

# Christmas boxes: a conjoint analysis application

Laboratory of Customer and Business Analytics

University of Trento - AA 2019/2020

Ambrosi Andrea - Balliu Bertiana - Marchesi Raffaele - Parolin Irene

<b>Introduction</b>	<b>2</b>
<b>Survey design</b>	<b>2</b>
Selection of attributes	2
Design model and generation	3
Distribution of the survey	4
<b>Preprocessing</b>	<b>4</b>
Data cleaning	4
Long format	5
<b>Analysis</b>	<b>5</b>
Descriptive analysis	5
Multinomial Logit model	6
Mixed MNL model	9
Willingness to pay	13
Simulating preference shares	14
Preference share simulation under the fixed effects MNL model	14
Preference share simulation under the Mixed MNL model	16
Individual level variables	17
<b>Conclusion</b>	<b>19</b>
<b>References</b>	<b>21</b>

# Introduction

In microeconomics, measurement of consumers' preferences is one of the most important elements of marketing research. It helps to explain the reasons of consumers' decisions and can lead to better managerial decision making, as a higher situational awareness is present. Using some statistical methods it is possible to quantify preferences and answer the question: what product or service will consumer choose?

To answer this question in a specific use case the present work aims at illustrating an application of Conjoint Analysis methods. The use case considers the product *Christmas box* and the type of specific Conjoint Analysis used is the choice-based.

The methods of conjoint analysis are based on the premise that individuals consider various aspects of a choice alternative. The methods then permit a decomposition of an individual's overall preference judgments about a set of choice alternatives into separate and compatible utility values corresponding to each attribute. These separate functions are called attribute-specific partworth functions. The choice-based conjoint methods (for stated choices) are based on the behavioral theory of random utility maximization. This approach decomposes an individual's random utility for an object into two parts: deterministic utility and a random component. Depending on the distributional assumptions for the random component, a number of alternative models are developed to describe the probability of choice of an object. The most popular one is the multinomial logit model that uses the extreme value distribution for the random term. These methods belong to the family of discrete choice analysis methods.

The business settings of this work consist in a company operating in the market of a specific product - *Christmas box*. The effort goes to the application and description of a path/heuristic on using the choice-based conjoint analysis (CBCA) on data collected through a survey in order to have indications on what product the respondents choose and how to construct a better product in a given simulated market. In the following sections the steps of the work are described and developed. Firstly, a survey was designed with the aim of collecting the needed data. Some data cleaning and preparation has been performed. Following, the analysis has consisted in considering two versions of the CBCA, the fixed effects multinomial (MNL) model and the mixed MNL one. In the next step, through these models, partworths/utilities of individual preferences has been estimated. Based on these estimates of partworths calculation of willingness to pay has been performed to have a more interpretable view. In addition, estimates of partworth has served also the purpose of simulating preference share (through a market simulator/optimizer) and how a better product can be designed. Successively, some experiments have been conducted to explore effects of individual levels on the heterogeneity of respondents. Lastly, some conclusions and some ideas about future development have been drawn.

## Survey design

### Selection of attributes

Let's suppose our firm want to sell a new Christmas box in a supermarket chain with products that we already sell in our store. Observing the market of Christmas boxes, we identified the main

products and price ranges that characterize the majority of them. Selected products, that became our attributes for the conjoint analysis, are the most classic ones to compose a box for Italian consumers, and we chose a bottle of wine and one of the most classic Christmas sweets as the base of box. Then we added products that fall into both the sweet and savory categories, that are a representative subgroup of all the possible alternatives.

These are attributes that we used to compose the questionnaire:

- *Dolce, Bottiglia*: binomial variables with levels *Pandoro* and *Panettone* for the first one and *Spumante* and *Passito* for the second one;
- *Grana, Cioccolatini, Miele, Salame, Torrone*: boolean variables indicating the presence or absence of products;
- *Prezzo*: multinomial integer variable with 4 levels of price 10, 13, 16 and 20 euros.

As it is the case of any statistical model, the more data we collect through the conjoint survey (given a number of attributes), the smaller the standard errors will be. Hence we tried to have the greatest number of respondents in our possibilities. On the other hand, the higher is the number of attributes and levels in a conjoint survey (given a number of respondents), the lesser precise the part worth estimates will be. For this reason we had to choose a limited number of attributes that at the same time constitute a good quantity of products to form a nice gift idea. Another important point was to choose a number that would make the list of attributes readable and usable to respondents.

## Design model and generation

Our choice based conjoint analysis will have to work with 8 variables that have a total of 18 levels that produce 512 unique profiles.

We have decided to use a *fractional design* that let us to work, as the name say, with a fraction of all possible profiles. The advantage is that it is possible to administer the same questionnaire to all the respondents but even more useful is the property of orthogonality: this means that each level of an attribute appears the same number of times within the questionnaire.

Our questionnaire will have 3 alternatives per question. To generate the alternatives we used the *Mix-and-match method*. This method uses the profiles, selected with a fractional factorial design, to generate 3 alternatives for each question. It works with a rotation method, to which is added a further randomization in the choice of profiles to better cover the design space.

At the end of this procedure we obtain a questionnaire with orthogonal design consisting of 16 questions with 3 alternatives each. We believe that in this format the number of questions to ask each respondent is excessive. Since it is a repetitive task with questions that are always similar to each other, the risk that respondents will lose their motivation before completing the answers.

When the experimental design contains too much questions it is possible to apply the *blocking technique*. The questions will be divided into blocks, possibly of equal size. It must be kept in mind that blocks alone do not respect the property of orthogonality, but by bringing together the answers, the complete design is again obtained. We decided to divide the questionnaire into 2 blocks with 8 questions each, so each respondent will have to answer only one block instead of the whole set of questions. This way the cognitive stress will be reduced ensuring a greater efficiency of the respondents.

## Distribution of the survey

To the design we have obtained we have decided to add individual level variables that will allow us to separate our respondents into classes of interest. The demographic questions are:

- *Età*
- *Genere*
- *Occupazione*
- *Provincia di residenza*

Then we have decided to add a question that will be asked at the beginning of the questionnaire. This will generate first of all an additional individual level variable and should help respondents contextualize the choices they are called to make: “*A chi regaleresti un cesto?*”.

Finally, the last question of the questionnaire will ask to directly give a preference for the attribute that we consider most significant in our design: “*Pandoro o Panettone?*”. The idea behind this question is to try to understand if personal taste influences the choice even when asked to choose a product intended for others.

The questionnaire was produced with Google Forms dividing it into two blocks. The two questionnaires were randomly administered to the respondents and distributing them simultaneously we obtained 244 answers for the first block and 264 for the second.

## Preprocessing

### Data cleaning

After the administration of the questionnaires, Google Forms allows to download the results in csv format. To each row corresponds one respondent and to each column corresponds one questionnaire question. At this point, a data cleaning operation was necessary to standardize the data collected.

