

Practicum 2 (part A)

Choice model specification and estimation (LCM)

1. Estimation and assessment of the Latent Class Model (LCM)

Last week, you gained first-hand experience building and estimating the Multinomial Logit (MNL) model. You applied the MNL model to the Colorectal Cancer (CRC) screening dataset, which investigated the influence of screening characteristics on non-participation in CRC screening in the Netherlands. Using the MNL model, you obtained insights into preferences by assuming all respondents are identical and have the same preferences (i.e., preference homogeneity). Assuming preference homogeneity is often an unrealistic assumption.

In this practicum, you will learn to build the Latent Class Model (LCM). The LCM assumes that there are (latent) classes of respondents with different preferences but that all respondents are identical within each class. Preference heterogeneity is identified from unobserved (random) variables.

Before moving on, please follow the steps below to set up your environment correctly:

1. Go to Canvas > Modules > Week 5 > Tutorial 2
2. Download the syntax file “**Practicum2_lcm.R**”
3. Go to download and find “**Practicum2_lcm.R**”
4. Move the syntax file to the folder created for the last week’s practical
5. Open the syntax file “**Practicum2_lcm.R**”, and let’s build an Latent Class model

When estimating an LCM, there is a trade-off between a good model fit (resembling the data as closely as possible) and the number of parameters (looking for the most parsimonious model). The more (latent) classes you have in your data, the better the model fit of the LCM will be (do you know why?), but the more parameters are in the model (do you know why?).

Before you start, let us remember Apollo’s basic structure that we introduced last week.

How does the Apollo package work?

In the Apollo package, all types of choice models use the same code structure. Each choice model is created by using the same building blocks. As the model complexity increases, components are added to the standard building blocks.



The standard building blocks to build a Latent Class model are:

- *Building block: Inputs*

In the “Inputs” block (see steps 1-9 in the syntax file), the set of information must be explicitly defined to set up the R-environment correctly for the Apollo package. As we will see, there is one additional component (step 9) in the “Inputs” block compared to the MNL syntax.

We start by setting up the R-environment.

- Step 1: clear the memory of the R-environment

```
36  ### Step 1: Clear memory  
37  rm(list = ls())
```

This step is optional but highly recommended (!). Clearing the memory removes all saved variables, data frames, and objects in your R-environment. When you do not clear the environment and you are (re)running a (adjusted) syntax file, it could result in error messages.

- Step 2: Set a working directory for the R-environment

```
39  ### Step 2: Set working directory for R initialization  
40  setwd(dirname(rstudioapi::getActiveDocumentContext()$path))
```

This step is optional but highly recommended (!). In this line, we instruct the R-environment to use the working directory (i.e., folder) in which the syntax file is stored. By doing so, we do not need to specify the full path to specific files, such as the dataset file.

- Step 3 + 4: Load and initialize the Apollo code

```

42  ### Step 3: Load Apollo library
43  library(apollo)
44
45  ### Step 4: Initialise Apollo code
46  apollo_initialise()

```

Here, we load and initialize the Apollo package into the R-environment.

- Step 5: Set the core controls

```

48  ### Step 5: Set core controls
49  ###      ATTENTION: Your inputs must be enclosed in quotes like "this"
50  apollo_control = list(
51    # USER ACTION: Specify model name
52    ## Note: Change the model name for every model that you run
53    modelName      = ,
54    # USER ACTION: Provide model description
55    ## Note: Change the model description to reflect the current model
56    modelDescr     = ,
57    # USER ACTION: Specify the column with the respondent id
58    indivID        = ,
59    # Define number of cores used during estimation (used to speed up estimation time)
60    nCores         = 5,
61    # Define seed used for any random number generation
62    seed            = 100,
63    # USER ACTION: Set path to the folder on your PC where the model results will be stored
64    ## Note: Use the "outputs" folder that was created by the pre-processing syntax
65    outputDirectory = paste(getwd(), 'outputs', "practicum_lcmmodel_3class_covar", sep=Platform$file.sep)
66  )

```

As we have loaded and initialized the Apollo package in steps 3 and 4, we need to instruct Apollo how to call our model (= *modelName*), what our model is (= *modelDescr*), which column of our dataset contains the respondent identifiers (= *indivID*), and where to save the model results (*outputDirectory*).

