# Performance comparison of the discrete choice models of consumer choice

**Exploration of the Econometrics and Machine Learning models' performances in the presence of heterogeneous preferences and random effects utilities**

**Research master thesis**

Nikita Gusarov

Master 2

MIASHS C2ES (UGA)

Under supervision of:

Iragaël Joly, HDR (GAEL, UGA, Grenoble INP)

Beatrice Roussillon, MCF (GAEL, UGA)

Université Grenoble Alpes

Faculté d'Economie et Gestion - FEG

2019 - 2020

**Abstract:** This works is a cross-disciplinary study of econometrics and machine learning (ML) models applied to consumer choice modelling. To breach the interdisciplinary gap an integrated simulation and theory-testing framework is proposed. It incorporates all essential steps from hypothetical setting generation to the comparison of various performance metrics.

The flexibility of the framework in theory-testing and models comparison over economics and statistical indicators is illustrated based on the work of Michaud, Llerena and Joly (2012). Two datasets are generated using the predefined utility functions simulating the presence of homogeneous and heterogeneous individual preferences for alternatives' attributes. Then, three models issued from econometrics and ML disciplines are estimated and compared.

This study shows the proposed methodological approach's efficiency, successfuly capturing the differences between the models issued from different fields given the homogeneous or heterogeneous consumer preferences.

**Key words:** Consumer Choice, Preference Studies, Willingness to Pay, Econometrics, Data Science, Machine Learning, Classification Techniques, Synthetic Datasets

**Author:** Nikita Gusarov (UGA)

**Under supervision of:** Iragaël Joly, HDR (GAEL, UGA, Grenoble INP); Beatrice Roussillon, MCF (GAEL, UGA)

---

**Abstrait:** Ce travail est une étude interdisciplinaire des modèles d'économétrie et d'apprentissage automatique (ML) appliqués à la modélisation des choix des consommateurs. Pour briser la frontière interdisciplinaire, un cadre intégré pour tester des théorie est proposé. Il intègre toutes les étapes essentielles de la génération de paramètres hypothétiques à la comparaison de diverses mesures de performance.

La flexibilité du cadre dans les tests de théorie et la comparaison de modèles par rapport aux indicateurs économiques et statistiques est illustrée à partir des travaux de Michaud, Llerena et Joly (2012). Deux ensembles de données sont générés à l'aide des fonctions d'utilité prédéfinies simulant la présence de préférences individuelles homogènes et hétérogènes pour les attributs des alternatives. Trois modèles issus des disciplines économétrie et ML sont ensuite estimés et comparés.

Cette étude montre l'efficacité de l'approche méthodologique proposée, en captant avec succès les différences entre les modèles issus de différents domaines compte tenu des préférences homogènes ou hétérogènes des consommateurs.

**Mots clés:** Choix du consommateur, Études de Préférences, Consentements à Payer, Économétrie, Science des Données, Apprentissage Automatique, Techniques de Classification, Données Synthétiques

# Acknowledgements

# Summary

# Introduction

The advances in statistical learning, data analysis and data science of the past decades have resulted in propagation of *Machine Learning* (ML) techniques to different applied fields, including social and human sciences. Nowadays, it is impossible to imagine a field of science that is not benefiting from the fruits of statistical learning. The works of De Palma et al. (2011) and Cascetta (2009) on transportation modelling, the publications of Molina and Garip (2019) dedicated to sociology problematic, the articles of Coussement, Benoit, and Poel (2010) concerning marketing decisions, actuary analysis studies (Denuit and Trufin (2019), Denuit and Hainaut (2019)) or even psychology with an example of Baayen et al. (2017) work reflect the literal omnipresence of the newly developed techniques.

However, there exist two completely distinct approaches to applying statistical learning, as described by Breiman and others (2001) and latter by Athey and Imbens (2019): the *Machine Learning* which focuses on the predictive qualities (figure 1c) and *Econometrics* which attempts to decipher the underlying properties of the data (figure 1b). In economics, where the research is focused on hidden patterns exploration, the scientific community prefers to implement the traditional econometrics techniques using the more advanced statistical models only in some special cases or as some assistance tools (Athey 2018). This discrepancy is explained by the fact that econometrics, contrary to traditional ML paradigm focusses on the accessibility of results. Consequently, many of the advanced ML techniques rarely appear in economics publications because of their believed lack of interpretability and excessive complexity in application. Nevertheless, some multidisciplinary scientists make attempts to breach this wall between *ML* and *Econometrics*: Varian (2014), Mullainathan and Spiess (2017) or, among the most recent, Athey and Imbens (2019). Their advances are mostly focused on resolving the general interdisciplinary tool-set integration questions, without considering the application specific details. Nevertheless, in the attempt to breach the interdisciplinary barrier the details reveal themselves to be of utmost importance in the solution of the problem.

Figure 1: The different paradigms



(a) Real world

(b) Econometrics

(c) Machine Learning

There have already been a multitude of studies comparing the performances of different econometric and ML models in various real world scenarios: the study of machine learning methods to model the

car ownership demand estimation of Paredes et al. (2017), for example; or the use of decision trees in microeconomics of Brathwaite, Vij, and Walker (2017). However, there's no known to us work incorporating at least all the baseline models, as it would require an unimaginable amount of efforts to accomplish. For instance, in the literature the performance of competing models are studied according to several absolutely alien criteria: in terms of the quality of data adjustments, in terms of predictive capacity, as well as in terms of the quality of the economic and behavioural indicators derived from estimates and, finally, according to their algorithmic efficiency and computational costs. None of the known to us articles manages to incorporate all these aspects into their benchmarks, limiting their studies only with several performance criteria.

These various aspects, greatly impact the performance of particular models or algorithms, although some of them are often ignored by the researchers. Not only there exists inconsistency in the targeted performance metrics in the contemporary models' comparisons, but there is also omnipresent problems of theoretical background choice, dataset selection or model's specifications. For example, speaking about the datasets used to support their findings, many researchers explore the impacts of different specifications on the same observed or simulated choice situation (Munizaga and Alvarez-Daziano (2005), Fiebig et al. (2010), McCausland and Marley (2013), Bouscasse, Joly, and Peyhardi (2019)) as it appears to be the most theoretically reliable procedure. However, there is still no established unified methodology documenting this field.

From this unambiguity in the scientific community the main problematic of this work arises. It is particularly important to establish a common framework for performance comparison of the discrete choice models be they from the econometrics or ML tool-set. However, this task cannot be accomplished outside a precise context, which will potentially impose some limitations over the models' structure, as well as influence the choice of performance metrics. In economics the discrete choice models are extensively used for consumer choice analysis (Anderson, De Palma, and Thisse 1992), willingness to pay derivation (Michaud, Llerena, and Joly 2012) and other preference studies. The field specific theories and traditional research objectives frame and define this study's scope.

From the economics perspective there exist three major points of interest to be taken into account. First, there is a strong interest in economics to explore the different behavioural set-ups, under different settings and assumptions. Secondly, given the different choice situations there is a potential need to test how the available mathematical models, potentially sensitive to the tested behavioural hypotheses or dependent on these hypotheses, perform in a given context. Last, but not least, a comprehensive implementation of a performance evaluation methodology, combining reproducibility and control of experimental conditions, should be introduced in the proposed framework.

**Consumer choice**

The economic decision theory derives mostly from the random utility theory (RUM) of McFadden (1974) and more recently of McFadden (2001), that were recently challenged by alternative visions such as random regret minimisation theory (RRM) of Chorus (2010), with a related relative advantage

maximisation theory (RAM) of Leong and Hensher (2015), or even quantum decision theory (QDT) of Yukalov and Sornette (2017), which offers a wide range of tools for modelling under uncertainty.

These different theories address various aspects of the decision making process, under different suppositions and incorporating different biases. For example, one of the basic assumptions of the traditional choice theory is the transitivity of choice, meaning there exists a strict hierarchy of individual preferences among alternatives. This assumption may be unsuitable for real world choice situation and lead to potential bias, which is addressed by quantum decision theory. QDT manages to bypass this shortcoming and incorporate non-transitivity of choices into the framework. There exist a multitude of other behavioural elements unexplained by the most traditional models that may be incorporated into the decision making framework, such as loss aversion for example, that could be addressed with random regret minimisation theory.

There is a particular interest in detecting the differences in the models' performances depending on the choice context and the assumed decision-making framework. It is important, because different consumer behaviour in the individual choice context result in different choice distributions, which may affect the models' performances. In economics RUM theory is nowadays one of the most used choice settings in the individual decision modelling. Nevertheless, there still exist some unexplored limitations, that such theoretical framework may impose over the estimation techniques, as well as to what potential biases a model's misspecification may lead.

**Mathematical models**

In general any classification technique may be used to model individual decisions, although nearly every model has some restrictions and limitations, which may largely affect its performances in a given context.