- Each respondent was assigned a unique *id*;
- Two blocks of questionnaires have been joined into a single data set because taken individually they do not constitute a complete orthogonal design. The number of respondents must be the same for both blocks so we have eliminated excess respondents from block 2. The final dataset, which we will use for analysis, will consist of 474 respondents, 237 per block;
- *Age* and *province* were open questions and it was necessary to standardize the format of the answers: ages have been turned into integer values and provinces into acronyms like “TN” for “Trento” etc.;
- *Gender* corrections: in the beginning we wrongly set choices between *M*, *F* and *Altro* for the question about sex. Then (luckily after only few answers) we changed the question in *Genere* with possible answers *Uomo*, *Donna*, *Preferisco non rispondere* in order to be more respectful;
- *Employment*: similarly to the above gender correction, we were not initially using the male/female version of answers so we had to change *Pensionato* into *Pensionato/a* etc., and then manage them during this phase;
- 9 respondents were eliminated because they did not respond to the gender question for various reasons, therefore being a minority, it would have turned into a huge mistake in estimating the model.

These corrections turned into some work to prepare the answers for the analysis but let us understand that it is better to double check before than going to cover after. In this case the fact that we were distributing the survey through Google form helped a lot because it has not been necessary to send the link again but only revise the problems.

## Long format

In order to proceed with the conjoint analysis, the dataset must have a different format from that of Google Forms: the long format. A dataset in this format has one alternative per row, therefore 24 rows per respondent (3 alternatives per 8 questions). Each row shows the respondent's individual data, an identifier per question and per alternative, and a variable for each attribute of the questionnaire. Finally, a boolean variable that specifies whether the alternative described by the row was chosen by the respondent or not.

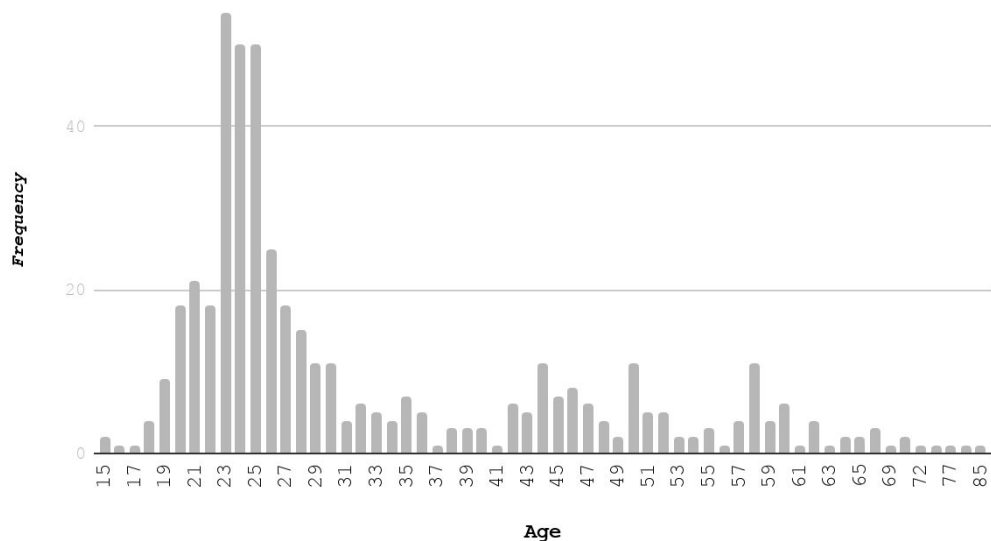
## Analysis

### Descriptive analysis

Before starting the conjoint analysis it is better to understand the properties of the data on which we will work:

- Individual level variables:

- **Età:** Min.: 15                      1st Qu.: 23                      Median: 26  
Max.: 85                      3rd Qu.: 42                      Mean : 32.62



Respondents cover a wide range of ages, mostly around the age of 25.

- **Occupazione:** Lavoratore : 229                      Studente: 192                      Pensionato: 20  
Disoccupato: 5                      Altro : 28

As expected, the employment is concentrated between workers and students but we've had also few unemployed and some retired.

- **Genere:** F: 274                      M: 200

The division is fairly balanced between men and women although there is a slight majority of women surveyed.

- **Provincia:** TN: 115 VI: 115 TV: 42 VR: 37 BO: 20  
BZ: 16 (Other): 129

Also for this attribute the result is as expected because it reflects our friendships and places of origin.

- **Destinatario:** Parente: 197 Famiglia del partner: 132  
Amico : 74 Collega : 71

This is the first unexpected fact as we had no idea what distribution would have been like. At first glance we can say that the interviewees tend to use this type of gift mostly for family members or relatives, and less for friends or colleagues.

- **Preferenza:** Pandoro: 272 Panettone: 202

Here we measure the preference between *Pandoro* and *Panettone* and the results show that there is an overall majority of votes for *Pandoro*, but it is not so remarkable.

- Frequency of the chosen attributes:

<b>Dolce</b>	Pandoro	→	2010,	Panettone	→	1782
<b>Bottiglia</b>	Passito	→	1880,	Spumante	→	1912
<b>Salame</b>	Si	→	2021,	No	→	1771
<b>Torrone</b>	Si	→	2008,	No	→	1784
<b>Cioccolatini</b>	Si	→	2155,	No	→	1637
<b>Miele</b>	Si	→	2178,	No	→	1614
<b>Grana</b>	Si	→	2287,	No	→	1505
<b>Prezzo</b>	10	→	1074,	13	→	1025,
	16	→	976,	20	→	717

These data provide a first insight into the choices made by the respondents. As it was expected, for all boolean variables the most chosen level is that which indicates that the attribute is present in the boxes and this means that respondents generally preferred boxes with multiple products. The price levels also follow the expected trend, with decreasing frequency for an increasing price. However the last level (20 euros) is significantly less chosen than the others. Finally, for the two non-boolean bivariate attributes, the choice of the bottle is more balanced than that of the desert, which shows a preference for pandoro over panettone.

## Multinomial Logit model

In our survey each respondent could choose among three alternatives, so we need to model a dependent variable that is a qualitative multinomial variable with three levels. The suitable model for a problem formulated in this way is the Multinomial Logit Model. It is a discrete choice model, that expand the binomial logit model to the case where the dependent variable has more than two levels. Let  $y_i$  denote the response variable taking on the finite numbers of mutually exclusive values 1, 2, 3 and let  $x_i$  denote the set of K alternative-specific attributes. The probability that respondent  $i$  chooses alternative  $j$  is given by:

$$p_{ij} = \frac{\exp(x_{ij}\beta)}{\sum_h \exp(x_{ih}\beta)}$$

With the MNL model, we can get a precise measurement of how much each attribute is associated to respondents' choices. To perform this estimate we have fitted our model in different ways: firstly we performed an estimate not only with all the attributes, but also with and an intercept for each alternative.

