in number of sampled respondents would lead to 2 times increase in variance.

$$\frac{Var(\widehat{y}_{SQD})}{Var(\widehat{y}_{SPD})} = \frac{(1 - \frac{0.5n}{N}) \frac{S^2}{0.5n}}{(1 - \frac{n}{N}) \frac{S^2}{n}} \approx 2 \tag{3}$$

Another implication that should be noted of in employing SQD where the estimation relies heavily on the randomised procedure of administering the micro-survey(s). The randomness is determined by the *decision rule* which is stipulated by the individual  $p_{ik}$ , (i.e  $p_{ik} = 0.5$  for SQD and unique  $p_{ik}$  in RSQD). This is the probability that a respondent  $i$  would be given the micro-survey. Since the random design controls the random behaviour of the sample and would directly influence any estimators that arrives from it, it makes sense to incorporate this design information in the computation of estimator. This is to ensure statistical validity of our estimate [3].

One way to incorporate design information is to employ Horvitz-Thomson estimator [14]. Horvitz Thomson estimator accounts for design information through inverse probability weighting. The Horvitz Thomson is given by equation (4). Where  $\widehat{y}_k$  is the Horvitz Thomson estimator for the mean of expenditure category  $k$ ,  $S_k$  is the subsample of respondents from the entire population who are assigned with micro-survey  $k$ .

$$\widehat{y}_k = (\sum_{i \in S_k} w_{ik})^{-1} (\sum_{i \in S_k} w_{ik} y_{ik}) \tag{4}$$

$$w_{ik} = (p_{ik}^{-1}),$$

An important consideration for the use of Horvitz Thomson estimator is the definition of  $S_k$ . The interpretation of mean of expenditure can be expressed in the following ways:

- 1) Conditional Mean.
- 2) Unconditional Mean.

*Conditional mean* refers to the computation of mean where we only include actual spenders from the all respondents receiving the micro-survey on  $k^{th}$  expenditure whereas *unconditional mean* refers to computation of mean regardless of whether respondents indeed incur expenses. We may refer  $S_{k,cond}$  to be the sub sample of all respondents who received the micro-survey on  $k^{th}$  category and indeed incur the expense whereas  $S_{k,uncond}$  to be the sub sample who are issued the micro-surveys on  $k^{th}$  category. It is straightforward to see that the following statement is always true.

- 1)  $S_{k,uncond} \geq S_{k,cond}$
- 2)  $\hat{y}_{k,uncond} \leq \hat{y}_{k,cond}$

### 5. Constructing Core

It is critical that we include relevant questions in the initial on-boarding section of RSQD which typically includes demographic and relevant auxiliary data. The responses to these questions are used towards the prediction of the expenditure behaviour pattern of the respondents which subsequently affect our *decision rule*.

We may use the Kruskal Wallis Test to understand the

association of variables used in the Core questions with each of the likely expenditure [16]. For each variable we construct the following;

$$H_0: M_1 = M_2 = \dots = M_k$$

$$H_a: M_1 \neq M_2 \neq \dots \neq M_k, \text{ for some } i \neq j$$

Where  $M_k$  represents group k of the variable, the null hypothesis states that all the groups come from the same distribution, whereas the alternate hypothesis states that some of the K groups differ in some ways. In the Kruskal Wallis Test, we reject the null hypothesis for small p-values.

Table 2. Breakdown of Auxiliary and Demographic Data.

| Characteristics  | Levels             | Counts | Percentage |
|--|--------------------|--------|------------|
| Nationality<br>(Top 5 Visitor-ship out of 122 Nationalities) | Indonesian         | 8246   | 17%        |
|  | Chinese            | 7010   | 14.5%      |
|  | Indian             | 3120   | 6.43%      |
|  | Australian         | 3068   | 6.32%      |
|  | Malaysian          | 2776   | 5.72%      |
| Marital Status   | Married            | 29426  | 60.9%      |
|  | Single             | 15659  | 32.3%      |
|  | Divorced           | 576    | 0.01%      |
|  | Widowed            | 282    | 0.01%      |
|  | Refused to provide | 55     | ~ 0%       |
|  | 12                 | 25     | 0.05%      |
|  | 13-14              | 140    | 0.2%       |
| Age Group  | 15-19              | 1443   | 3.0%       |
|  | 20-24              | 3848   | 7.9%       |
|  | 25-29              | 7103   | 14.6%      |
|  | 30-34              | 7381   | 15.2%      |
|  | 35-39              | 6598   | 13.6%      |
|  | 40-44              | 5683   | 11.7%      |
|  | 45-49              | 4184   | 8.6%       |
|  | 50-54              | 3309   | 6.8%       |
|  | 55-59              | 2824   | 5.8%       |
|  | 60-64              | 2084   | 4.3%       |
| Gender   | 65                 | 1376   | 2.8%       |
|  | Male               | 26336  | 57%        |
|  | Female             | 19662  | 43%        |