Usually the choice of model is rarely discussed in applied studies, as the researchers tend to use either the simplest model possible or attempt to implement one particular model of interest ignoring some times the other possible choices. For example, many traditional econometrics studies, given a multiple choice problem context, use a multinomial logistic regression (MNL) or even simplify the problem to a binary case, allowing to implement even more traditional models such as binary logit or binary probit models. However, there exists a multitude of particular cases in modelling individual choices, that require specific techniques to be implemented. A family of duration models may be used to model the individual decisions over time (Vitetta (2016)); network modelling that allows to incorporate spatial and social dependencies for the explored data (Brock and Durlauf (2003)); preference learning techniques aiming to explore the positioning of different alternatives by an individual (Tsoukiàs and Viappiani (2013), Pigozzi, Tsoukiàs, and Viappiani (2016)) and many other advanced techniques from *machine learning* field such as neural networks or support vector machines.

An incorrect choice of the modelling technique may have a strong impact on the derived target values leading to some erroneous conclusions in the end. For example, an incorrectly estimated willingness to pay for a particular product may lead to significant losses. When conducting an applied research study

one should always be conscious of the eventual biases introduced by the choice of the model and the eventual consequences of these choices. Some models are not suitable to be implemented on a particular set of data, while others are unable to provide necessary information about the relationships within a particular dataset or derive the particular target values of interest.

Taking into account the implications of RUM theory, there exists a particular interest to make the focus on the state of art econometric discrete choice models (Agresti (2013), Agresti (2007), Baltagi (2008), Train (2009), McFadden (2001), McFadden (1974)) as well as their counterparts used in ML (Hastie, Tibshirani, and Friedman (2009), Kotsiantis, Zaharakis, and Pintelas (2006)). A comparison of some simple models against more complex ones may reveal the trade-off between precise estimates and the resources invested.

**Data**

Different sources of data are available for a researcher, that could be divided into two groups (Japkowicz and Shah 2011): *field datasets*, which are gathered through an experiment or collected from the real world observations or real world uncontrolled experiment; and *synthetic datasets*, which are artificially generated by the researcher to suit his needs and respect some particular limitations. Although this variability of dataset choices not that evident in the context of applied studies, there is an ongoing debate concerning the eventual impacts of data choice on the models' performances and resulting metrics.

Given a task of performance evaluation and comparison for different algorithms or mathematical models there is always a difficult choice of the data type to be used in the study. Both of the mentioned above dataset types have their advantages and disadvantages and require a particular attention. However, having for objective the theory- and model-testing framework construction there is a strong interest to use the artificially generated data in order to have as much control as possible over the situation.

**The framework and context**

Given these three key elements we propose an integrated simulation and theory-testing framework which will encompass all the different aspects of the model comparison task. The steps to be integrated into such framework encompass many theoretical questions starting with the underlying theoretical assumptions and ending with the choice of correct performance metrics. Consequently, this work attempts to fill the gap between two statistical paradigms: *econometrics* and *machine learning*, taking into account the key elements, among which the different combinations of decision theory assumptions, dataset generation procedures, mathematical models and target performance measures. The problematic arises from the insufficient points of contact among researchers from different fields of applications, as well as insufficiently unified methodology to put into relations the different approaches. A work that uses unified knowledge from several disciplines might be highly beneficial for the scientific community as it will lie a foundation and provide support for future applied studies. Following the logic of Athey (2018) and Mullainathan and Spiess (2017) the project will attempt to merge the essentials of ML and econometrics paradigms, retaining their key concepts in the context of consumer choice problem.

We propose to use an applied paper in econometrics of choice modelling to facilitate understanding

of the field of application and tools. This means not that we will attempt to replicate the results, but rather to use the context provided in the work for demonstration of the proposed hypothesis-testing framework. We select the article of Michaud, Llerena, and Joly (2012) as our reference paper, because of the advantages to work directly with the authors of the paper. The work of Michaud, Llerena, and Joly (2012) is focused on investigation of consumers' willingness to pay (WTP) for environmental attributes of a non-food agricultural products, taking roses as example. Authors constructed an experimental framework to derive the premium the testing subjects were ready to pay for such environmental attributes as lower carbon imprint and ecological labelling, certifying the source of the environmentally friendly practices. That study explored individual preferences for roses with an eco-label and a carbon footprint using discrete choice modelling techniques and real economic incentives resulting in real purchases of roses. The gathered dataset was analysed with a mixed logit model demonstrating notorious premiums for both attributes. We will benefit of the obtained results to demonstrate all of the complexity of a proposed theory-testing framework, its functionality and perspectives.

The present report is divided into two main parts. The first section presents the chosen context for this work followed by short presentations of all the theoretical aspects which play their major roles in this study, tracing at the same time parallels with the context. The second part presents the results of all the results step-by-step, demonstrating the functionality of the designed framework. Each of the sections has an identical logical structure of presentation of the framework's components in the successive order: starting with the behavioural modelling and data related questions, directly followed by the models' presentation and the performance measures. The final section concludes.

# 1 The framework design

This section introduces the design and provides an example of available functionality of an integrated experimental framework for model performance exploration. In doing so, we strive to reduce and simplify the framework, illustrating the theoretical discussion of the eventual questions that arise during the model evaluation.

There exist multiple ways to provide an illustration for the generalized framework due to its extended flexibility on different levels of scientific procedure. Nevertheless, in our work we are attempting to extend this illustrative objective to all the possible levels available by the devised tool-set. The idea is to demonstrate all the features of different frameworks' layers in the context of a performance comparison. First of all, there is a particular interest to demonstrate the advantages of possibility to test different choice settings, providing different artificial datasets for exploration. What is more, it would be interesting to contrast different mathematical models and algorithms used to study these datasets and evaluate their performance using different criteria, which will allow for more flexibility.

The work of Michaud, Llerena, and Joly (2012) investigates the impacts of the environmental characteristics in the context of a consumer choice of non-alimentary agricultural goods taking roses as an example. We will inspire ourselves with the context, assumptions and findings of this study and build our work around these pre-sets. We may be interested to observe how some minor changes in the model may affect the results, which pushes us to consider some simple, yet educative changes in the model.

The organisation of this section is as follows. First of all, we introduce in detail the context and discuss which features and characteristics to retain give the Michaud, Llerena, and Joly (2012) work. We will provide a description of the procedure adopted for this illustration procedure as well. After brief overview of the original article and delimitation of the general assumptions we will provide a detailed discussion over every single major part of testing framework with extensive argumentation. Starting with the presentation of the underlying concepts of the decision theories and dataset generation procedure we will continue with a discussion of different modelling techniques and a detailed description of the models to be tested over the artificial dataset. Finally, we will provide a panorama of the performance assessment metrics, before switching over to application.

## 1.1 Context: Willingness to pay for environmental attributes of non-food agricultural products

We choose to use the estimated results of Michaud, Llerena, and Joly (2012) as a starting point for our work, copying the context of the study with some minor adjustments. In this part we will provide only a general overview of the assumptions made in the article "as is". This description will serve us as a reference for future discussion, because afterwards we will be presenting our changes, modifications and additions to these given our needs.

In the article of Michaud, Llerena, and Joly (2012) the choice of roses as the non-food agricultural product was determined by several criteria. Initially, roses were supposed by authors to have characteristics that respect the limitations imposed by the experimental economics. These are popular widespread products known to all the test subjects, being not easily available at the same time. What is more, the production of roses have been the object of a growing attention because of potential environmental damages inflicted in the process. This last feature made them a perfect product to explore the impacts of the environmental properties on the consumer choice.

Two environmental aspects of roses' production were explored by Michaud, Llerena, and Joly (2012). The first one, eco-labelling, described the cultivation environment and conditions, including the use of pesticides, fertilizers, as well as reasonable consumption of water and energy. This labelling was adopted shortly before 2010 by some of the producers, who attempted to reduce the harm to environment, to signal their eco-responsible position to consumer. Authors mention such dedicated eco-labels as the American *VeriFlora* "Certified Sustainably Grown" label guaranteeing the low environmental impact of roses' production, or the European equivalent: "*Fair Flowers Fair Plants*" (FFP) label certifying the environmental performance of agricultures by several criteria such as the "*fertilizer use, crop production, energy efficiency, waste management and a number of social requirements*". The second chosen environmental feature of roses was their carbon footprint, measured by the greenhouse gases emissions during the cultivation and transportation. This criteria being particularly important because of an increase of roses production in developing countries in Africa, South America and Asia, which are later sold on the European market, resulting in immense amount of CO2 emissions during the transportation.