Coefficients :	Estimate	Std. Error	z-value	Pr(> z )	
2: (intercept)	0.217	0.042	5.175	2.3e-07	***
3: (intercept)	0.094	0.044	2.153	0.031	*
panettone	-0.149	0.048	-3.144	0.002	**
spumante	0.036	0.049	0.738	0.460	
salame	0.201	0.046	4.386	1.2e-05	***
torrone	0.242	0.047	5.130	2.9e-07	***
cioccolatini	0.475	0.044	10.872	< 2.2e-16	***
miele	0.393	0.041	9.597	< 2.2e-16	***
grana	0.582	0.045	13.072	< 2.2e-16	***
€ 13	0.045	0.065	0.695	0.487	
€ 16	-0.109	0.067	-1.642	0.101	
€ 20	-0.502	0.064	-7.904	2.7e-15	***

(baseline for binary variables is FALSE)

The *Estimate* column provides the estimated average part worths for each level; they have to be interpreted in respect to the reference level of each attribute. Significant statistical values for this model are *prezzo20* and all boolean variables. It could be trivial that the presence or the absence of a product in our boxes plays an important role and gives us significant values, but we can already see that not all boolean variables have the same weight in the model. Surprisingly, the interception related to the second alternative is also significant. Probably because, even if only slightly, it is the alternative that has been chosen several times by the respondents (37% of cases). Secondly, we tried to fit the model with a fixed intercept, in this way we did not consider alternative position like an attribute.

Coefficients:	Estimate	Std. Error	z-value	Pr(> z )	
panettone	-0.145	0.047	-3.064	0.002	**
spumante	0.044	0.048	0.907	0.364	
salame	0.217	0.046	4.705	2.5e-06	***
torrone	0.247	0.047	5.261	1.4e-07	***
cioccolatini	0.473	0.043	10.903	< 2.2e-16	***
miele	0.401	0.041	9.852	< 2.2e-16	***
grana	0.571	0.044	12.899	< 2.2e-16	***
€ 13	0.033	0.065	0.507	0.612	
€ 16	-0.113	0.066	-1.709	0.087	.
€ 20	-0.506	0.064	-7.943	1.9e-15	***

(baseline for binary variables is FALSE)

In our framework, the position of alternatives is random and did not have any specific value. Despite this, with the model that estimates different coefficients for each alternative, a better goodness of fit is obtained. We can see the Likelihood Ratio test: with this method we compare the log-likelihood of two models to see if one has significantly better goodness of fit than the other (the higher log-likelihood, the better).

#### Likelihood ratio test:

Model 1: scelta ~		dolce + bottiglia	+ salame + torrone + cioccolatini +	
		miele + grana	+ prezzo   -1	
Model 2: scelta ~		dolce + bottiglia	+ salame + torrone + cioccolatini +	
		miele + grana	+ prezzo	
#Df	LogLik	Df	Chisq	Pr(>Chisq)
1	10	-3850.5		
2	12	-3836.8	2	27.386 1.13e-06 ***

We know that the position of alternatives has no relevant value in our survey. So, the model with which we will make our analyses will be the one without intercepts. However, we wanted to make

sure that the significance of the intercepts is due to a favorable random distribution of alternatives. To verify this, we checked that none of the respondents had responded to the survey at random, choosing the same alternative every time. Fortunately this case did not occur.

Finally, a third MNL model was created, in which the price variable was considered numerical, instead of categorical. In doing so, a single coefficient for price is estimated. This new model significantly worsens the goodness of fit. From the distribution of the choice for each price level, seen during the descriptive analysis, it is possible to observe the choices do not follow a linear trend, but that the higher level is clearly preferred less than the others. For this reason, considering separately the price levels, it is possible to achieve a better goodness of fit.

#### Likelihood ratio test:

```
Model 1:  scelta ~ dolce + bottiglia + salame + torrone + cioccolatini
          + miele + grana + as.numeric(as.character(prezzo)) | -1
Model 2:  scelta ~ dolce + bottiglia + salame + torrone + cioccolatini
          + miele + grana + prezzo | -1
```

	#Df	LogLik	Df	Chisq	Pr(>Chisq)
1	8	-3859.8			
2	10	-3850.5	2	18.595	9.165e-05 ***

## Mixed MNL model

With an MNL model, the coefficients of the attributes are estimated as the mean of the utilities of the respondents. In order to take into account the heterogeneity of the respondents we use a Mixed MNL model. The mixed model respondent-level coefficients are called random effects or random coefficients, because they are calculated as random variables with an a priori defined density function. These coefficients can be obtained thanks to repeated observations for each respondent, typical of conjoint surveys. Our respondents may have different preferences, so we can expect that a model capable of capturing consumer-level coefficients should have better goodness of fit than the MNL model. Unlike the previous MNL model, in this case the coefficients  $\beta$  are not fixed, but vary across respondents. Their variation is described by a density function  $f(\beta, \theta)$ . For our analysis we assume that all coefficients follow a Normal distribution, and that Teta includes means, variances and covariances parameters. The unconditional choice probability  $p_{ij}$  that a respondent  $i$  will choose the alternative  $j$  is obtained by integrating the conditional probability to the value of  $\beta_i$  for all possible values of  $\beta_i$ .

$$p_{ij} = \int_{\beta} (p_{ij}|\beta) f(\beta, \theta) d\beta = \int \frac{\exp(x_{ij}\beta_i)}{\sum_h \exp(x_{ih}\beta_i)} f(\beta, \theta) d\beta$$



Estimating a mixed MNL model with covariances among random effects set to zero, we get the following results:

Coefficients:	Estimate	Std. Error	z-value	Pr(> z )	
panettone	-0.224	0.063	-3.558	0.001	***
spumante	0.069	0.061	1.131	0.258	
salame	0.271	0.057	4.737	2.2e-06	***
torrone	0.374	0.061	6.173	6.7e-10	***
cioccolatini	0.538	0.055	9.774	< 2.2e-16	***
miele	0.551	0.054	10.241	< 2.2e-16	***
grana	0.756	0.057	13.202	< 2.2e-16	***
€ 13	0.029	0.080	0.359	0.719	
€ 16	-0.253	0.082	-3.096	0.002	**
€ 20	-0.726	0.083	-8.743	< 2.2e-16	***
sd. panettone	1.371	0.092	14.830	< 2.2e-16	***
sd. spumante	0.629	0.098	6.446	1.2e-10	***
sd. salame	0.948	0.092	10.322	< 2.2e-16	***
sd. torrone	-0.451	0.107	-4.237	2.3e-05	***
sd. cioccolatini	-0.160	0.118	-1.356	0.175	
sd. miele	0.595	0.090	6.598	4.2e-11	***
sd. grana	0.523	0.103	5.087	3.6e-07	***
sd. € 13	-0.072	0.159	-0.455	0.649	
sd. € 16	-0.389	0.143	-2.718	0.007	**
sd. € 20	0.721	0.114	6.336	2.4e-10	***

(baseline for binary variables is FALSE)

#### Random coefficients:

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
panettone	-Inf	-1.148	-0.224	-0.224	0.701	Inf
spumante	-Inf	-0.355	0.069	0.069	0.494	Inf
salame	-Inf	-0.369	0.271	0.271	0.910	Inf
torrone	-Inf	0.069	0.374	0.374	0.678	Inf
cioccolatini	-Inf	0.430	0.538	0.538	0.646	Inf
miele	-Inf	0.150	0.551	0.551	0.952	Inf
grana	-Inf	0.403	0.756	0.756	1.109	Inf
€ 13	-Inf	-0.020	0.029	0.029	0.078	Inf
€ 16	-Inf	-0.515	-0.253	-0.253	0.010	Inf
€ 20	-Inf	-1.212	-0.726	-0.726	-0.239	Inf