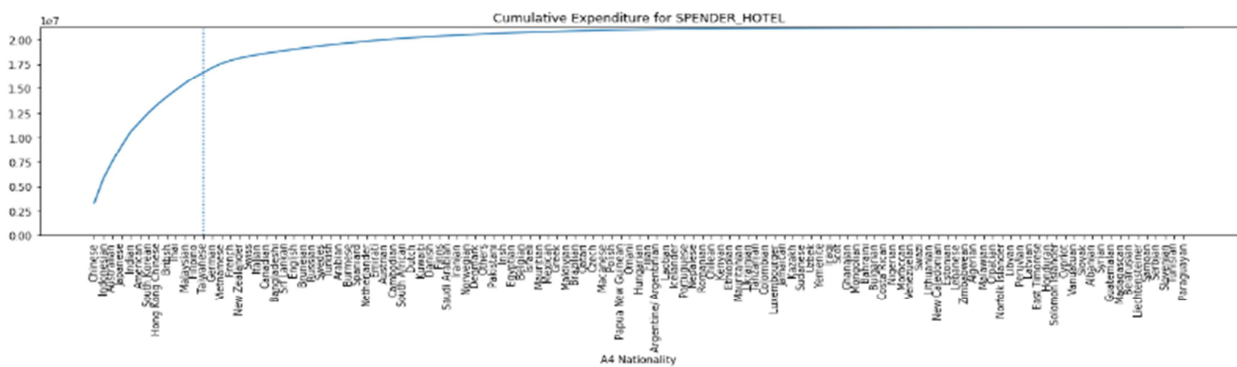


Figure 2. Cumulative Hotel Expenditure sorted by country which contributed the most expenses.

### 6. Data Processing

Feature engineering and data cleaning was carried out to

ensure scalability and meaningful data interpretation. Data pre-processing steps such as top-coding was employed to reduce the influence of extremely high expenditures and outlier management by setting extremely high expenditure

values to values at the 97.5 percentile [11].

To prevent issues of high dimension, variable reduction was practiced. Variable reduction was achieved through keeping only category levels that were significant to the expenditure, for modelling purposes.

In Figure 2, categorical levels were recoded into ‘relevant levels’ and ‘others’. Relevant levels include levels that contributed significantly to the volume of total expenditure. For example, we retained 13 variables and recoded the remaining variables as ‘Others’ for levels for *Nationality* to explain the Hotel category expenditure.

## 7. Methodology

A simulation study was conducted to compare the results of variants of RSQD with SQD. Below are the methods we assessed to determine whether a responsive design provided improvements over non-responsive designs.

### Random Assignment

- 1) Probability Proportional to Size (PPS)
- 2) Logistic Regression (Log-FW1)
- 3) Logistic Regression with more detailed Core (Log-FW2).

$$\begin{aligned} \text{logit}(p_{ik}) = & \beta_0 + \sum_j \beta_{1jk} \times \text{AGEGROUP}_{ij} + \beta_{2k} \times \text{GENDER}_i + \sum_j \beta_{3jk} \times \text{NATIONALITY}_{ij} + \sum_j \beta_{4jk} \times \\ & \text{MARITALSTATUS}_{ij} + \sum_j \beta_{5jk} \times \text{INCOME}_{ij} + \sum_j \beta_{6jk} \times \text{No.WORKINGADULT}_{ij} + \beta_{7k} \times \text{FIRSTVISIT}_i + \\ & \sum_j \beta_{8jk} \times \text{MONTH}_{ij} + \beta_{9k} \times \text{VISITOTHERCOUNTRY}_i + \sum_j \beta_{10jk} \times \text{MODEOFARRIVAL}_{ij} + \beta_{11k} \times \\ & \text{INTENDEDLENGTHofSTAY}_i + \beta_{12k} \times \text{LENGTH of STAY}_i + \sum_j \beta_{13jk} \times \text{TRAVELCOMPANION}_{ij} + \sum_j \beta_{14jk} \times \\ & \text{PURPOSEOFVISIT}_{ij} \end{aligned} \quad (5)$$

To analyse the results of the 4 designs, a simulation run of  $M = 1000$  was performed. With each run, random administration of the micro-surveys was performed via a Bernoulli trial where the parameters of the Bernoulli was determined by  $p_{ik}$  [6].