The authors assumed that the individuals had heterogeneous preferences for the environmental attributes of roses. In other words, it was assumed that each individual had his personal attitude to the eco-label and carbon footprint of the roses, determined by their awareness of the environmental issues. The experimental design took into account this assumed dimension through observation of multiple simultaneous choices for each of the subjects in order to capture individual specific elements. To model such complex repeated choice framework, authors used well developed RUM behavioural theory (McFadden 2001) paired with the power of the mixed logit model, which is a generalisation of a simple logit model, allowing for more flexibility, such as random effects modelling.

The assumptions made by the researchers may be roughly divided into two categories, which will define the structure of this section. First one comprises the behavioural assumptions concerning the decision-making procedure, which encompasses the different restrictions on the experimental design, individual's behavioural strategy and choice preferences, which aim at elimination of various behavioural biases and simplification of future mathematical analysis and data treatment. The second regroups the assumptions related to the modelling process. It encompasses theoretical assumptions imposing restrictions on the mathematical model, its choice and estimation techniques. Finally, we present the target effects computed by the researchers in the context of the study, as the main objective was not the general approximation and modelling of a consumer choice, but rather extraction of particular values of interest

7

such as willingness to pay for the alternatives' attributes.

### 1.1.1 Experimental design

First of all, we should start with a description of the experimental design framework introduced by the authors in order to obtain valid results. This would allow to correctly implement such complex econometric model as mixed logit on the next stage. The experimental design assumed that individuals make their decisions based on the perceived utility of a particular alternative, following the traditional restrictions described by McFadden (2001).

Because the study collected data through a controlled experiment setting, some restrictions were imposed on the observed characteristics in order to simplify the analysis. The roses, as available alternatives, were defined by three attributes observed by subjects:

- the FFP EU eco-label (*Label*)
- the carbon footprint (*Carbon*)
- the price (*Price*)

These attributes varied across the available options of the alternatives present in different choice sets. Precise written instructions were transmitted to the subjects making available information about the criteria certified by the FFP labelling as well as some briefing about the organization issuing this labels (the Horticultural Commodity Board). These data-sheets provided as well a summary of Cranfield University's report about roses' carbon footprint. Both these attributes (eco-label and carbon footprint) were understood as a binary variables valuing 0 or 1 depending on the presence of a particular attribute for a particular rose. Finally, in addition to the two environmental attributes, a price was introduced into experimental design, which varied by 0.50€ between 1.50€ and 4.50€, creating this way a seven level factor. The table 1 regroups the main characteristics for these variables.

Table 1: Alternatives' attributes

| Statistic | Levels | Min | Max | Step |
|-----------|--------|------|------|------|
| Price | 7 | 1.5€ | 4.5€ | 0.5€ |
| Label | 2 | 0 | 1 | 1 |
| Carbon | 2 | 0 | 1 | 1 |

In order to avoid the substitution bias[1] as the subjects might have decided to purchase a rose for the experiment somewhere else for a lower price rather than in the laboratory a special measure was introduced into experimental design. In the experimental literature the implemented method is known as the

---

[1]The substitution bias occurs when the individuals tend to switch to the less expensive alternative available, given the relative prices changes. In the experiment context, the customers might have preferred to buy identically priced roses in a better placed store, instead of waiting for bought in experiment roses to be delivered.

"*field price censoring*", which means that the values used in the laboratory are censored according to the field market price (Harrison, Harstad, and Rutström 2004).

The elicitation of the individual preferences for the different roses' attributes was ensured through a combination of discrete choice questions and real economic incentives. The stated choice surveys are a popular choice for study of consumer preferences for public and private goods. The discrete choice methodology and experimental design setting provides the advantage to vary several attributes of a particular product and to estimate the marginal rates of substitution between these attributes. A particular accent was made on the derivation of the willingness to pay (WTP) for specified features of interest. This tool provides a great flexibility, allowing to test different scenarios all of which could be presented in a single study, although, there is always a danger that the choices made by consumers in experimental surveys might not reflect their real preferences. The participants to hypothetical surveys were generally stating higher WTP values for private and public goods, leading to a potential bias in the estimates when compared to the real world. Following this reasoning the authors have introduced incentives into their choice experiment linking this way the participants' decisions to real consequences by resulting in acquisition of randomly chosen alternative from the pool of chosen alternatives.

The choice set generation was devised with intention to resemble to maximum the actual purchase decisions with the inclusion of a "*do not buy*" option, in order not to force the subjects to buy anything. In other words, the presence of such alternative ensured that subjects were never pushed to purchase a rose, imitating this way a real shopping situation, when consumers always have the possibility of not purchasing any roses if none of the alternatives suited them in a particular choice set. Consumers were asked to make twelve different choices displayed.

In the case of prices allocation random design techniques were used to configure the subsets of choice sets among subjects. The two level factors standing for the roses' environmental attributes could be regrouped into four different combinations defining four types of roses. The experimental design introduced the roses in pairs to subjects, creating this way several three alternatives choice sets. Even though the different combinations of two roses potentially create sixteen different alternative pairs, the authors limited their search to six completely different pairs of roses. The resulting experimental sets of six choice sets were repeated twice resulting in twelve cards, which were then introduced to subjects. All the cards were distributed simultaneously so that consumers could make their choices in any order. Individuals were informed from the beginning that one of their decisions would be randomly drawn at the end of the experiment. Finally, the random draw resulted in the purchase of a real rose offered against payment, this condition ensured that the subjects considered each choice made during the experiment as a real purchase decision, weighting carefully the available alternatives.

Generic titles were randomly allocated to the roses within choice sets: rose A and rose B respectively. Such "*unbranded*" alternatives' titles allowed to ensure that they can only be differentiated according to their attribute combinations. This way the choice between a "*Rose A*" and a "*Rose B*" can only be defined by their attributes alone (Label, Carbon footprint and Price), but not by their label. The same strategy applied to prices, which were randomly assigned to the alternatives within the choice sets by a

random number generator setting prices within the defined limits.

Taking into account the experimental design we are going to follow the authors' ideas in simulating an identical experimental design with statistical methods available. The artificial choice situation will assume three alternatives: two unlabelled ones doted with a common utility function, while the third is the baseline alternative of "no choice" option. In order to study the heterogeneity of the individual preferences the subjects should be placed in a situation of repeated choice, facing several choice situation. The alternatives will be described by three attributes, while individuals will be distinguished by four characteristics.

### 1.1.2 Econometric model

Consumers' decisions are analysed with the discrete choice framework based on the utility maximisation assumption. This framework assumes that consumers associate each alternative in a choice set with a utility level and choose the option, which maximises this utility. The general estimation framework of the Random Utility Model (RUM) proposed by McFadden (1974) provides the opportunity to estimate the effects of product attributes and individual characteristics and to compute willingness to pay indicators.

Authors implemented the mixed logistic regression with random, correlated attributes' effects to estimate the willingness to pay of the individuals for each of the explored attributes of a rose. The mixed logit model takes into account the repeated nature of the choices made by the respondents. This model relaxes the Independence from Irrelevant Alternatives (IIA) hypothesis of the more traditional multinomial logit, allowing the random components of the alternatives to be correlated, at the same time the error terms are still considered to be identically distributed (Greene 2008). The alternative specific parameters are assumed to be randomly distributed across the population contrary to the fixed parameters specification for a traditional multinomial logit model. In other words, the mixed logit model provides the opportunity to consider heterogeneous effects among individuals by allowing taste parameters to vary in the population. The authors suppose that the random taste heterogeneity should be evident in response to the eco-label and the carbon footprint attributes of the roses, because of different level of environmental awareness across population. Following the ideas of Bernard and Bernard (2009) the authors introduce the cross-product for eco-labelling and carbon footprint as a random parameter as well in attempt to test the effect of the simultaneous presence of both of these attributes on consumer choice. This addition results in a total of four random parameters to be estimated: the two parameters describing roses attributes, their cross-product and the "*Buy*" option dummy variable which captures heterogeneity in consumers' preferences for a rose. All of the random parameters associated with the roses' attributes are assumed to follow normal distribution, which is traditional for the procedure of mixed logit modelling. Given that the normal distribution is symmetric and unbounded, the resulting model allows for both positive and negative effects to exist inside population. To simplify the analysis and assuming the reasoning of Revelt and Train (1998), the authors restrict the price coefficient to be

fixed in the population. Such choice of price's effects specification ensures that all respondents have a negative price coefficient, leading to a normally distributed estimate of willingness to pay.