(baseline for binary variables is FALSE)

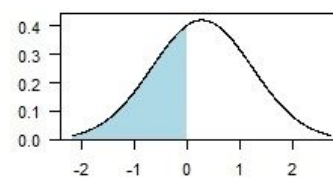
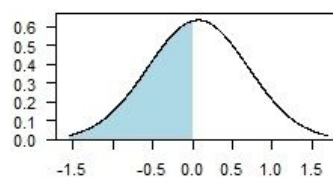
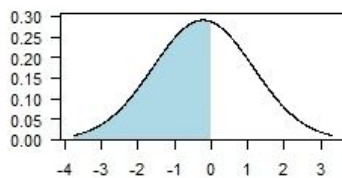
In addition to the estimate of the average part worth coefficients, an estimate of their standard deviations is calculated. These new coefficients describe variability across respondents. The table with the summary measures for the random coefficients shows in further detail the degree of separation from the average of each individual level coefficient. When the first and third quartiles have values of opposite sign, a significant part of the respondents is not correctly represented by the average value of the fixed effects.

For our analysis, the most striking data is that of the *dolce* attribute; with the MNL model it could seem a variable that was not very relevant for the respondents, while the random effects show that it is the attribute with the greatest variance. This means that for the respondents the presence of panettone or pandoro significantly influences the choice of the alternative, but they are divided in an almost balanced way in preferring one or the other, with a slight majority preferring *pandoro*. This separation is also present for the *bottiglia* attribute, but with a reduced standard deviation. It means that respondents are less influenced by the two levels of this variable. Boolean attributes all

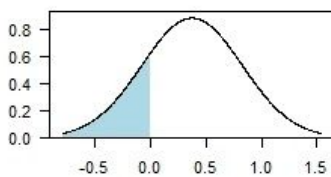
have more consistent random effects across the population. The only exception is *salame*, with a negative first quartile; it could be the manifestation of a minority of vegetarian respondents who prefer to avoid this product. Finally, regarding the price, the random effect *price13* is distributed with little variance around the intercept represented by a price of 10 euros. For our respondents this price change is generally not very relevant.

The following graphs show the distribution of random effects, specifying in percentage how the respondent level coefficients are divided with respect to the average. It can be seen that some variables, such as *dolce* and *bottiglia* separate almost equally the tastes of the respondents. On the contrary, everyone seems to have a positive preference for the variable *cioccolatini*, which manages to represent the whole sample.

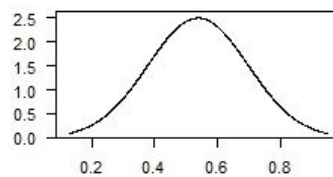
Distribution of dolcepanettone : 56 % of 0    Distribution of bottigliaspumante : 46 % of 0    Distribution of salameTRUE : 39 % of 0



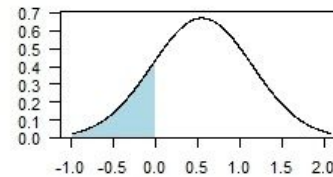
Distribution of torroneTRUE : 20 % of 0



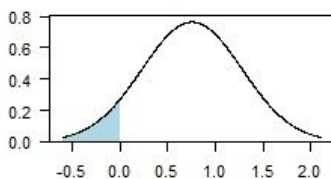
Distribution of cioccolatiniTRUE : 0 % of 0



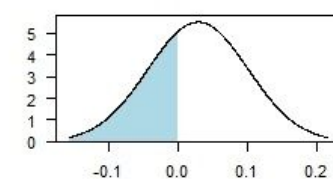
Distribution of mieleTRUE : 18 % of 0



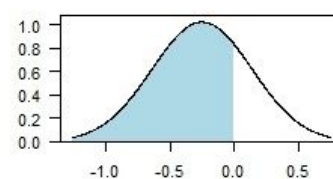
Distribution of granaTRUE : 7 % of 0



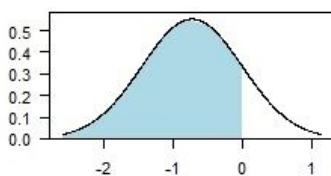
Distribution of prezzo13 : 35 % of 0



Distribution of prezzo16 : 74 % of 0



Distribution of prezzo20 : 84 % of 0



The mixed MNL model can also be estimated by admitting correlation between random effects. In addition to mean and variance, the covariances between the coefficients are obtained. With these new coefficients we have built the following correlation matrix:

	panettone	spumante	salame	torrone	cioccolat.	miele	grana	€ 13	€ 16	€ 20
panettone	1.000	-0.488	0.111	0.004	0.179	-0.107	-0.031	-0.230	0.306	0.217
spumante	-0.488	1.000	0.052	-0.207	-0.771	-0.147	0.081	0.243	-0.341	-0.234
salame	0.111	0.052	1.000	0.346	0.074	0.153	0.329	0.095	-0.248	0.064
torrone	0.004	-0.207	0.346	1.000	0.745	-0.472	0.201	-0.002	-0.237	0.004
cioccolat.	0.179	-0.771	0.074	0.745	1.000	-0.104	0.144	-0.199	0.046	0.099
miele	-0.107	-0.147	0.153	-0.472	-0.104	1.000	0.556	-0.248	-0.062	-0.156
grana	-0.031	0.081	0.329	0.201	0.144	0.556	1.000	-0.291	-0.443	-0.365
€ 13	-0.230	0.243	0.095	-0.002	-0.199	-0.248	-0.291	1.000	0.600	0.762
€ 16	0.306	-0.341	-0.248	-0.237	0.046	-0.062	-0.443	0.600	1.000	0.907
€ 20	0.217	-0.234	0.064	0.004	0.099	-0.156	-0.365	0.762	0.907	1.000

(baseline for binary variables is FALSE)

Some correlations were expected: it is usual to find them between the levels of the same attribute, so we are not surprised by the high values between the different prices. In addition to these, less obvious relationships also appear. For example that of 0.74 between *cioccolatini* and *torrone* indicates that respondents treat them similarly, and lovers of sweets will want to have both. A similar but less pronounced correlation (0.56) is observed for the pair *grana* and *miele*. Instead it is less easy to interpret the negative correlation between *cioccolatini* and *bottigliaSpumante* of -0.77. With a likelihood ratio test we evaluate the goodness of fit of the uncorrelated mixed MNL model, compared to that of the MNL model of the previous section.

#### Likelihood ratio test:

```
Model 1: scelta ~ dolce + bottiglia + salame + torrone + cioccolatini + miele
          + grana + prezzo | -1
Model 2: scelta ~ dolce + bottiglia + salame + torrone + cioccolatini + miele
          + grana + prezzo | -1
```

	#Df	LogLik	Df	Chisq	Pr(>Chisq)
1	10	-3850.5			
2	20	-3660.3	10	380.36	< 2.2e-16 ***

As expected, the new model significantly improves the goodness of fit. We can also compute the test between the uncorrelated model and fully correlated one.