**Table 3.** Generation of Data from a Bernoulli Distribution.

| Algorithm to generate from $X \sim \text{Bernoulli}(p)$                                 |
|---|
| for $i$ from 1 to $N$ :   |
| Step (i) Generate $u_i$ from Uniform (0, 1)   |
| Step (ii) If $u_i \leq p$ :   |
| set $x_i = \text{Success}$  |
| Else set $x_i = \text{Fail}$  |
| Then $X = \{x_1, x_2, \dots, x_N\}$ follows a Bernoulli distribution with parameter $p$ |

And the following statistics are computed

$$\bar{\theta}_k = M^{-1} \sum_{m=1}^M \hat{y}_{mk} \quad (6)$$

$\bar{\theta}_k$  represents the overall simulation mean of expenditure  $k$  where  $\hat{y}_{mk}$  is the Horvitz Thomson estimator for expenditure  $k$

$$V_k = (M - 1)^{-1} \sum_{m=1}^M (\hat{y}_{mk} - \bar{\theta}_k)^2 \quad (7)$$

Random assignment serves as the base case for comparison. It represents the situation where no information regarding the respondents are used. As such  $p_{ik} = 0.5$  for all expenditures and for all respondents.

PPS serves to improve over the random assignment by altering the respective  $p_{ik}$  for each of the categories based on the proportion of spenders for each expenditure group in the historical data. Higher  $p_{ik}$  is given to expenditure  $k$  where historically expenditure  $k$  is more often incurred. Note that  $p_{ik}$  remains the same for all respondents.

The Logistic Regression methods use a data driven approach that utilises prior information in the form of the Core questions.  $p_{ik}$  values were generated uniquely based on the Core questions that were reported. Higher  $p_{ik}$  was generated for expenditure categories where the respondent was likely to incur. A more detailed Core questionnaire was used in Log-FW2 and had additional attributes such as *Airline*, *Flight class* and *date of last visit*.

Equation (5) shows the logistic regression model used in Log-FW1 that was fitted onto expenditure  $k$ . A total of 15 logistic regression models were fitted, one model for each expenditure category.

$V_k$  is the simulation variance for expenditure  $k$

$$SE_k = \sqrt{V_k} \quad (8)$$

$SE_k$ , the simulation standard error for expenditure  $k$

$$CV_k = SE_k / \bar{\theta}_k \quad (9)$$

$CV_k$  the simulation coefficient of variance for expenditure  $k$ .

$$RB_k = (\bar{\theta}_k - \hat{y}_k) / \hat{y}_k \times 100\% \quad (10)$$

$RB_k$ , the simulation relative bias for expenditure  $k$

$$RBSE_k = SE_k / \hat{y}_k \times 100\% \quad (11)$$

$RBSE_k$ , the simulation standard error for the relative bias  $RB_k$

$$RMSE_k = \sqrt{V_k + RB_k^2} \quad (12)$$

$RMSE_k$ , root mean square error for expenditure  $k$ . Sensitivity, specificity, positive prediction value (PPV) and negative prediction value (NPV) are also computed.

$$\text{Sensitivity} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}} \quad (13)$$



$$Specificity = \frac{True\ Negative}{True\ Negative + False\ Positive} \quad (14)$$

$$PPV = \frac{True\ Positive}{True\ Positive + False\ Positive} \quad (15)$$

$$NPV = \frac{True\ Negative}{True\ Negative + False\ Negative} \quad (16)$$

A metric termed root variance was used to compare all the models against the baseline model. The root variance is the ratio of the variance of the estimate under alternative methods compared to the Random Assignment design. A value less than unity represents a gain in information since we are able to estimate with less variability.