The systematic part of the utility relatively to the "No buy" option was expressed through a linear in parameters form:

$$V_{ij} = \alpha_{i,Buy} + \beta_{Buy,Sex}Sex_i + \beta_{Buy,Age}Age_i + \beta_{Buy,Income}Income_i + \beta_{Buy,Habit}Habit_i +$$
$$+ \gamma_{Price}Price_{ij} + \gamma_{i,Label}Label_{ij} + \gamma_{i,Carbon}Carbon_{ij} + \gamma_{i,Label \times Carbon}Label \times Carbon_{ij} \quad (1)$$

Where $j$ was an alternative among the available choice set of three options: buy rose A, buy rose B or do not buy anything. The dummy variable $\alpha_{i,Buy}$ was introduced to capture the effect of a decision to buy a rose, while the vectors of $\beta$ and $\gamma$ regrouped the effect of individual characteristics and the attributes of alternatives respectively.

In their article Michaud, Llerena, and Joly (2012) did not provide an extensive demonstration or description of the model selection procedure. What is more, we have little information as to what model comparison and validation techniques were implemented. Only the final model, chosen by authors, was presented to us, which brings some limitations for our study.

In the end of this subsection, it is important to highlight that the mixed logit models are usually specified with uncorrelated random effects, although it's not the case in the context of this particular study. The authors introduce correlation between the normally distributed alternative specific coefficients: $\alpha_{i,Buy}$, $\gamma_{i,Label}$, $\gamma_{i,Carbon}$ and $\gamma_{i,Label \times Carbon}$.

### 1.1.3 Willingness to pay and premiums

The only target metrics present in the article were the willingness to pay (WTP) and premiums for particular attributes. The former could be read as the value the consumers are willing to pay for a rose. The latter may be translated as how much consumers are ready to pay for a unit change of a given attribute of the product. Both the WTP for a product and the premiums can be computed as the marginal rates of substitution between the quantity expressed by the attributes and the price (Louviere, Hensher, and Swait 2000). The WTP for a rose in this case could be expressed as:

$$WTP = \frac{\frac{\Delta V}{\Delta BUY}}{\frac{\Delta V}{\Delta Price}} = \frac{-\alpha_{Buy}}{\beta_{Price}} \quad (2)$$

Where $\frac{\Delta V}{\Delta BUY}$ is the difference in the relative utility $V$ associated with the "Buy" and "No buy" choices. The premiums for the particular attributes $Z_k$ of a given product could be identically expressed as:

$$WTP = \frac{\frac{\Delta V}{\Delta Z_k}}{\frac{\Delta V}{\Delta Price}} \tag{3}$$

Since the random parameters of the utility function were assumed to be correlated, authors used Krinsky and Robb parametric bootstrapping method (Krinsky and Robb 1986) with 1000 draws to estimate the standard deviations and confidence intervals for these parameters.

## 1.2 Theories of consumer choice

Once we have presented the assumed context for this study, we will dive further into details and present all the key behavioural elements of this work one by one. In this section we are going to present in detail the questions and problematic associated with the behavioural modelling of the consumer choice. Particularly, we are going to introduce the terminology to be used in this work, some of which was already partially presented in the previous section.

### 1.2.1 General terminology

For the presentation of general methodology we are going to adopt the ideas of De Palma et al. (2011), introducing this way the principal concepts and main components of the decision theory. Traditionally it comprises several components: the decision makers or *individuals*, described by their characteristics; a set (or sets) of available *alternatives*, defined by their attributes; and a decision rule or set of rules, describing the procedure adopted by the individuals to make actual decisions.

The individuals are supposed to have different tastes, and therefore we must explicitly treat the differences in the decision-making processes among individuals, doted with different characteristics. Therefore the characteristics $X_i$ of the decision maker $i$ constitute an important part of the problem.

The decision maker chooses from a finite and countable set of alternatives $\{\omega_i, \ldots, \omega_j\}$, which consists of the entire *universal set of alternatives* $\{\omega_1, \ldots, \omega_r\} \in \Omega$ as defined by the particular choice environment. A decision maker $i$ may only consider a subset of this universal set $\Omega$, and this consideration set is conventionally named *a choice set* $\Omega_i$. In discrete choice analysis, each alternative $\omega_j$ is characterized by its attributes $Z_j$. For example, in the particular case study the observed attributes of roses are their price, the eco-label and the relative carbon footprint. Decision makers evaluate the attractiveness of an alternative based on these attribute values before making their choice.

Finally, the decision rule describes the process by which the decision maker $i$ evaluates the available information $Z_j \forall \omega_j \in \Omega_i$ and arrives at a unique choice. There is a wide range of available decision rules, including dominance, satisfaction, lexicographic, elimination by aspect, habitual, imitation, and utility (De Palma et al. 2011). However, only the latter class is most often associated with discrete choice analysis because to its extensive use in the consumer choice behaviour modelling. The utility

theory takes its roots from the microeconomic consumer theory and is adjusted according to the needs of the modeller. A utility $U_{ij}$ represents the attractiveness of a particular alternative $\omega_j$ for a particular individual $j$ in a scalar form.

### 1.2.2 Random utility maximisation models

The random utility maximisation models (RUM) were introduced and developed by McFadden (1974). The theory of optimization implies that this is a classical indirect utility function, with the following properties: "*it has a closed graph and is quasi-convex and homogeneous of degree zero in the economic variables*" (McFadden 2001). The last element in applying the standard model to discrete choice is to require the consumer's choice among the feasible alternatives to maximize conditional indirect utility based on some reference alternative, rather than absolute utility.

In our work we use the notation introduced by Bhat (1995) and later adopted by Cascetta (2009) when representing the utility functions as they are more simple and easy to understand compared to initial McFadden (1974) specification. The functional form of the canonical indirect utility function depends on the structure of preferences, including the trade-off between different available alternatives. The perceived utility $U_{ij}$ can be expressed as the sum of two terms: a systematic utility and a random residual term:

$$U_{ij} = V_{ij} + \eta_{ij} \tag{4}$$

Where $U_{ij}$ stand for utility, $V_{ij}$ at the same time represent its deterministic part defined by some fixed deterministic function and $\eta_{ij}$ reflects some unobserved random effects. The latter having being a random variable following Gumble distribution, parametrized with $(\mu = 0, \theta = 1)$, which may be interpreted as:

$$\eta_{ij} = -log(-log(\epsilon_{ij})) \tag{5}$$

With $\epsilon_{ij}$ a variable uniformly distributed and independent across alternatives, the disturbances are independently identically distributed Extreme Values (EV). This produces a MNL model in which the systematic utility has a linear in parameters form for each alternative $\omega_j \in \Omega$. The systematic utility $V_{ij}$ represents the mean utility perceived by all decision-makers having the same choice context decision-maker.

$$V_{ij} = f(X_i, X_j) + \eta_{ij} \tag{6}$$

Traditionally in the most simple models this deterministic utility part is represented by some linear in parameters function:

$$f(X_i, Z_j) = \alpha_j + \beta_j X_i + \gamma Z_j \tag{7}$$

One family of RUM-consistent discrete choice models that is very flexible is the random parameters or mixed multinomial logit (MMNL or more often denoted as ML) model, which is used in the Michaud, Llerena, and Joly (2012) work. The random parameters set-up assumes $\gamma$ effects to be randomly distributed across individuals, usually following normal random distribution. In some cases, these parameters may be assumed to be correlated, which potentially reflects better the real world.

In our study we are going to explore two equally possible in real life specification for data generation procedure: one assuming random effects for alternative specific variables and another keeping these parameters fixed. Speaking about the utility definition, we assume, that the work of Michaud, Llerena, and Joly (2012) managed to obtain correct estimates for a relative utility function of roses and we take this particular function structure in order to generate utilities for a given dataset. This assumption will offer us a baseline and target effects' values to compare our estimation with.

## 1.3 Different datasets available in research

There exist numerous difficult questions related to the models' comparison task such as performance measures' choice or models' specification, but beforehand there always stand the data related questions. It is due to the fact that all the other questions and the validity of the obtained responses rely entirely on the choice of the inputs and the data available. Many of the existing applied econometrics papers use the most simple specification of the Multinomial Logistic Regression (MNL), that may lead to erroneous results and conclusions.

Many of the models' performances and performance measures depend on the dataset properties and the particular application case. This means that in comparison of different mathematical models, implementing some complex tools such as a neural network models (NN), for example, we should pay attention to use appropriate data-model to estimate such model. This particular problem, as many others related to the models' performance evaluation, was extensively described by Japkowicz and Shah (2011).

When it comes to model comparison, the additional requirements arise to the validation datasets and we should find answers to several questions:

- What datasets should be used?
- Should the model be validated on one dataset or several several?
- Should a synthetic or real-world data be used?
- If several dataset are chosen, which ones should be used on different validation steps?
- How the algorithms should be tuned face to the dataset selection?
- What properties the studied data should have?

Moreover, different models may require different data adaptation methods to be implemented. For example, the popular multinomial logistic regression allows to take into account the individual characteristics as well as the attributes of the various alternatives issued from some limited set. The ML approaches, such as Support Vector Machines or Linear Discriminant Analysis does not allow such flexibility. For these models, even if we can represent each point in the modelled space as a combination of individual characteristics and attributes of alternatives, we can only classify the instances by iterative binary separation (Tsoumakas and Katakis 2007). Consequently, the questions of the dataset properties arises, which are tightly intertwined with the available models choice and the implemented learning techniques.