#### Likelihood ratio test:

```
Model 1: scelta ~ dolce + bottiglia + salame + torrone + cioccolatini + miele
          + grana + prezzo | -1
Model 2: scelta ~ dolce + bottiglia + salame + torrone + cioccolatini + miele
          + grana + prezzo | -1
```

	#Df	LogLik	Df	Chisq	Pr(>Chisq)
1	20	-3660.3			
2	65	-3563.0	45	194.74	< 2.2e-16 ***

The fully correlated random effects are those that produce the model with the best goodness of fit, although an increase in the number of degrees of freedom is paid.

## Willingness to pay

In order to have a better interpretation of part-worths provided by the MNL model, we use the Willingness to pay (WTP) that “is the maximum price at or below which a consumer will definitely buy one unit of a product.”<sup>1</sup>

“The approaches to measure consumer WTP can be differentiated whether they measure WTP directly or indirectly and whether they measure consumer hypothetical or actual WTP.”<sup>2</sup>

With the direct approach consumers are asked to directly state their WTP for a specific product through an open-ended question. On the opposite, with the indirect approach (such as our choice-based conjoint analysis) WTP is calculated on the basis of consumers' choices among several product alternatives and a “none” choice option.

$$WTP_{\text{level of an attribute}} = \frac{\beta_{\text{level of an attribute}}}{\beta_{\text{price}}}$$

However, many studies have shown that both direct and indirect approaches can generate inaccurate results for various psychological and technical reasons. More fundamentally, both approaches measure consumers' hypothetical, rather than actual, WTP and thus can generate hypothetical bias, which the economics literature defines as the bias induced by the hypothetical nature of a task.

Our results, computed with the choice-based conjoint analysis approach, are shown in the table below, ordered by absolute value. The WTP was obtained from fixed effects by an MNL model that uses the price as a numerical variable.

Product	Willingness to pay (€)
Grana	-10.33
Ciocolatini	- 8.46
Miele	- 7.28
Torrone	- 4.94
Salame	- 3.99
Panettone (in respect to Pandoro)	+ 2.35
Spumante (in respect to Passito)	- 0.93

The above values endorse that there are some products respondents will like to prefer in the box such as Grana, Ciocolatini and Miele. Since our variables represent the presence of a product, results are mainly negative because interviewees would be divided equally between the two choices of buying a box with Grana or spending € 10,33 less for a box without Grana.

If, on the other hand, we carry out this analysis with the fully correlated random effects of the mixed MNL model, the WTP values we obtain are also distributed according to Normal distributions (in this case the sign is reverse).

	1st Qu.	Mean	3rd Qu.
panettone	-11.39	-1.70	7.99
spumante	-4.26	0.85	5.95
salame	-3.69	3.31	10.30
ciocolatini	3.06	5.11	7.16
torrone	1.06	4.78	8.50
miele	1.37	6.22	11.06
grana	2.98	8.46	13.93
(baseline for binary variables is FALSE)			

<sup>1</sup> *Microeconomic Analysis*

<sup>2</sup> *How Should Consumers' Willingness to Pay be Measured? An Empirical Comparison of State-of-the-Art Approaches*

As with the part worth coefficients, the wide variability of the effect associated with *dolce* can be seen. Compared to the previous table, the mean values of the coefficients are slightly lower, because considering also the standard deviations and the correlations, the model reduces the utility associated with the single attributes.

## Simulating preference shares

Our conjoint study led to a set of utilities or part-worths that quantify respondents' preferences for each level of each attribute. An important tool to assess the role product attributes consists on using the model to obtain preference share predictions. In other words, simulating the preference shares (with a market simulator). The simulator converts raw conjoint (part-worth utility) data into something much more managerially useful: simulated market preferences/choices <sup>3</sup>.

So, we can use the estimated model to predict preference share for our firm's new option of Christmas box against a set of other options offered by the main competitors. By changing the levels' attributes of the planned option, we can see how variations in the option influence the preference share.

### Preference share simulation under the fixed effects MNL model

Let's suppose that the company had in mind to market a new option of Christmas box with *panettone*, *spumante*, *torrone*, *miele e grana* for a *price* of 16 euros. The estimated MNL model can be used to predict the preference share for this option against other options that represent properly the market. In the table below we can see the preference share predictions after simulating the market.

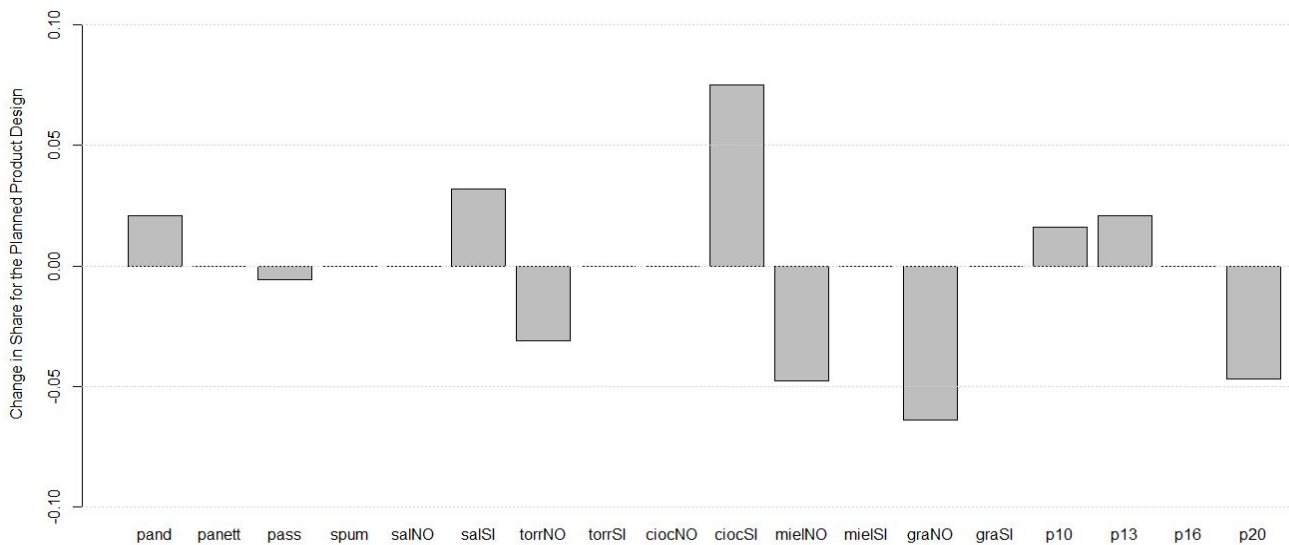
	share	2.5%	97.5%	dolce	bottiglia	salame	torrone	cioccolat.	miele	grana	prezzo
364	0.163	0.145	0.192	panettone	spumante	FALSE	TRUE	FALSE	TRUE	TRUE	16
7	0.077	0.066	0.091	pandoro	spumante	TRUE	FALSE	FALSE	FALSE	FALSE	10
209	0.175	0.149	0.207	pandoro	passito	FALSE	FALSE	TRUE	FALSE	TRUE	13
226	0.141	0.125	0.162	panettone	passito	FALSE	FALSE	FALSE	TRUE	TRUE	13
302	0.109	0.094	0.125	panettone	passito	TRUE	TRUE	FALSE	TRUE	FALSE	16
343	0.196	0.173	0.213	pandoro	spumante	TRUE	FALSE	TRUE	FALSE	TRUE	16
445	0.137	0.116	0.152	pandoro	passito	TRUE	TRUE	TRUE	TRUE	FALSE	20

The first row refers to our company's planned design, while the other six rows represent the simulated/potential competitors' options of Christmas boxes. The column share provides the model predicted preference shares. We find that, among this set of Christmas boxes, we would expect consumers to choose our product basically 16% of times. Our company could use the estimated model to assess how modifying the attributes of the planned product would affect the preference shares. We are aware of the fact that these predicted shares are made relative to this specific given set of potential competitors, thus implying that the predicted preference share for our planned christmas would change if this set was different.