Note that, compared to the MNL, one control (“*seed*”) was added to the Latent Class model. It controls a random number generation process. Since the log-likelihood function of the LCM is not concave, there may exist multiple optima in the function that “trick” the optimization of the coefficients to an early stop. To avoid this issue, Apollo runs the Latent Class model several times using different random starting coefficient values to identify the best model. To ensure the results are reproducible, this control fixes the seed of the random number generation process.

- Step 6: Load the dataset into the R-environment

```

64  ### Step 6: Load data
65  # Set path to the folder on your PC where the dataset is stored
66  path_data = paste(getwd(), 'data', "dataset.csv", sep=Platform$file.sep)
67  # Load dataset into global environment
68  database = read.csv(path_data, header=TRUE)

```

Because we have set the working directory (see step 2) in which the syntax file and dataset are stored, we do not need to specify the path to the dataset. We only need to check the name of the dataset whether it is correct or not (see purple names).

- Step 7: Initialize all parameters that need to be estimated in the choice model

Here, taste (e.g., attribute coefficients and alternative specific constants) and class membership parameters (e.g., deltas and gammas) must be initialized with starting values prior to the model estimation procedure.

```

74  ### Step 7: Initialise all parameters that needs to be estimated in your MNL model
75  # USER ACTION: Define the (1) class-specific and (2) class membership parameters
76  #           followed by assigning a starting value. The class-specific
77  #           alternative specific constant for the opt-out option and the
78  #           constants for the class membership models are already defined.
79  #           Please, complete the list with the parameters that are missing.
80  #           Provide names for each parameter following by assigning a
81  #           starting value.
82  apollo_beta=c(# Class 1
83          asc_out_1      = 0,
84
85          # Class 2
86          asc_out_2      = 0,
87
88          # Class 3
89          asc_out_3      = 0,
90
91          # Class membership - class 1
92          delta_1        = 0,
93
94          # Class membership - class 2
95          delta_2        = 0,
96
97          # Class membership - class 3
98          delta_3        = 0)

```

Note the differences between this code and the MNL.

- Taste coefficients: in the Latent Class model you specify one vector of parameters for each desired latent class.
- Class membership coefficients: these are new parameters specific to the Latent Class model. For each latent class there should be:

- A constant “*delta*” that supports the identification of class membership (and estimated class sizes) from unobserved heterogeneity. Alone, this constant assigns the same class membership probabilities to every respondent.
- You can add additional coefficients (let us say *gamma*) to be estimated with support of observed variables (e.g., sociodemographic, attitudes). Combined with the constant, these coefficients allow each respondent to have different membership probabilities conditional on their observed characteristics.

Action required! Go to your **Practicum5_lcmodel.R** file, read the user actions, and complete the code in step 7. But do not run the syntax file yet (!).

▪ Step 8: Define which parameters should kept fixed during the estimation

```
100  ### Step 8
101 ## USER ACTION: Complete the list with parameters (as initialised above) that
102 ##           should be kept fixed during estimation (in quotes)
103 apollo_fixed = c("delta_3")
```

Remember that “*apollo_fixed*” is a list of parameters that should be kept fixed at their initial values. Just like with one level of each categorical attribute, the membership coefficients of one latent class also can't be estimated.

Action required! Go to your **Practicum5_lcmodel.R** file, read the user actions, and complete the code in **step 8**. But do not run the syntax file yet (!).