In this section we will discuss the different existing approaches to data management in theory testing and hypothesis verification. Firstly, we will present the general questions and problematics. Then a solution to be implemented in this particular study will be described and discussed.

### 1.3.1 Theoretical concerns in dataset selection

The data related problematic arise firstly during model generation step of the standard statistical learning procedure and persists till the stage of the model comparison. Speaking about the model validation, the usual *rule of thumb* approach is the cross-validation technique, although some advanced users suggest that this method may not be always appropriate (Japkowicz and Shah 2011). In econometrics, for example, as well as in many other applied disciplines, researches tend to oversimplify the validation step by completely avoiding this important step, or by performing only *single-fold* validation. On the other hand, many advanced statistical model and ML methods require a separate tuning step during model set-up, which alone requires verification and validation on some dataset. It remains questionable whether the overall model validation dataset and the dataset used for fine tuning should be the same or not.

Of particular interest for our study is the ongoing discussion between two sides of the statisticians' community, mentioned by Japkowicz and Shah (2011), about whether the algorithms and statistical models should be compared over the real world datasets or using some synthetically generated data. On one hand, the datasets composed of the observations or obtained through controlled experiments perfectly reflect the real world situation, being at the same time too case specific. In other words, it is always dubious that a model or a theory verified for one particular real dataset has any external validity. The obtained insights can rarely be extended over a larger population. Artificial data can be designed in a controlled manner to study specific aspects of the performance of algorithms and models. Moreover, the artificial data is highly useful for testing particular theories, for example, the behavioural theories or their impact on different models. Consequently such data may allow for tighter control, which gives rise to more carefully constructed and more enlightening experiments. Although, the real data are hard to obtain and are difficult to analyse, the artificial data introduces the danger of the problem's oversimplification. In our case study these features are of utmost importance, because, given

the framework, the artificial data enables us to test desired hypothesis in a controlled environment.

Generation of synthetic datasets is a common practice in many research areas. Such data is often generated to meet specific needs or certain conditions that may not be easily found in the original, real data. The nature of the data varies according to the application filed and includes text, graphs, social or weather data, among many others. In this particular work we, for example, face the consumer choice data, which describes individuals and their choice sets.
The common process to create such synthetic datasets is to implement small scripts or programs, restricted to limited problems or to a specific application.

As Garrow (2010) points it out, even observing the growing use of artificial data in discrete choice and classification analysis, "little is known about how the methodology used to generate synthetic datasets influences the properties of parameter estimates and the validity of results based on these estimates". That is, there are two potential sources of biases when using synthetic discrete choice data:

- The unknown effect of the dataset generation method;
- The parameter estimation bias.

The first one is rather complex and has many different element, that could potentially affect the estimated results. There exist different methods for artificial dataset generation, starting with use of *robots* (artificial observation instances) and ending with Markov Chains Monte Carlo simulation and Neural Network use. One of the most evident errors in this case could arise from the fact, that the closer the estimated model is to the model implemented to generate the dataset, the better would be the observed results, which may not be true in the real world.

The second bias arises in the situation where the real world parameters are used to generate artificial dataset, exactly as in this particular work. The potential difference between the ideal simulated situation and the real world situation lead to different choice structures. The theoretical model supporting the data-generation process may be potentially erroneous, leading to erroneous conclusions if only such dataset was used for incorrect purpose.

### 1.3.2 Artificial dataset generation procedure

For the objectives of this study we assume the best option is to generate our own artificial dataset based on a predefined utility function and given a predetermined statistical properties for individual characteristics and alternatives' attributes. Such set-up ensures that we know exactly the data generation process and have all the control over the parameters and experimental design. As was mentioned above, this choice may be dangerous in terms of justification of the resulting external validity of obtained results in application to any other real world dataset. However we ensure this way, that the obtained results could be potentially compared with the baseline target parameters and the initial effects are observed to us.

First step in the dataset generation is the generation of the experimental design framework, imitating the original choice set-up, as described in the Michaud, Llerena, and Joly (2012) article. Our first steps are identical to the original work, as we start with the generation of all possible combinations of binary factors for our alternatives: roses described by two binary attributes and their price. There exists only four different roses types, if described by their binary attributes alone, as can be seen in the table 2.

Table 2: Possible attributes of roses

| Type | Eco-label | Carbon footprint |
|------|-----------|------------------|
| 1    | 0         | 0                |
| 2    | 0         | 1                |
| 3    | 1         | 0                |
| 4    | 1         | 1                |

Given a multiple choice context when an individual is choosing among three alternatives: two different roses, defined by labels A and B; and a "No buy" option. Consequently there exist multiple possibilities to regroup two roses into a choice set, for instance, in Michaud, Llerena, and Joly (2012) are generating six choice sets ensuring that roses in a given choice set always have different attributes, while in practice there exist sixteen possible combination of two roses given they are described by two binary factor variables. The choice of the choice set delimitation in the article could be understood as the individuals participating in the stated choice experiment are scarcely interested in answering multiple questions, while six or twelve choices to consider appear to be a reasonable number. On the contrary, our experimental artificial set-up allows to ask as many questions to as many individuals as we want. For example we can generate $7 \times 4 \times 7 \times 4 = 784$ choice sets for each individual, containing all the possible combination of two different roses, each described by two binary factor attributes as well as their price, which has 7 different levels (varying by 0.50€ in a range from 1.50€ to 4.50€). However, such excessive set-up can have its toll on the computation times, being in the same time absolutely unreasonable and unrealistic, were we to replicate our results in a stated choice experiment. Consequently, for price allocation we are going to implement the same strategy as the authors of the article, meaning that the prices will be randomly assigned inside the choice sets, while the choice sets will follow a complete full-factorial design given two alternatives with attributes. The following table 3 demonstrates this idea.

The prices are randomly allocated within given choice sets, although there are some subtleties, which were discovered in attempt to replicate the variability achieved in the original work. The main idea is to ensure that both groups of roses (A and B) will have identical characteristics, which is important for the later model estimation. At the same time, we are interested in providing the test subjects with identical choice sets to avoid eventual bias, which may be important if we were facing a small number of observed individuals. Consequently, we randomly allocate prices within a given choice set and distribute these identical choice sets to all of the individuals. The variability in the prices across alternatives is achieved through a replication of this procedure $n$ times. The resulting statistics and distribution will be discussed

17

Table 3: Choice sets attributes' combinations

| | Rose A | | Rose B | |
|---|---|---|---|---|
| Choice set | Eco-label | Carbon footprint | Eco-label | Carbon footprint |
| 1 | 0 | 0 | 0 | 0 |
| 2 | 0 | 0 | 0 | 1 |
| 3 | 0 | 0 | 1 | 0 |
| | | ... | | |
| 15 | 1 | 1 | 1 | 0 |
| 16 | 1 | 1 | 1 | 1 |

in the second part of the work, where we will focus our attention on the applied part.

On the next step we generate a population of "*robots*", or artificial individuals, who will be making their choices provided the described above choice sets. It is important as well to mention, that the distributions we use to generate the data are theoretical rather than empirical ones. The individuals are generated based on the descriptive statistics for population available in the reference paper. This choice is done based on the final objective of the proposed testing framework to allow the researchers to test and verify their hypothesis related to the behavioural assumptions, modelling and performance estimation in the consumer choice experimental context. We assume that characteristics of the individuals are normally distributed, which is rarely the case in practice, where skewed distributions are dominant. Such choice imitates a replication attempt of a given empirical paper given the information available in the article only, which are usually the means and variances, rather than complete empirical distribution descriptions.

Finally, having at our disposal a set of individuals as well as a number of choice sets for the individuals to consider, we define the utility function based on the estimates of the authors. Such choice implies, that we assume all the hypothesis made when treating the original dataset to be verified for the artificial model. The utility functions are assumed as described in the preceding subsection to conform with the standard random utility maximisation (RUM) definition as the individuals are striving to maximise their perceived utility given their characteristics and the observed attributes of the alternatives. The utility is linear in parameters with additive error term.

Following this procedure we generate two synthetic datasets: one the most basic one with only fixed effects present, while the other includes random effects for the alternative specific attributes. These datasets are then used to estimate, test and compare the models' performances.

To summarise, this section we will once again list the key hypothesis we make in the artificial dataset creation:

- The dataset comprises:

- 4 individual characteristics ($Sex$, $Age$, $Habit$ and $Salary$)
- 3 alternative's attributes ($Price$, $Label$ and $Carbon$)
- 2 product variables ($Buy$ dummy variable and $LC = Label \times Carbon$ cross-product)

- The individuals are assumed to maximise their utility, when making their choices, which corresponds to RUM behavioural framework;
- The utility functions are linear, additive in parameters with an additive error term $\epsilon$;
- The error term is assumed to be iid. across population and follow a Gumble distribution: $\epsilon \sim G(0, 1)$;
- The individuals may (or may not) express heterogeneous preferences for the environmental attributes (eco-$Label$ and $Carbon$ footprint), which results in two different artificial datasets;
- In the case of heterogeneous preferences a total of four random parameters are assumed to be correlated ($Buy$ dummy, $Label$, $Carbon$ and their cross-product $LC$) and respect a multivariate normal distribution.

The detailed procedure of the choice modelling, as well as the exact values of the parameters and some eventual difficulties in the dataset generation are described in the applied section of this work.

## 1.4 Statistical tools for choice modelling

As it was mentioned, there are different fields of application ranging from *econometrics* (Agresti 2013) to *machine learning* (Zielesny 2011), encompassing eventually such fields as transportation systems analysis (Cascetta 2009) and logistics (De Palma et al. 2011), actuarial science (Denuit and Trufin 2019), preference learning (Fürnkranz and Hüllermeier 2010), psychology, sociology and more). The more generalised models are regrouped under the *statistical models* label (Hastie, Tibshirani, and Friedman 2009), but nevertheless they are mostly limited and are not taking into account many of the field specific questions. Taking into account that our study is mostly axed towards the study of the consumer choice data and related discrete choice problems it is important to somehow limit the study's scope to a number of selected models, without loosing the context.

Speaking about the econometrics models, this field of applied statistics alone has a number of questions to answer before proceeding. For example, we may question the particular task that we are performing while applying the econometric models to some *discrete choice* problematic. Usually the economists are interested in deciphering and understanding the underlying process (Athey and Imbens 2019), even though there is a long lasting debate on the validity of obtained measures as well as causality implications (Chen and Pearl 2013): "*The source of confusion surrounding econometric models stems from the lack of a precise mathematical language to express causal concepts.*" This results in completely different cultures of the data exploration and study objectives. This particular problem was largely addressed by different researches, among which: Athey and Imbens (2019), Mullainathan and Spiess (2017), Agrawal, Gans, and Goldfarb (2019), Varian (2014) and Breiman and others (2001). Even as

there are some attempts to merge all the existing branches and approaches to statistical modelling into some sort of a uniform culture (Donoho 2017), the scientific community has a long route to make in order to achieve this objective. There exist as well many other more subtle problems in the econometric field. For example, different error term and different link function specifications (Bouscasse, Joly, and Peyhardi 2019) in econometrics models rise the question of what exactly we may consider as single *entry* to our list of models to evaluate.

On the other hand, speaking about the ML counterpart, the focus is generally made on the predictive precision if we were to focus our attention on the supervised ML sub-field (Mullainathan and Spiess 2017). In their quest to achieve the best predictive precision with a particular model, the *machine learning* scientists study not only the theoretical models themselves, but the algorithms used to estimate these models (Zielesny 2011), that potentially augments the dimensions to take into consideration in this particular work. Moreover, not only there exist a confusion on what algorithms are to be associated with each particular model (or potentially a number of models defined by model/algorithm pairs), but many models are specified using a set of hyper-parameters, which are to be chosen by the researcher. This aspect immensely complexifies the task for us, as it is uncertain how exactly should we define the values of these arbitrary chosen parameters. It's worth mentioning that in many cases these parameters are case specific and may vary from one application to another, resulting in different performances over different datasets.
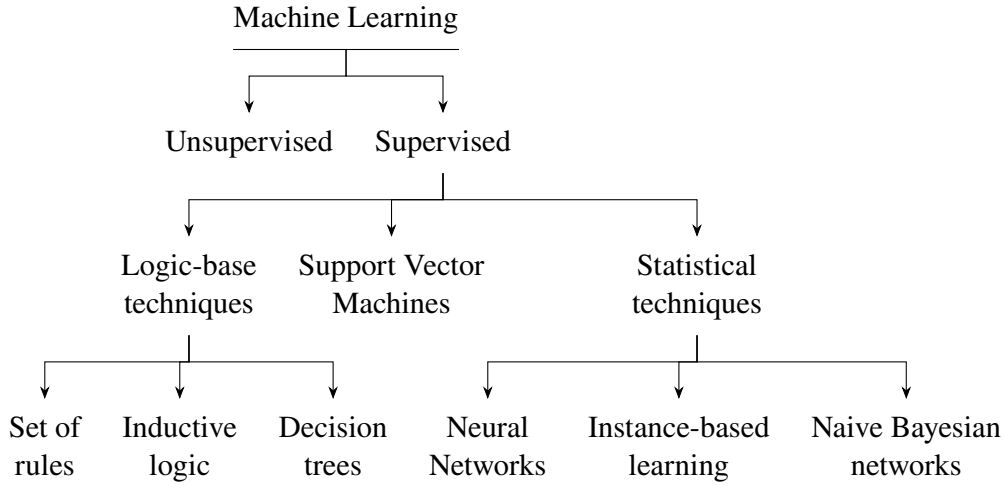
As it is mentioned by Kotsiantis, Zaharakis, and Pintelas (2006) the choice of which specific learning algorithm to be implemented is a critical step for any work, and a separate subset of training dataset is usually used for this task. The classifier's evaluation is most often based on prediction accuracy, which describes the percentage of correct predictions among their total number, which requires some unrelated data to be calculated as out of sample estimates provide more reliable information about the performance of a particular algorithm.

This section will be opened by a brief introduction to the multitude of the existing models, which is a particularly important point, given the scope of the study. Each and every dataset, each and every relationship between several variables may be modelled with different techniques and different assumptions. There is a tremendous amount of work to be done in order to systematise all the existing mathematical models, not speaking about their extensions or their numerical implementations. The first part of this section will demonstrate the complexity of the models' choice given an application context. Only then, we are going to present the selected models and their mathematical formulation: the MNL model, the MMNL model and their artificial NN counterpart.

### 1.4.1 Taxonomy of statistical models

Before proceeding with a discussion concerning eventual problems and difficulties affecting the modelling part of every empirical study, we will provide an overview of different families of models, encompassing both the *econometrics* and *machine learning* fields. The following presentation is a gen-

20

Figure 2: Taxonomy as proposed by Kotsiantis (2006)



eralised vision of the existing discrete modelling techniques, which can be used for classification tasks. As general as it is, this part respect the setting of the discrete choice behavioural modelling.

There exist several possibilities to divide ML algorithms into groups in order to provide an exhaustive and complete taxonomy of this field and the same reasoning may be applied to econometric models. However, the existing taxonomies are rarely complete and focus mostly on one or several grouping aspects. They define the general structure of a particular taxonomy, but rarely take into account a sufficient number of different descriptive features, which may vary across statistical models. For example, we may take a look at Kotsiantis, Zaharakis, and Pintelas (2006) work attempting to provide an overview of different classification techniques on figure 2.
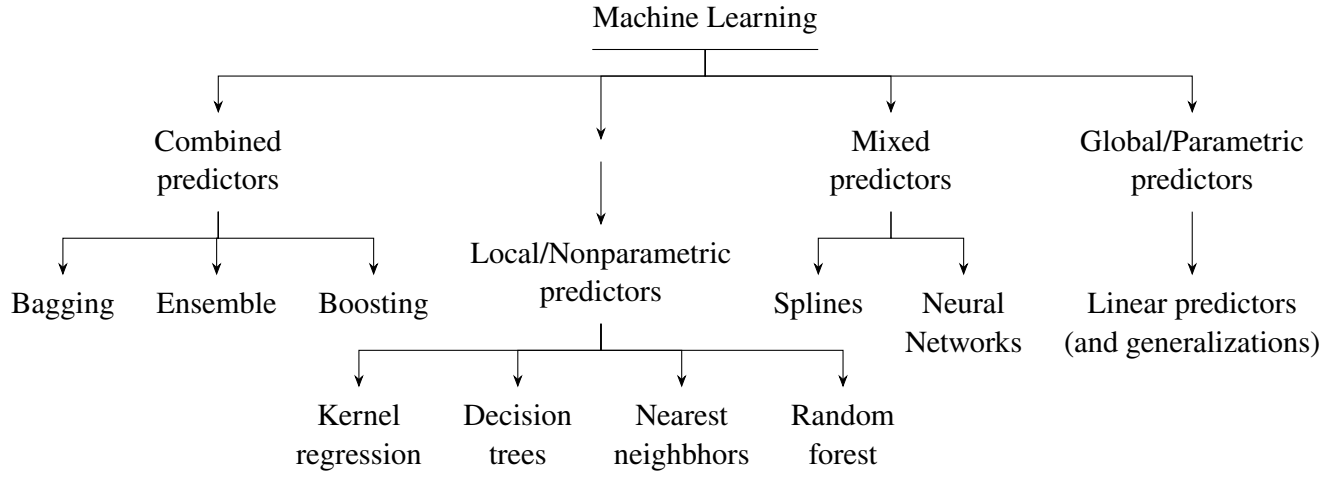
This taxonomy is fairly simple and encompasses a large number of models' families specifically designed for classification.
In the works of Hastie, Tibshirani, and Friedman (2009), Cascetta (2009) and Ayodele (2010) we may see some more recent attempts to organise the existing models into a single hierarchically related structure, although neither of known to the author works offers sufficiently extended reasoning over the relations between different classification techniques (several of the resulting taxonomies could be seen in the Appendix A). Moreover, not only the taxonomies may be based on the models' themselves, but it can be constructed around their algorithmic properties, as in Mullainathan and Spiess (2017). The resulting tree is represented on the figure 3.

In attempt to generalize the existing taxonomies and unite somehow the different classification models and techniques, we may roughly divide them in categories by different criteria. Usually there is no evident hierarchical dependency between the different criteria, which immensely complexifies the task of unified taxonomy construction.

First of all we may divide the models onto *supervised* and *unsupervised* learning techniques (Hastie,

Figure 3: Taxonomy as proposed by Mullainathan (2017)

```
                              Machine Learning
          ┌───────────────────────┼───────────────────────┐
     Combined                                   Mixed              Global/Parametric
     predictors                               predictors              predictors
  ┌──────┼──────┐       Local/Nonparametric   ┌─────┴─────┐              │
  │      │      │           predictors        │           │         Linear predictors
Bagging Ensemble Boosting                  Splines     Neural      (and generalizations)
          ┌──────┬───────┼───────┐                     Networks
        Kernel  Decision  Nearest  Random
      regression  trees  neighbhors  forest
```

Tibshirani, and Friedman 2009), which is the most widely used model separation in ML field. Sometimes this separation is complimented by various intermediate combinations of these two. The supervised methods have the goal to predict the value of an outcome measure based on a number of given input measures, the outcome variable is available through the learning process to guide the researcher and algorithm providing some baseline for testing. In the statistical literature the inputs are often called the predictors, the inputs, the features, or the independent variables. In the econometrics the terms explicative or endogenous variables are more popular. The outputs are denominated as responses, or, in econometrics, the dependent or endogenous variables. The unsupervised learning is used without any outcome measure available, with a main objective being to describe the associations and patterns among a set of inputs. Such formulation of a learning task is rather implemented to describe how the data is organized or clustered, find the underlying patterns and dependencies. As for the intermediate models' families, we may address the article of Ayodele (2010), where authors present different mixed types of learning tasks, although this particular classification is not widely used. Among these models we find: *semi-supervised* learning, combining both labelled and unlabelled examples to generate an appropriate function or classifier; *reinforcement* learning, in which algorithm learns to interact with the data generating source, given an observation of the world, in this context every action of model has some impact in the environment, and the environment provides feedback that guides the learning algorithm; The *transduction* is nearly identical to supervised learning, although instead of an attempt to construct a function it tries to predict new outputs based on training inputs, training outputs, and new inputs; and finally *learning to learn*, when the algorithm learns its own inductive bias based on previous experience, which is a more advanced reinforcement learning problem.

Depending on the output variable structure we attempt to model we may examine the taxonomy proposed by Agresti (2013). This taxonomy is based on the output variable format: it may be either discrete or continuous. The *continuous* variables are the simplest case, where the output is assumed to be continuous on a given interval and in the statistical society is usually addressed as "*regression*" task. It's

counterpart, the discrete dependent variable is sometimes addressed as "*classification*" task and it is the focus of this particular work. The *categorical* variable has a measurement scale consisting of a set of categories and these variables are of many types: binary variables, nominal data, ordinal data or count variables. The *binary* data assumes that there exist only two categories, often given the generic labels "success" and "failure" numerically represented as 0 and 1. In the context of the undertaken study we may imagine a binary variable representing the individual choice of "Buy" against "No buy". The *nominal* variables represent categories without a natural ordering and are measured on a nominal scale. The perfect example for this data type is our choice set delimitation with several unordered and independent options for individuals to consider: buy rose A, buy rose B or do not buy anything. For nominal variables, the order of listing the categories is irrelevant to the statistical analysis, and the main importance is given by the choice of baseline option, which is important for some of the statistical models. *Ordinal* data or ordered discrete data is an advanced representation for nominal data, where many categorical variables do have ordered categories, representing some given preferences order, for example. For these variables, the distances between categories are usually unknown and these intervals may be uneven between different categories. An *interval* variable is one that does have numerical distances between any two values. For most variables of this type, it is possible to compare two values by their ratio, in which case the variable is also called a ratio variable. The final class if the *count* data, which is specific for special cases of discrete-continuous data treatment.

By their structure the models may be separated into *additive* and *non-additive* as described in Hastie, Tibshirani, and Friedman (2009), both of which could be understood either as additive (non-additive) in error term or having a full additive (non-additive) structure. The first group encompasses different regression and classification models where either the main function has additive structure:

$$f(X) = E(Y \mid X) \tag{8}$$

Or the error term is additive defining the following model:

$$Y = f(X) + \epsilon \tag{9}$$

The *non-additive* models, also denominated as *multiplicative* models, include all other eventual specifications which could not be viewed or approximated by the additive relations. This particular separation could be extended even further, as the models could be viewed as *linear* and *non-linear* in their parameters, or in their overall functional form. The former either assume that the regression function $E(Y \mid X)$ is linear, or that the linear model is a reasonable approximation for the particular situation. The non-linear models usually regroup the various extensions and generalisations for the linear models integrating various non-linear transformations.

One more possibility to separate different discrete choice models in particular is by taking into account the probability structure they are attempting to model as mentioned in Jebara (2004). The models are

separated into two major groups: generative and discriminative models, to which sometimes a third ambiguous group of non-model techniques is added. The *generative* algorithms model the full structured joint probability distribution over the examples and the labels given by $P(Y, X)$. The models in this context are typically cast in the language of graphical models such as Bayesian networks. The joint distribution modelling offers several attractive features such as the ability to deal effectively with missing values, for example. On the other hand, the *discriminative* methods such as support vector machines or boosting algorithms focus only on the conditional relation of a label given the example, the probability being written as $P(Y \mid X)$. Their parametrized decision boundaries are optimized directly according to the classification objective, encouraging a large margin separation of the classes. They often lead to robust and highly accurate classifiers.

The estimates structure differs across model families as well, as described in Hastie, Tibshirani, and Friedman (2009). There are two principal approaches to modelling given by *parametric* estimators, which are usually easy to read and interpret, and their *non-parametric* counterpart, offering the best results in terms of precision in most cases. The multitude of non-parametric regression techniques or learning methods can be separated into a number of classes by the nature of the restrictions imposed, although we are not going to provide an extensive description of all of them. What is more important, that there exist different families of mixed models, profiting from both the parametric and non-parametric feature. They are traditionally regrouped into a single family of *semi-parametric* models.

In this work we face a classification task which can be understood, given the context, as consumer choice modelling. In order to correctly model the consumer choice structure we will need to use the models allowing to work with nominal discrete data, because the consumer choices can not be positioned in some logical order defining a continuous variable. The desire to obtain some explanatory results leads us to restrict our choice to some additive and, moreover, linear models, which would identify the parameters of a given relative utility function. The latter argument implies that the models should be parametric, producing some exact estimates for given set of parameters.

### 1.4.2   Description of models to be compared

For our particular demonstrative task, which is restricted by the context of the study of Michaud, Llerena, and Joly (2012), we have already described the advantages and reasons behind the unrelenting theoretical assumptions concerning the behaviour of individual, as well as the dataset generation procedure. The two resulting datasets allow us to explore the effects of the random effects of the alternatives' attributes on the modelling. This possibility is particularly important, as usually researchers ignore the possibility of random effects presence in the population and use more simple and conventional multinomial logistic models to model various discrete choice situations. However, we are not going to test only one model over the obtained dataset, but rather introduce several models with different specifications in order to demonstrate a vast potential of our testing framework and its advantages for research.

As we are exploring an over-simplified framework, we are going to study first two different traditional

models each perfectly adapted to model one of the two generated datasets respectively. We are speaking about the multinomial logistic regression, which should yield perfect fit results on a fixed effects dataset and its counterpart - the mixed multinomial logistic regression, which should be the most performant in the presence of random effects in the utility functions. Many of the existing applied econometrics papers use the most simple specification of the Multinomial Logistic Regression (MNL), that may lead to erroneous results and conclusions in the presence of random coefficients. Eventually these models will allow us to verify, whether or not we are able to obtain the same results as at the input.