To get a feeling for the degree of uncertainty due to sampling and measurement error associated with a given share of preference we estimated a confidence interval. We used the bootstrap option to have a robust calculation of confidence interval. From the table above we can see that the 95% confidence intervals for our simulated market are quite narrow and consequently our simulation quite accurate.

<sup>3</sup> Orme, B. "Getting started with conjoint analysis: strategies for product design and pricing research second edition." *Madison: Research Publishers LLC* (2010).

### Sensitivity chart:



Using the estimated MNL model we were able to predict how the preference shares would change if variations on the levels of the attributes were considered, without changing the market (the set of competing products). The figure above shows the plot of the expected changes in preference shares because of changes in each of the attributes of our planned design (*panettone*, *spumante*, *torrone*, *miele*, *grana* and a price of 16€), one at a time. From our sensitivity chart for our planned product we can say that changing, ceteris paribus, the design:

- from *panettone* to *pandoro* would increase the preference share by almost 0.025
- from *spumante* to *passito* would slightly decrease the preference share, by almost 0.01
- from not having *salame* to having *salame* would increase the preference share by almost 0.03
- from having *torrone* to not having it would decrease the preference share by almost 0.03
- from not having *cioccolato* to having it would increase the preference share by almost 0.08
- from having *miele* to not having it would decrease the preference share by almost 0.05
- from having *grana* to not having it would decrease the preference share by almost 0.07
- from price 16 euros to lower prices 10euros and 13euros it would increase the preference share by respectively almost 0.02 and 0.025. While changing the price from 16 euros to higher price of 20 euros it would decrease the preference share by 0.05.

Considering the results of the sensitivity chart we constructed the product with the preferences of the consumers. We took care to not change the attributes with high decreasing effect on preference and change attributes with high increasing effect on the preferences. We kept the attributes *miele* and *grana* as not having them would decrease importantly the preference for our product. We added the attribute *cioccolatini* as it has the highest increase in preference amongst the boolean attributes not included. We took care to keep a reasonable price for the modified option of the product. In the table below we can see the market simulation after substituting the initial product with that constructed.

	share	2.5%	97.5%	dolce	bottiglia	salame	torrone	cioccolat.	miele	grana	prezzo
371	0.220	0.201	0.241	pandoro	spumante	FALSE	FALSE	TRUE	TRUE	TRUE	16
7	0.072	0.060	0.084	pandoro	spumante	TRUE	FALSE	FALSE	FALSE	FALSE	10
209	0.163	0.143	0.184	pandoro	passito	FALSE	FALSE	TRUE	FALSE	TRUE	13
226	0.131	0.116	0.151	panettone	passito	FALSE	FALSE	FALSE	TRUE	TRUE	13
302	0.102	0.087	0.117	panettone	passito	TRUE	TRUE	FALSE	TRUE	FALSE	16
343	0.183	0.163	0.203	pandoro	spumante	TRUE	FALSE	TRUE	FALSE	TRUE	16
445	0.128	0.112	0.143	pandoro	passito	TRUE	TRUE	TRUE	TRUE	FALSE	20

From the market simulation with the new Christmas box, we can note that the preference share increased with almost 0.06.

## Preference share simulation under the mixed MNL model

Using the MNL model to predict and simulate the preference share we simplify the behaviour of consumers because we consider their preferences as their average value. In the following tables we report the results of the preference share simulation for the selected product (the same as in the fixed effects MNL model) in the simulated/potential competitors' set of products with the mixed MNL model that considers the heterogeneity of consumers. In the first table we can see that the preference share of our designed product is slightly higher than in the case of fixed. This can be interpreted with the fact that in the mixed MNL model are considered even the consumers preferring "niche" products where *panettone* is present. So, we can infer that our initial design of the product having *panettone* present is preferred from a fraction on consumers which were not represented in the fixed effects MNL model simulation.

	colMeans(shares)	dolce	bottiglia	salame	torrone	cioccolat.	miele	grana	prezzo
364	0.171	panettone	spumante	FALSE	TRUE	FALSE	TRUE	TRUE	16
7	0.092	pandoro	spumante	TRUE	FALSE	FALSE	FALSE	FALSE	10
209	0.194	pandoro	passito	FALSE	FALSE	TRUE	FALSE	TRUE	13
226	0.144	panettone	passito	FALSE	FALSE	FALSE	TRUE	TRUE	13
302	0.123	panettone	passito	TRUE	TRUE	FALSE	TRUE	FALSE	16
343	0.135	pandoro	spumante	TRUE	FALSE	TRUE	FALSE	TRUE	16
445	0.140	pandoro	passito	TRUE	TRUE	TRUE	TRUE	FALSE	20

We considered the sensitivity chart for our designed product, and in the mixed MNL it produced the same results as in the fixed effects version of the model. So, the improved product results the same as in the fixed effects MNL model. After considering this product to simulate the preference share in the same set of competitors' products, we can see in the table below that our modified product have a preference share of around 0.20, which is lower than that of the fixed effects model. This is explainable with the fact that considering the heterogeneity of consumers, niche products preferences of the small part different from the average one (which is considered in the fixed effects) becomes visible to the model and the predictions about the preference share.

	colMeans(shares)	dolce	bottiglia	salame	torrone	cioccolat.	miele	grana	prezzo
371	0.197	pandoro	spumante	FALSE	FALSE	TRUE	TRUE	TRUE	16
7	0.086	pandoro	spumante	TRUE	FALSE	FALSE	FALSE	FALSE	10
209	0.169	pandoro	passito	FALSE	FALSE	TRUE	FALSE	TRUE	13
226	0.154	panettone	passito	FALSE	FALSE	FALSE	TRUE	TRUE	13
302	0.144	panettone	passito	TRUE	TRUE	FALSE	TRUE	FALSE	16
343	0.120	pandoro	spumante	TRUE	FALSE	TRUE	FALSE	TRUE	16
445	0.130	pandoro	passito	TRUE	TRUE	TRUE	TRUE	FALSE	20

Both models, fixed effects and mixed one are used with the caution to not consider identical or very similar options of products while predicting the preference share. This, in order to prevent (at least partially) the violation of the IIA property.

## Individual level variables

With our survey, in addition to the choices of alternatives, we also collected various data on an individual level on our respondents. It is possible that individuals with different demographic characteristics may have different preferences. To verify this, we introduce respondent information as individual level variables in our MNL model. An individual level variable interacts with all those at the choice level, so the model being calculated will have to estimate many more coefficients. The variables that we have included in the model are gender, age divided into four classes, occupation, and recipient of the christmas box. For each of these we have created a different model, to understand with which interaction the best goodness of fit is obtained. Below we see the likelihood ratio tests for each new model, using the MNL model as baseline.

### Likelihood ratio test:

```
Model 1:  scelta ~ dolce + bottiglia + salame + torrone + cioccolatini + miele
           + grana + prezzo | -1
Model 2:  scelta ~ dolce + bottiglia + salame + torrone + cioccolatini + miele
           + grana + prezzo | dolce * genere + bottiglia * genere
           + salame * genere + torrone * genere + cioccolatini * genere
           + miele * genere + grana * genere + prezzo * genere
```

	#Df	LogLik	Df	Chisq	Pr(>Chisq)
1	10	-3850.5			
2	63	-3769.7	53	161.65	6.396e-13 ***

### Likelihood ratio test:

```
Model 1:  scelta ~ dolce + bottiglia + salame + torrone + cioccolatini + miele
           + grana + prezzo | -1
Model 2:  scelta ~ dolce + bottiglia + salame + torrone + cioccolatini + miele
           + grana + prezzo | dolce * classeEta + bottiglia * classeEta
           + salame * classeEta + torrone * classeEta
           + cioccolatini * classeEta + miele * classeEta
           + grana * classeEta + prezzo * classeEta
```

	#Df	LogLik	Df	Chisq	Pr(>Chisq)
1	10	-3850.5			
2	125	-3720.1	115	260.83	2.461e-13 ***

### Likelihood ratio test:

```
Model 1:  scelta ~ dolce + bottiglia + salame + torrone + cioccolatini + miele
           + grana + prezzo | -1
Model 2:  scelta ~ dolce + bottiglia + salame + torrone + cioccolatini + miele
           + grana + prezzo | dolce * occupazione
           + bottiglia * occupazione + salame * occupazione
           + torrone * occupazione + cioccolatini * occupazione
           + miele * occupazione + grana * occupazione
           + prezzo * occupazione
```

	#Df	LogLik	Df	Chisq	Pr(>Chisq)
1	10	-3850.5			
2	125	-3722.8	115	255.32	1.177e-12 ***



#### Likelihood ratio test:

Model 1: scelta ~ dolce + bottiglia + salame + torrone + cioccolatini + miele  
+ grana + prezzo | -1

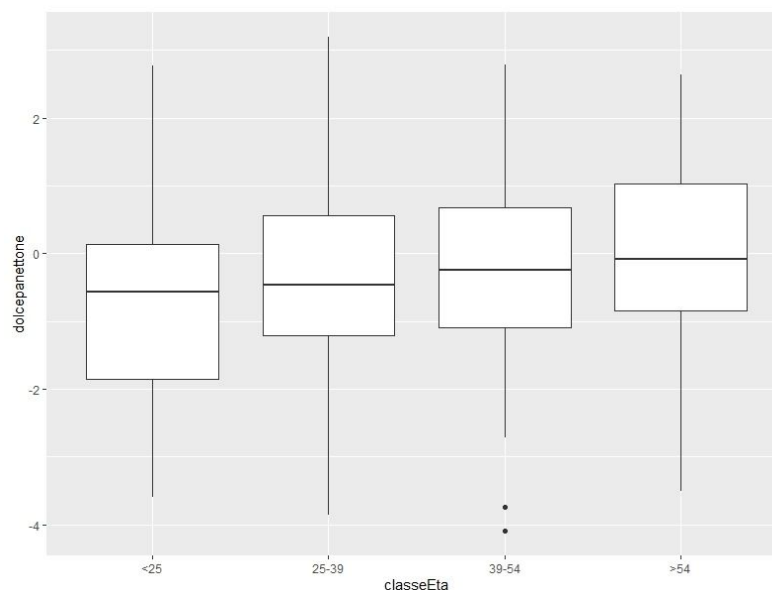
Model 2: scelta ~ dolce + bottiglia + salame + torrone + cioccolatini + miele  
+ grana + prezzo | dolce \* destinatario  
+ bottiglia \* destinatario + salame \* destinatario  
+ torrone \* destinatario + cioccolatini \* destinatario  
+ miele \* destinatario + grana \* destinatario  
+ prezzo \* destinatario

	#Df	LogLik	Df	Chisq	Pr(>Chisq)
1	10	-3850.5			
2	125	-3729.4	115	242.17	4.437e-11 ***

All individual level variables significantly improve the goodness of fit. The age groups variable produces the best result, while with gender the worst one is obtained.

It is important to underline one detail: although there is an improvement in the goodness of fit, the price to pay is a considerable increase in the number of degrees of freedom. In fact, the new models also estimate a coefficient for each interaction between each level of the individual level variable and the choice attributes. Furthermore, it is not possible to exclude the intercepts associated with the position of the alternatives. If we estimated a model with the contemporary interactions of all four individual level variables, the goodness of fit would still improve; but the number of degrees of freedom would increase to 279. The problem with this approach is therefore that it is difficult to interpret the coefficients in detail, given their very high number. With this analysis, however, we understood that the use of individual level variables can help us formalize the heterogeneity of respondents, and that age groups should be the most significant variable. While gender is the least able to capture the differences within our sample.

To better evaluate whether the heterogeneity of the consumer can be explained by his individual characteristics, we use the random coefficients of the mixed MNL model. We associate each individual level part worth with the corresponding value of the individual level variable. As we have seen, the attribute that varies the most across respondents is "sweet". So we analyzed its relationship with the individual level variables to understand if the latter are able to explain the heterogeneity of the respondents. The result is that individual level variables allow identifying characteristics that separate respondents based on their preferences. For example, with age it is possible to identify a pattern that shows how the preference for pandoro decreases with increasing age, this preference is instead very marked in younger people.