- **Step 9:** Detail the class membership model

This is the new component of the “*Inputs*” block. Apollo does not recognize the structure of the coefficients initialized in **step 7** (“*apollo_beta*”). This component is where the coefficients are mapped into the Latent Class model structure.

```

105  ### Step 9: Define class membership model
106  apollo_lcPars=function(apollo_beta, apollo_inputs){
107    lcpars = list()
108    ## USER ACTION: Complete the empty lists by specifying the missing class-specific parameters
109    ## which are needed for the class-specific utility functions
110    lcpars[["asc_out"]] = list(asc_out_1, asc_out_2, asc_out_3)
111    lcpars[["b_eff"]] = list()
112    lcpars[["b_fneg"]] = list()
113    lcpars[["b_freq_2yr"]] = list()
114    lcpars[["b_freq_3yr"]] = list()
115    lcpars[["b_wdiag_2wks"]] = list()
116    lcpars[["b_wdiag_3wks"]] = list()
117    lcpars[["b_wfup_4wks"]] = list()
118    lcpars[["b_wfup_8wks"]] = list()
119
120    ## List of class-membership functions:
121    ## These must use the same names as in classAlloc_settings (see below), order is irrelevant
122    V=list()
123    # USER ACTION: Define class-membership function for class 1
124    V[["class_1"]] = delta_1
125
126    # USER ACTION: Define class-membership functions for class 2
127    V[["class_2"]] = delta_2
128
129    # USER ACTION: Define class-membership functions for class 3
130    V[["class_3"]] = delta_3
131
132    ## Define settings for class-membership model
133    classAlloc_settings = list(
134      # USER ACTION: Attach class-membership functions to the respective classes
135      classes      = c(class_1=, class_2=, class_3=),
136      # USER ACTION: Define which classes are "available" in our study, all classes are "available"
137      avail        = ,
138      # USER ACTION: Attach list of class-membership functions
139      utilities    =
140    )
141
142    lcpars[["pi_values"]] = apollo_classAlloc(classAlloc_settings)
143    return(lcps)
144  }

```

The “*lcPars*” list maps the taste coefficients into the attributes. You can recognize the attributes from last week’s practical. In the Latent Class model, for every attribute/attribute level there will be a sub-list of the coefficients defined in **step 7** associated with the attribute/attribute level. Each sub list will contain as many coefficients as the desired latent classes.

The “*V*” list maps the coefficients of the membership function into each latent class. When the membership class has only constants the mapping is as you see in the figure above. To add observed variables, the mapping should be completed multiplying one coefficient by one observed variable as in the utility function.

Action required! Go to your **Practicum5_lcmodel.R** file, read the user actions, and complete the code in **step 10**. But do not run the syntax file yet (!).

- *Building block: Validate inputs*

The final step in preparing the code and data for model estimation is to make a call to:

- *Step 10: Validation of the inputs*

```
146  ### Step 10: Checkpoint for model inputs
147  apollo_inputs = apollo_validateInputs()
```

The function runs several checks and produces a consolidated list of model inputs. This function takes no arguments but looks at the R-environment for the various inputs required to build and estimate the model.

- *Step 11: Define the choice model and likelihood function*

```
149  ### Step 11: Define model and likelihood function
150  apollo_probabilities=function(apollo_beta, apollo_inputs, functionality="estimate"){
151
152  ### Attach inputs and detach after function exit
153  apollo_attach(apollo_beta, apollo_inputs)
154  on.exit(apollo_detach(apollo_beta, apollo_inputs))
155
156  ### Create list of choice probabilities P
157  P = list()
158
159  ### Define settings for MNL model component that are generic across classes
160  mnl_settings = list(
161    # USER ACTION: Attach utility functions to the alternatives in your dataset
162    alternatives = c(ALT1=, ALT2=, ALT3=),
163    # USER ACTION: Define which alternatives are "available" in each choice task; in our study, all alternatives are "available"
164    avail      = list(ALT1=, ALT2=, ALT3=),
165    # USER ACTION: Specify the column containing the chosen alternative; beware, no dummies are used (!)
166    choiceVar  =
167  )
168
169  ### List of utility functions for each latent class: these must use the same names as in mnl_settings (see above), order is irrelevant
170  # USER ACTION: Set number of latent classes you are estimating in the model
171  ## Note: You can call class-specific parameters by NAME_PARAM[[s]]; see example in ALT3 for the class-specific
172  ## alternative specific constant
173  for(s in 1:'REPLACE THIS WITH THE NUMBER OF CLASSES YOU WANT TO ESTIMATE'){
174    V=list()
175    # USER ACTION: Define utility function for alternative 1 for class "s"
176    V[["ALT1"]] =
177
178    # USER ACTION: Define utility function for alternative 2 for class "s"
179    V[["ALT2"]] =
180
181    # USER ACTION: Define utility function for alternative 3 for class "s"
182    V[["ALT3"]] = asc_out[[s]]
183
184    mnl_settings$utilities = V
185    mnl_settings$componentName = paste0("Class_",s)
186
187    ### Compute within-class choice probabilities using MNL model
188    P[[paste0("Class_",s)]] = apollo_mnl(mnl_settings, functionality)
189
190    ### Take product across observations for same individual (i.e., considering the panel structure of the data)
191    P[[paste0("Class_",s)]] = apollo_panelProd(P[[paste0("Class_",s)]], apollo_inputs ,functionality)
192  }
193
194  ### Compute latent class model probabilities
195  lc_settings = list(inClassProb = P, classProb=pi_values)
196  P[["model"]]= apollo_lc(lc_settings, apollo_inputs, functionality)
197
198  ### Prepare and return outputs of function
199  P = apollo_prepareProb(P, apollo_inputs, functionality)
200  return(P)
201 }
```