### 8. Evaluation

With the main goal of reducing respondent burden, the focus was on the classification statistics. Sensitivity and specificity values above 0.5 meant that the model was detecting more actual spenders better than a coin flip and hence was preferable. High positive predictive values

(PPV) and negative predictive values (NPV) suggested better precision at detecting actual spenders among those predicted to be a spender and not detecting spenders among those predicted to be not a spender, respectively. A decision is made that a PPV and NPV greater than 1.2\* prevalence rate and 1.2\*(1-prevalance rate) for each expenditure category was preferred. This indicated that the model was detecting 20% more than the prevalence rate. We achieved desirable values for the classification metrics and this suggested high responsiveness for the survey design.

The accuracy of estimation which we term as precision is also an aspect that should be considered. The Coefficient of Variance  $CV_k$  is a normalized measure of dispersion where smaller values signify more precise estimates. A general rule of thumb, is that a  $CV_k$  less than 0.1 is desirable. To have any practical evidence for RSQD we require a root variance value less than 1, since it represents an improvement over the baseline model.

Table 4. Comparison of the Results of the Four Models.

| Expenditure      | Precision |         |         | Responsiveness |         |         |
|------------------|-----------|---------|---------|----------------|---------|---------|
|                  | PPS       | Log-FW1 | Log-FW2 | PPS            | Log-FW1 | Log-FW2 |
| Hotel            | **        | ***     | ***     | *              | ****    | ****    |
| Homestay         |           |         |         | *              | **      | **      |
| Hostel           | *         |         |         | *              | **      | **      |
| Serviceapartment |           |         |         | *              | **      | **      |
| Hawker           | **        | ****    | ****    | *              | *       | **      |
| Casualdining     | **        | ****    | ****    | *              | *       | **      |
| Finedining       | **        | **      | **      | *              | **      | **      |
| Transport        | **        | ****    | ****    | *              | **      | **      |
| Sightseeing      | *         | **      | **      | *              | **      | **      |
| Attractions      | **        | **      | **      | *              | ***     | ****    |
| Entertainment    |           | **      | **      | *              | *       | **      |
| Medical          |           | *       | *       | *              | ***     | ***     |
| Shopping         | **        | ****    | ****    | *              | **      | **      |
| Business         |           | **      | **      | *              | **      | **      |
| Education        |           |         |         | *              | **      | **      |

Table 4 above, shows the results of the 3 models compared to the baseline model. An “\*” was assigned to each model for both the precision and responsiveness whenever the metrics fulfilled the desirable requirements. The shaded cells under the expenditure columns show the proportion of spenders in categories that are significantly small and can be regarded as rare events.

An overall observation that we observed, is that, responsive design does offer the advantage of being more responsive to the tourist with minimal trade-off in the precision domain. The shaded expenditure cells correspond to an expenditure that had low prevalence rate in the training data, where prevalence rate is less than 5%. Further, we observed that, responsive split questionnaire design offers the advantage of

being able to detect these rare expenditures. This suggests that responsive design can provide improvements in issues regarding missing rare events.

#### Effectiveness

To evaluate the effectiveness of the responsive design in reducing respondent burden, we compared the number of irrelevant micro-surveys administered. Shown in Table 5 below, is the average number of irrelevant micro-surveys administered by each model broken down by the expenditure. We observed that the PPS model attained the least number of irrelevant micro-surveys issued (31561). This is a 87% reduction, a significant improvement compared to the Random case. We also noted that our two logistic regression models attained a respectable improvement (75%) over the random design.

**Table 5.** Summary of the Average Number of Irrelevant Micro-Surveys Issued.

| Average number of Irrelevant Micro-Survey Issued |        |       |         |         |
|--|--------|-------|---------|---------|
|  | Random | PPS   | Log-FW1 | Log-FW1 |
| Hotel  | 7220   | 1362  | 4984    | 4907    |
| Homestay   | 22872  | 389   | 347     | 354     |
| Hostel   | 22595  | 6125  | 787     | 763     |
| Serviceapartment                                 | 22916  | 230   | 293     | 296     |
| Hawker   | 9412   | 436   | 9680    | 9481    |
| Casualdining                                     | 9568   | 282   | 9620    | 9460    |
| Finedining                                       | 21535  | 5912  | 2596    | 2672    |
| Transport  | 5038   | 50    | 6707    | 6469    |
| Sightseeing                                      | 18673  | 187   | 6082    | 7013    |
| Attractions                                      | 13643  | 4349  | 7998    | 7857    |
| Entertainment                                    | 20704  | 207   | 3668    | 3783    |
| Medical  | 22592  | 225   | 427     | 396     |
| Shopping   | 4119   | 1571  | 5568    | 5588    |
| Business   | 22158  | 1937  | 1368    | 1366    |
| Education  | 22889  | 8300  | 359     | 407     |
| TOTAL  | 245934 | 31561 | 60486   | 60813   |