<b>dolcepanettone:</b>	classeEta: <25	classeEta: 25-39	classeEta: 39-54	classeEta: >54
	-0.664	-0.312	-0.246	-0.012

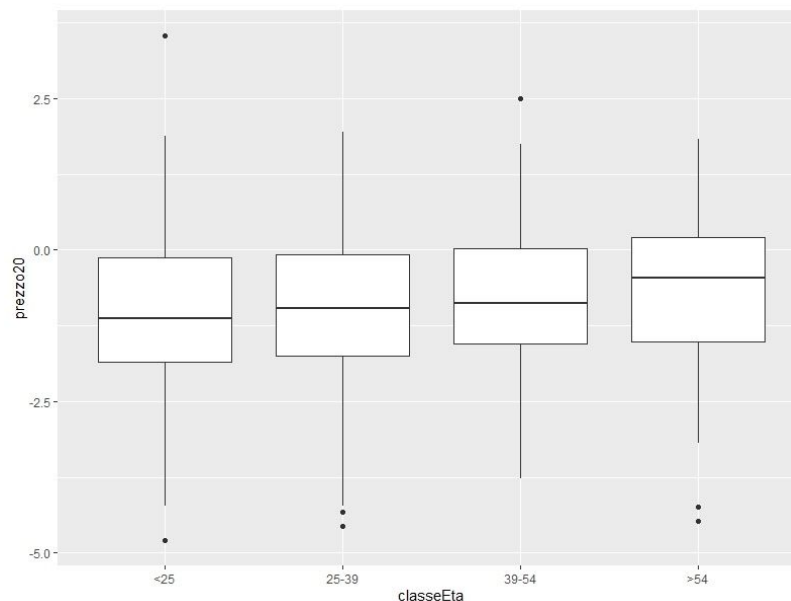
The same pattern can be identified with the occupation variable:

<b>dolcepanettone:</b>	occupazione: Disoccupato/Altro	occupazione: Lavoratore
	-0.130	-0.344
	occupazione: Pensionato	occupazione: Studente
	0.233	-0.583

While with gender the separation is much less pronounced:

<b>dolcepanettone:</b>	genere: F	genere: M
	-0.420	-0.376

Age also has an interesting relationship with the utility that respondents associate with price. If we take into consideration the random effect that represents the highest price, *price20*, we discover that the youngest respondents are those who give the most negative utility. With increasing age the higher price range becomes less rejecting.



<b>prezzo20:</b>	classeEta: <25	classeEta: 25-39	classeEta: 39-54	classeEta: >54
	-1.051	-0.965	-0.818	-0.699

Our analysis therefore shows that the age of customers would be the most useful data for those who decide to produce a new product to be offered for sale. If you had access to the average age of customers, this information could prove useful for placing on the market products aimed at specific age groups.

## Conclusion

In this work we presented a choice-based conjoint analysis for the *Christmas boxes* market. We first identified the most important attributes that define a Christmas box. The variables chosen were *dolce* (pandoro or panettone), *bottiglia* (spumante or passito), *salame*, *torrone*, *cioccolatini*, *miele*, *grana*, and the price (10, 13, 16, 20 euros). We built a choice based survey with fractional factorial design, choosing the attributes with a mix and match method, and dividing the set of questions into two blocks of 8. We conducted the experiment, collecting the answers of 508 participants.

After a data cleaning and preparation phase, the choice-based conjoint analysis was performed on 474 respondents. We calculated the part worths with different Multinomial Logit Models, and then we refined the estimate with the random effects of the Mixed MNL Model. With the Willingness To Pay we have given a more tangible interpretation of the estimated partworth, normalizing them for the price. The result was that not all boolean variables have the same validity for the respondents; for example *grana* has more than double the importance of the choice compared to *salame*. The most interesting variable turned out to be the alternative between *Pandoro* and *Panettone*. Our respondents split almost equally between the two options, with an advantage for the *pandoro*; the same result we obtained by asking to openly provide the preference in the questionnaire. This variable is also the one with the greatest variance, meaning that it is decisive in choosing one basket over the other, regardless of preference (*panettone* or *pandoro*).

We then estimated the preferences of our respondents by simulating a market of competitive products, showing how the information provided by our survey could provide a competitive advantage in identifying a design capable of earning most of the preference shares. In this phase of the analysis we were also able to show how the mixed MNL model is able to identify and represent market niches in the estimates.

Finally, we tried to use the individual level variables collected with the questionnaire to explain the heterogeneity of the respondents. The variable that most of all helped us to group respondents' behaviors was age. We were able to identify different preferences based on the age class both in the choice between *pandoro* and *panettone*, and in the degree of aversion for the higher price level. This suggests that knowing the age of consumers could further assist in choosing the most competitive design.

Interesting areas for exploration and/or extension of this work would be evaluating the survey designing part, using options of more robust of modelling with respect to the IIA assumption, and, while having real sales data (so, knowledge about the market share), considering the CBCA a tool for a better managerial decision making.

Regarding the evaluation of survey: as data quality is crucial, using an evaluation system to assess the survey design which collects the data is important. A system evaluation is the one that considers four criteria for evaluating choice set designs: (1) statistical efficiency; (2) minimal overlap; (3) utility balance; and (4) managerial efficiency.

Regarding alternatives of a/some more robust techniques of conjoint analysis with respect to the IIA assumption, options as Multinomial Probit Model or/and Nested Logit Model would be interesting to apply.

Regarding the difference between preference share and market share: while comparing preference share with market share, if our product has a preference share that is lower than its respective market share our product might be on a vulnerable position share because distribution or awareness of our products is particularly strong, or a market shift is occurring and other products have not yet reached their potential. When preference share does not equal market share, the conjoint model won't necessarily tell us whether distribution, awareness, loyalty, or other issues are at work. But, with an important information about whether our product is in a vulnerable position (preference is less than market share) or whether we have an opportunity (preference is greater than market share), even before any product or price changes are made. This way, being more aware of the situation and our position, better decisions could be made.

# References

1. Microeconomic Analysis, Vol. 3 - Varian, Hal R., 1992
2. How Should Consumers' Willingness to Pay be Measured? An Empirical Comparison of State-of-the-Art Approaches, Journal of Marketing Research - Miller, Klaus M., Hofstetter, Reto, Krohmer, Harley, Zhang, John Z., 2011
3. Orme, B. (2010, 2019) Getting Started with Conjoint Analysis: Strategies for Product Design and Pricing Research. Fourth Edition, Madison, Wis.: Research Publishers LLC
4. <https://cran.r-project.org/package=mlogit>
5. Train, Kenneth E. 2009. Discrete Choice Methods with Simulation. 2nd ed. Cambridge University Press.
6. McFadden, Daniel, and Kenneth Train. 2000. "Mixed Mnl Models for Discrete Response." Journal of Applied Econometrics 15 (5): 447–70.
7. Enneking, Ulrich & Neumann, Claudia & Henneberg, Sven. (2007). How important intrinsic and extrinsic product attributes affect purchase decision. Food Quality and Preference. 18. 133-138
8. Raghavarao, D. & Wiley, J.B. & Chitturi, P.. (2010). Choice-based conjoint analysis: Models and designs.
9. Lebeau, Kenneth & Van Mierlo, Joeri & Lebeau, Philippe & Mairesse, Olivier & Macharis, Cathy. (2012). A choice-based conjoint analysis on the market potential of PHEVs and BEVs in Flanders. World Electric Vehicle Journal. 5. 871-880
10. Rao, Vithala. (2013). Applied Conjoint Analysis.
11. Johnson, F. & Lancsar, Emily & Marshall, Deborah & Kilambi, Vikram & Mühlbacher, Axel & Regier, Dean & Bresnahan, Brian & Kanninen, Barbara & Bridges, John. (2013). Constructing Experimental Designs for Discrete-Choice Experiments: Report of the ISPOR Conjoint Analysis Experimental Design Good Research Practices Task Force. Value in health : the journal of the International Society for Pharmacoeconomics and Outcomes Research. 16. 3-13.
12. <https://www.sawtooth.com/index.php/blog/archives/conjoint-discrete-choice-model-output-whats-all-the-share-about/>