In the figure above, the code structure *for(s in 1:S) {}* is a loop. It means that the command inside the loop will be repeated *S* times and *[[s]]* will be replaced by the number controlling the loop lap (1, 2, or 3 if *S*=3).

Therefore, there will be 9 utility functions or 3 (as in the MNL) for each class. The loop structure helps to automate the process of writing the utility functions.

Question 1.1 Consider an LCM with 2 and 3 classes. Write down the class-specific utility functions for all three choice alternatives, assuming that the attributes of “Effectiveness” and “Risk false negative” are coded as linear.

Answer:

Question 1.2 What is the difference in the number of parameters between the model with 2 and 3 classes, considering that no explanatory variables are included?

Answer:

Before moving on to the following questions, please carefully read the following points:

- The LCM probabilistically allocate respondents to classes during the estimation procedure. Because we add and/or remove variables in the model and because of the probabilistic nature of the class-membership model, the LCM can use different classes as a reference class each time you (re-)estimate the model. Although the classes are the same in terms of class shares and class-specific preference parameters, the order of the classes can differ from already estimated LCMs.
 - When (re)estimating a choice model in Apollo, please do not forget to change the model’s name and output directory in step 5 of the R-code.
- *Building block: Estimation*

When you have completed step 11, you have initialized all essentials and built your MNL model. Now, it is time to start the estimation procedure.

- Step 12: Search for starting value

```
203  ### Step 12: Searching for starting value (recommended to ensure model convergence!)
204  apollo_beta = apollo_searchStart(apollo_beta,
205      apollo_fixed,
206      apollo_probabilities,
207      apollo_inputs,
208      searchStart_settings=list(nCandidates=2))
```

This is another new step compared to the MNL model. In contrast with the MNL, the likelihood function of the MNL is not concave. It means that it might not have a global optimum. To avoid the risk of being trapped in a local optimum, Apollo initiates the estimation from several starting values and selects the most promising one.

- Step 13: Model estimation

```
210  ### Step 13: Model estimation
211  model = apollo_estimate(apollo_beta, apollo_fixed, apollo_probabilities, apollo_inputs)
```

By running step 13, the estimation procedure will be started by Apollo. There are no actions required. It will be all handled automatically for you.

- *Building block: Output*

The last building block is to instruct Apollo to print the model output and save the results in the folder that we have specified in step 5 (see “*outputDirectory*”).

- Step 14: Print model output with two-sided p-values

```
213  ### Step 14: Print model output with two-sided p-values
214  ### Note: if one-sided p-values are needed, set "printPVal=1" (p-values are not reported if set to "0")
215  modelOutput_setting=list(printPVal=2)
216  apollo_modelOutput(model, modelOutput_setting)
```

After completing the estimation procedure, we instruct Apollo on how to process the model results. For now, you do not need to adjust anything. However, in some situations, you might want additional information. The Apollo manual explains what information you can extract and how you can do this (see [here](#), page 37).

Now let us interpret the results of the Latent Class model that you estimated.

Question 1.3 Complete the R-code by following the instructions indicated as # **USER ACTION** and estimate the LCM with a different number of classes to see which model fits best. You may limit your estimations to 2-3 classes to save time. Based on the model results, how many classes would you select for your analysis and why?

Answer:

In the following questions, we will use the LCM with 3 classes.

Question 1.4 Which class is larger, and how many respondents are assigned each class?

Answer:

Question 1.5 What do you conclude when you compare class 2 with the MNL model you estimated in Question 2.7 in last week's practicum?

Answer:

Question 1.6 What do you conclude about class 3, compared to class 2?

Answer:

Question 1.7 The CRC dataset also includes socio-demographic variables. Include gender, age, and educational background in the LCM with 3 classes as explanatory variables for the class membership. Estimate the model again. Does gender, age, and educational background relate to class membership and if so, how?

Answer: