

Individual Coursework Assignment

Introduction

The individual coursework assignment (100% of the course unit marks) will seek to demonstrate your practical understanding of the marketing analytics concepts and techniques discussed throughout the semester.

This document consists of 9 pages.

Learning Objectives

The main objectives of this assessment are:

- Possess an applied understanding of different analytics methods
- Demonstrate the ability to use the R programming environment to manipulate data and implement statistical modelling and analysis in order to address marketing problems
- Break complex tasks into parts and steps
- Plan and manage time
- Communicate findings and ideas professionally and creatively

Overview

You will take over the role of *Market and Insights Analyst* at the consulting services department of a multinational professional services firm. As part of this role, you are asked to work across the following three (3) different client engagement projects.

All projects carry the same weighting. Each project includes a series of tasks summing up to 100 points. The descriptions of the three (3) client engagement projects are available at the end of this document.

Submission and Word Limit

You will need to produce one (1) written **report**. The word limit for the report is 3,500 words ($\pm 10\%$). Your report must be submitted as a Microsoft Word document (.docx).

As an **appendix**, you will also need to provide either an Rmarkdown file (.Rmd) or an R Script file (.R) that includes all code commands necessary to reproduce your analyses and graphs. You will not get any marks without the relevant code commands.

Layout and Style

The report should follow the suggested format guidelines: Font Arial 11pt, single line spacing, page numbers, numbered headings, numbered and labelled figures and tables.

Collaboration

Please ensure that language, thoughts, ideas, code commands, and any other forms of intellectual output developed by you are not shared with other students as this constitutes a form of plagiarism. Please familiarise yourselves with the Plagiarism information on MySurrey.

Marking Structure

You may also wish to view the University grade descriptors to review the generic statements that describe achievement in terms of the range and breadth of knowledge and abilities a student is required to achieve. These can be accessed in the Assessments section of MySurrey.

Submission Deadline and Method

You need to submit your report and appendix by **Wednesday 25 May 2022, 4pm**. You will need to submit via the dedicated folder on the course SurreyLearn page. No email submissions will be accepted.

Project 1

For this project you will analyse data from a survey in which 200 respondents were asked to rate the importance of a number of store attributes when choosing where to buy office equipment. The file `office.csv` contains data for the project. For each respondent, we have the following variables:

Variable	Description
<code>respondent_id</code>	An identifier for our observations
<code>variety_of_choice</code>	Importance of this attribute on a 0-10 scale
<code>electronics</code>	Importance of this attribute on a 0-10 scale
<code>furniture</code>	Importance of this attribute on a 0-10 scale
<code>quality_of_service</code>	Importance of this attribute on a 0-10 scale
<code>low_prices</code>	Importance of this attribute on a 0-10 scale
<code>return_policy</code>	Importance of this attribute on a 0-10 scale
<code>professional</code>	Whether the respondent is a professional or not (e.g., student)
<code>income</code>	Gross annual income expressed in thousands of pound sterling
<code>age</code>	Respondents' age in years

Task	Points
1 Read and inspect the data set. Provide a descriptive analysis for each of the variables in the data set.	5
2 Make a new data object (e.g., a <code>data.frame</code> or <code>tibble</code>) for clustering that includes only the attitudinal variables from the original data set. Then normalise (use z-score standardisation) all variables in this new data object. Which variable has the smallest minimum value and which variable has the largest maximum value in the normalized data set?	5
3 Run the hierarchical clustering algorithm using <code>method = "ward.D2"</code> on the normalised data and use <code>set.seed(123)</code> for reproducibility. Plot the dendrogram.	5
4 Suppose that after looking at the dendrogram and discussing with the marketing department, you decide to proceed with a 6-cluster solution. Divide the data points into 6 clusters. How many observations are assigned to each cluster?	5
5 Use the normalised data to calculate the means for each of the attitudinal variables per cluster. Use the <code>flexclust</code> package to generate a segment profile plot. Comment on whether any cluster memberships	10

have changed, if any. Check the concordance between the `hclust` and `as.kcca` procedures.

- | | | |
|----|--|----|
| 6 | Describe the 6-cluster solution using the cluster numbers corresponding to the hierarchical clustering procedure. | 10 |
| 7 | Comment on why you may decide to NOT proceed with this 6-cluster solution. | 10 |
| 8 | Generate a 5-cluster solution. How many observations are assigned to each cluster? | 5 |
| 9 | Repeat the steps performed previously to describe the clusters for the 5-cluster solution (i.e., calculate cluster means and segmentation plot). Describe the 5-cluster solution using the cluster numbers corresponding to the hierarchical clustering procedure. Give “expressive” labels to the clusters. | 10 |
| 10 | Comment on why you may find this 5-cluster solution better than the previous 6-cluster solution. | 10 |
| 11 | Use all the variables not included in the clustering procedure to evaluate whether the 5-cluster solution is meaningful. Generate ideas on how to target each segment (at least one idea per segment). | 15 |
| 12 | Run the k-means clustering algorithm on the normalised data, creating 5 clusters. Use <code>iter.max = 1000</code> and <code>nstart = 100</code> and <code>set.seed(123)</code> for reproducibility. How many observations are assigned to each cluster? | 5 |
| 13 | Check the concordance between the <code>hclust</code> and <code>kmeans</code> procedures. What is the Hit Rate? | 5 |

Project 2

For this project you will model website user conversion. You will be working on a dataset with more than 20 thousand unique users of a website based in four countries. The file `ecommerce.csv` contains data for the project. For each user, we have the following variables:

Variable	Description
<code>country</code>	The country the user accessed the site from (France, Germany, Ireland, or UK)
<code>source</code>	The source through which the user accessed the site (ads, search, or direct link)
<code>total_pages_visited</code>	The number of pages visited by the user
<code>visit_duration</code>	The amount of time the user spent in the site (in seconds)
<code>discount</code>	Whether the user was offered a discount (10% off first order; yes, no)
<code>conversion</code>	Whether the user converted, or made a purchase (yes, no)

Task	Points
1 Read and inspect the data set. Provide a descriptive analysis for each of the variables in the data set.	5
2 Build a simple logistic regression model of <code>conversion</code> on <code>discount</code> . Call this model <code>m1</code> . Comment on the coefficient estimate of <code>discountyes</code> . What is the sign of the coefficient? Is the effect statistically significant?	5
3 Calculate the odds ratio for <code>discountyes</code> . What does this mean?	5
4 Calculate the 95% confidence interval for the odds ratio for <code>discountyes</code> . What does this mean?	5
5 Generate a double-decker mosaic plot (using the <code>ggmosaic</code> package) to visualise the count of the combinations of the following variables: <code>discount</code> (on x-axis), <code>conversion</code> (as fill colour), and <code>source</code> (as facets). Use the plot to describe whether (and how) the effect of <code>discount</code> on <code>conversion</code> is different for the three <code>source</code> channels.	5
6 Build a logistic regression model that predicts <code>conversion</code> from <code>discount</code> and <code>source</code> . Call this model <code>m2</code> . Comment on the coefficient estimates of <code>sourcedirect</code> and <code>sourcesearch</code> .	5

- | | | |
|----|---|----|
| 7 | Calculate the odds ratios for <code>sourcedirect</code> and <code>sourcesearch</code> . What do these mean? | 5 |
| 8 | Build a logistic regression model that predicts <code>conversion</code> from <code>discount</code> and <code>source</code> and also includes their interaction. Call this model <code>m3</code> . Comment on the coefficient estimates of the interaction terms. | 5 |
| 9 | Calculate the 95% confidence intervals for the odds ratios for the interaction terms. What do these mean? | 5 |
| 10 | Build a logistic regression model that predicts <code>conversion</code> from all available variables in the data set. This model should also include an interaction of the <code>discount</code> and <code>source</code> variables. Call this model <code>m4</code> . Which variables are significant at the 95% level? | 5 |
| 11 | Calculate the correlation between the two numerical variables in the data set (<code>total_pages_visited</code> and <code>visit_duration</code>). Comment on the result. How may this affect <code>m4</code> ? | 5 |
| 12 | Build another logistic regression model from <code>m4</code> by removing the <code>total_pages_visited</code> variable. Call this model <code>m5</code> . How has the effect of <code>total_pages_visited</code> changed compared to <code>m4</code> ? | 5 |
| 13 | Make a plot that visualises the odds ratios (as points) of the variables in <code>m5</code> as well as their confidence intervals (as error bars). | 10 |
| 14 | Use model <code>m5</code> to predict the conversion probabilities for each user in the data set. Store these probabilities in the data set, in a variable called <code>base_prob</code> . What is the mean value of <code>base_prob</code> ? | 5 |
| 15 | Calculate an indicator variable for whether individuals will convert or not, based on their predicted probabilities from the previous task, using a threshold value of 0.5. Call this variable <code>pred_conversion</code> . How many users so we predict to convert? | 5 |
| 16 | What is the accuracy or hit rate? | 5 |
| 17 | What is the area under the curve? | 5 |
| 18 | Predict new probabilities under a hypothetical scenario that the values variable <code>total_pages_visited</code> were increased by one unit (i.e., one page) for all users. Store these probabilities in the data set, in a variable called <code>new_prob</code> . What is the mean value of <code>new_prob</code> ? | 5 |
| 19 | Calculate the lift metric for the hypothetical scenario from the previous task (i.e., Task 18). | 5 |

Project 3

For this project you will run a Choice-Based Conjoint study in the Cloud Services Platform market (e.g. Amazon Web Services, Google Cloud, Microsoft Azure). The client wants to make some product design decisions such as core feature-sets, pricing, and tiers of service to optimise revenue or new sign-ups.