## 9. Conclusion

We have explored the use of the split questionnaire design to obtain the tourist mean expenditure across 15 different categories. The objective of this study was to reduce the respondent burden and to increase the response rate by using the split questionnaire design. We leveraged on the idea of issuing only relevant surveys to each respondent as this directly reduces the connotation of being probed with irrelevant questions which is regarded as a primary reason for non-response. To achieve this, we developed a responsive design by utilising the prior information that we collected as part of the questionnaire, and the created decision rules for the administration of the relevant micro-surveys. The results that we have attained are promising in terms of the trade-off between precision and responsiveness. Besides meeting these objectives, we also demonstrated that a secondary advantage of responsive split questionnaires also brings about significant cost savings to the survey organisation.

To improve upon our results, we feel that more sophisticated classification models such as classification trees or deep learning models can be explored in order to increase the classification power in detecting expenditure patterns among a highly heterogeneous population. In addition, we also feel that variables that might explain different expenditure patterns can be explored and be included in the core questionnaire.

## References

- [1] Adiguzel F. and Wedel M. (2008). Split Questionnaire Design for Massive Surveys. *Journal of Marketing Research*, 45 (5), 608-612.
- [2] Ali, M., Kauermann, G. (2021). A split questionnaire survey design in the context of statistical matching. *Stat Methods and Applications*. <https://doi.org/10.1007/s10260-020-00554-2>
- [3] Breidt J. and Opsomer J. (2017). Model-Assisted Survey Estimation with Modern Prediction Techniques. *Statistical Science*, 190-205.
- [4] Chun A. Y., Heeringa G., Schouten B. (2018). Responsive and Adaptive Design for Survey Optimization. *Journal of Official Statistics*, pp 581-597. doi: 10.2478/jos-2018-0028.
- [5] Chipperfield, J. O., Barr, M. L. and Steel, D. G., (2018). Split Questionnaire Designs: collecting only the data that you need through MCAR and MAR designs. *Journal of Applied Statistics*, 1465-1475. doi: 10.1080/02664763.2017.1375085.
- [6] Christian P. Robert, G. C. (2004). *Monte Carlo Statistical Methods*. Springer. doi: 10.1007/978-1-4757-4145-2.
- [7] Cochran, W. G. (1977). *Sampling techniques 3rd Edition*. New York: Wiley.
- [8] Early K, Mankoff J, Fienberg S. (2017). Dynamic Question Ordering in Online Surveys. *Journal of Official Statistics*, 625-657. doi: 10.1515/jos-2017-0030.
- [9] Gonzalez, M. (2012). *The use of responsive split questionnaires in a panel survey*.
- [10] Groves M. and Heeringa G. (2006). Responsive design for household surveys: tools for actively controlling survey errors and costs. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 169 (3), pp. 439-457. doi: 10.1111/j.1467-985X.2006.00423.x.
- [11] Hundepool, A. (2012). *Statistical Disclosure Control - Wiley Series in Survey Methodology*. John Wiley & Sons Inc.
- [12] Murphy J., Biemer P. and Berry C. (2018). Transitioning a Survey to Self-Administration using Adaptive, Responsive, and Tailored (ART) Design Principles and Data Visualization. *Journal of Official Statistics*, Vol. 34, 625-648. doi: 10.2478/jos-2018-0030.
- [13] Rolstad, S. (Dec, 2011). Response Burden and Questionnaire Length: Is Shorter Better? A Review and Meta-analysis. *Value in Health*, 1101-1108. doi: <https://doi.org/10.1016/j.jval.2011.06.003>
- [14] Thompson, D. G. (1952). A Generalization of Sampling Without Replacement From a Finite Universe. *American Statistical Association*, 663-685.
- [15] Wagner, J. R. (2008). Adaptive survey design to reduce nonresponse bias. *PhD Thesis University of Michigan*.
- [16] Wallis, William H. Kruskal and W. Allen. (1952). Use of Ranks in One-Criterion Variance Analysis. *Journal of the American Statistical Association*, 583-621. doi: 10.2307/2280779.