You will work with the `cloud.csv` file. The file contains data on choices made by 200 respondents. Each respondent evaluated 15 choice sets. Thus, the file contains data on $200 \times 15 = 3000$ choice sets. Each choice set had three alternatives. A respondent's task was to choose one alternative from a choice set. The following table describes the variables in the dataset:

Variable	Description
<code>respondent_id</code>	Identifier for each respondent (1 to 200)
<code>choiseset_id</code>	Identifier for each choice set for each respondent (1 to 15)
<code>alternative_id</code>	Identifier for each alternative in a choice set (1 to 3)
<code>choice_id</code>	Identifier for each choice set in the entire study (1 to 3000)
<code>cloud_storage</code>	Attribute cloud storage with three levels: 30GB / 2000GB / 5000GB
<code>customer_support</code>	Attribute customer support with two levels: Yes / No
<code>cloud_services</code>	Attribute cloud services with three levels: Email / Email + Video / Email + Video + Productivity
<code>price</code>	Attribute price with three levels: £6 per month / £12 per month / £18 per month
<code>choice</code>	Shows which alternative was chosen in each choice set (Dummy coded: 1 if alternative was chosen; 0 otherwise)

Task	Points
1 Read and inspect the data set. Provide a descriptive analysis for each of the variables in the data set. Make sure you provide an analysis that is meaningful for each variable type (e.g., factors, identifiers).	5
2 Convert the attribute variables <code>cloud_storage</code> and <code>price</code> so that the factor reference levels are the levels representing the smallest values (i.e., 30GB for <code>cloud_storage</code> and p6 for <code>price</code>). Why there is no need to perform this step on the rest of the attribute variables?	5
3 Create a new variable in the data set that turns <code>price</code> into numeric class (do not overwrite <code>price</code>). Call this new variable <code>price_n</code> . What is the mean of variable <code>price_n</code> ?	5

- 4 There are 3000 choice sets in the data set. Therefore, there were 3000 choices made. Out of these 3000 choices, how many times did respondents choose a 30GB cloud storage? What is the percentage of respondents who chose email only as cloud service? 5
- 5 Use the `dfidx()` function from the `dfidx` package to create a specially formatted data object that will be used in the process of estimating a multinomial conjoint model. In the argument `idx`, use a `list` of the two indexes (`choice_id` and `respondent_id`) that define unique observations. Also use `alternative_id` as the variable defining the levels of the alternatives. Call this data object `m_data`. How many variables (i.e., columns) does `m_data` have? 5
- 6 Use `m_data` to build a multinomial logit model that predicts `choice` from `cloud_storage`, `customer_support`, `cloud_services`, and `price`. Make sure that you tell the `mlogit()` function to exclude the intercept term. Call this model `model1`. Use `set.seed(123)` right before running the command that builds the model. Comment on the coefficient estimates of `cloud_storage5000gb` and `pricep12`. 5
- 7 Now follow the same process as in Task 6 to build a multinomial logit model that uses `price_n` instead of `price`. Call this model `model2`. Again use `set.seed(123)` right before running the command that builds the model. Comment on the coefficient estimate of `price_n`. What does this mean? 5
- 8 Use a likelihood ratio test to test the `model2` against `model1`. What is the outcome of the test? Are `model2` and `model1` significantly different? Which model we should choose between the two and for what reason(s)? 5
- 9 Use `model2` to predict the choice probabilities for different alternatives in the data. What is the predicted probability of choosing the third alternative in the first choice set? 5
- 10 Use the predicted probabilities from Task 9 to compute the predicted alternatives using the maximum choice probabilities. Which is the predicted alternative in the third choice set? 5
- 11 Then we can extract the selected alternatives from the original data. Which is the selected alternative in the fifteenth choice set? 5
- 12 Compute the confusion matrix for `model2`. What is the accuracy (or hit rate) of `model2`? How does `model2` compare to the baseline method (i.e., making random predictions)? 10

- 13 Now let us see how we can use the `model2` parameters to predict market shares under hypothetical market scenarios for an arbitrary set of products. First, build a custom function to predict market share for an arbitrary set of alternatives available in a data set `d`. You can find the commands for building the custom function in the “Multinomial Choice Modelling Practical”. Call the custom function `predict.share`. 0

- 14 Create a data object (i.e., `data.frame` or `tibble`) with the following hypothetical market consisting of five alternatives: 0

<code>cloud_storage</code>	<code>customer_support</code>	<code>cloud_services</code>	<code>price_n</code>
30gb	no	email	6
30gb	no	email, video	12
30gb	yes	email	12
5000gb	yes	email	18
5000gb	no	email, video, productivity	18

Call this data object `d_base`.

- 15 Run the customer function `predict.share` using `model2` and `d_base` as input arguments. What is the predicted market share for alternative four of this hypothetical market? 5
- 16 Now consider a modification on the previous hypothetical market, in which the level of the `cloud_services` attribute changes for the fifth alternative to “email, video”. What is the predicted market share for alternative four of this new hypothetical market? 5
- 17 Which alternative was affected the most from this modification of the hypothetical market, and by how much (in percentage terms)? 10
- 18 Use the `model2` coefficients to calculate how much a consumer would be willing to pay (in £ per month) for customer support. 5
- 19 Use the `model2` coefficients to calculate how much a consumer would be willing to pay (in £ per month) for an upgrade from 30GB to 2000GB cloud storage. 5
- 20 Use the `model2` coefficients to calculate how much a consumer would be willing to pay (in £ per month) for an upgrade from 2000GB to 5000GB cloud storage. 5

END OF DOCUMENT