What is more, as the main objective of this work is to demonstrate proposed framework's flexibility, we are going to show how a completely alien model to econometrics, such as neural networks model, may be explored and compared with more traditional tools. More precisely, we are going to use a neural networks imitating the procedure of the multinomial logistic regression, while the other will be more traditional multilayer neural network. It is because this model can be viewed as an even wider generalisation of the generalised additive models (GAM), that it is possible to simulate a model similar to MNL and MMNL models. This choice was made because the seemingly identical model by its structure may produce different results, depending on the implemented estimation technique. The NN techniques offer us a great number of different algorithms which are more advanced than the algorithms traditionally implemented in econometrics, which make us wonder, whether the changes in the estimation algorithm will allow us to achieve better results.

In this part we will attempt as well to introduce some common notation for the different models' families, issued from different disciplines.

### 1.4.2.1 Logistic regressions

Multi-category logit models simultaneously use all pairs of categories by specifying the odds of outcome in one category instead of another (Agresti 2007). As described in Agresti (2013), many applications of multinomial logit models relate to determining effects of explanatory variables on a subject's choice from a discrete set of options.

**Multinomial Logit**

Even if in the original article of Michaud, Llerena, and Joly (2012) a Mixed Logit model is used, here we start our study with an introduction of the multinomial logistic regression (MNL) model, assuming the fixed effects presence. This model will allow us to contrast the performances in case of both fixed and random effect theoretical assumptions and compare them with a more advanced version of mixed multinomial logistic regression and NN model. This assumption is relaxed in the Mixed Logit model (ML or MMNL), where coefficients (or some of them) vary by individual (Agresti 2013). The logistic regression models are derived from GLM specifications (Agresti 2007):

$$g(\mu_i) = \sum_r \beta_r x_{ir} \tag{10}$$

25

Where $g(.)$ is a link function, which is a logistic transformation for binary logistic model. It is important to say that in this theoretical introduction we ignore in some extent the previously introduced terminology: $i$ still denotes the individual observations, laying in range of $\{1, \ldots, N\}$ in this case; the $r$ index here stands for different variables, because we do not use matrix notation for the reasons of simplicity.

Here we propose the econometric specification of a *multinomial logit (MNL)* model as described by Cascetta (2009). The MNL model is one of the simplest *random utility model (RUM)* (McFadden 1974). This class of models relies on the hypothesis, that an individual $n$ maximises his perceived utility over a set of alternatives $\Omega$, his utility determined by a fixed and a random parts, as described earlier:

$$U_{ij} = V_{ij} + \eta_{ij} \text{ where } V_{ij} = \alpha_j + \beta_j X_i + \gamma Z_j \tag{11}$$

Both $\beta$, representing the alternative specific individual coefficients, and $\gamma$, standing for population-wide attributes effects, are assumed to be fixed across population, meaning that all the individuals have identical preferences and are subject to identical effects. As precise in Agresti (2013) this approach enables discrete-choice models to contain characteristics of the chooser and of the choices. It offers the model an immense flexibility. The MNL is based on the assumption that the residuals $\eta_{ij}$ are identically and independently distributed (iid.) as Gumbel random variables with zero mean and scale parameter $\theta$, which is usually equal to 1 ($\theta = 1$). This calibration is done due to computational reasons, which will be explained later in this part.

One of the key concepts when it comes to modelling of the described above process is the *latent variable* notion. The latent variable $Y$ corresponds to its more meaningful counterpart $V$ and is sometimes understood as probability to choose a particular alternative. Obviously, as in the experimental context we are unable to observe the real choice probabilities, this variable takes values 0 or 1 depending on whether or not a particular alternative was chosen:

$$Y_i j = I(V_{ij} > V_{il} | j \neq l, \forall l \in \Omega_i) \tag{12}$$

Under the assumptions made here, the probability of choosing alternative $\omega_j$ from among those available $\{\omega_1, \ldots, \omega_k\} \in \Omega$ by individual $i$, can be expressed in closed form as:

$$P_{ij} = \frac{e^{V_{ij}/\theta}}{\sum_{l=1}^{k} e^{V_{il}/\theta}} \tag{13}$$

The probability structure incorporates the theoretical assumptions of the finite choice set, the uniqueness of the chosen alternative and the idea of utility maximisation. In a more comprehensive form, we may say that an individual chooses a particular alternative $\omega_j$ or simply $j$ among all available for him alternatives $\Omega_i$ only if its utility is higher than any others' alternative utility:

$$P_{ij} = P(\eta_{il} - \eta_{ij} < V_{ij} - V_{il}) \forall l : l \neq j, l \in \Omega_i \tag{14}$$

Knowing the structure of $V_{ij}$ and assuming the $\theta$ parameter for Gumble distribution of $\eta$ is 1 we may rewrite the probability as:

$$P_{ij} = \frac{e^{\alpha_j + \beta_j X_i + \gamma Z_j}}{\sum_{l=1}^{k} e^{\alpha_l + \beta_l X_i + \gamma Z_l}} \tag{15}$$

The alternative $\omega_j$ in such case is denoted as reference alternative or baseline alternative and is subject to several restriction for the sake of identifiability. The most important one is that we can not identify all the parameters in the probability function, which require us to impose some restrictions over effects structure. Traditionally (Agresti 2013) the reference level coefficients are assumed to be 0, reducing this way the number of parameters to estimate. This choice has some important consequences for the models' interpretation, because the estimated effects for other alternatives in this case should be treated as differences between the actual effects for the baseline alternative and other alternative respectively. The estimated parameters are in fact:

$$V_{ij} - V_{il} = (\alpha_j + \beta_j X_i + \gamma Z_j) - (\alpha_l + \beta_l X_i + \gamma Z_l) \tag{16}$$

Where $l \neq j$ and $j, l \in \Omega_i$. Which could be transformed into:

$$V_{ij} - V_{il} = (\alpha_j - \alpha_l) + (\beta_j - \beta_l)X_i + \gamma(Z_j - Z_l) \tag{17}$$

At this stage an important remark should be made, which concerns the understanding of individual characteristic effects and alternatives' attributes effects. It is theoretically possible to estimate a common individual effect for all the alternatives should we only wish to. The main idea lies in the correct parametrisation of the initial framework. To achieve identifiability for the individual characteristic specific effects we should observe enough within choice set variance, as otherwise the resulting singularity will incapacitate us to perform the estimation. In other words, we can understand this procedure as manually setting the individual effects to 0 for our baseline alternative and estimating the resulting model. Speaking about the changes in the dataset, the described above procedure is strictly equivalent to setting the baseline alternative's individual characteristics vector to zeros and estimating the resulting feature matrix as alternative specific attributes.

The traditional vision of alternative specific individual characteristics effects, assuming $\beta_j = 0$, is:

$$(\beta_j - \beta_l)X_i = -\beta_l X_i \text{ if } \beta_j = 0 \tag{18}$$

The analogous vision for alternatives' attributes effects, when reference attribute $Z_j$ is set to 0 is:

$$\gamma(Z_j - Z_l) = -\gamma Z_l \text{ if } Z_j = 0 \tag{19}$$

As we can see $\beta_l$ and $\gamma$ parameters are roughly equivalent in these two cases, assuming we are interested in means over the set of individuals $N$ and alternatives $\Omega$.

$$E_{il}(-\beta_l X_i) = E_{il}(-\gamma Z_l) \forall i \in N, \forall l \in \Omega \tag{20}$$

Which under transformation equals to:

$$- E_l(\beta_l) E_i(X_i) = -\gamma E_l(Z_l) \tag{21}$$

Assuming $X$ and $Z$ here is the same variable, varying across individuals and characteristics ($Z_j = 0$), we obtain that:

$$- E_l(\beta_l) X = -\gamma Z \Rightarrow E_l(\beta_l) = \gamma \tag{22}$$

This could be empirically confirmed through estimation of two different specifications and aggregation of obtained results.

However, were we in need to estimate an individual for all the alternatives except the baseline one, we could benefit from this transformation to do so. Such transformation allows us to take the multiple choice context of the expiremental setup.

**Mixed Multinomial Logit**

Following Agresti (2007) presentation, generalized linear models (GLMs) extend ordinary regression by allowing non-normal responses and a link function of the mean. The generalized linear mixed model, denoted by GLMM, is a further extension that permits random effects as well as fixed effects in the linear predictor. We begin with the most common case, in which is an intercept term in the model.

$$g(\mu_i) = \sum_r \beta_{ir} x_{ir} \tag{23}$$

Where $\beta_i$ is issued from some multivariate distribution. Traditionally this distribution is assumed to be a multivariate normal distribution (MNV) giving:

$$\beta_i \sim MNV(\beta, \Sigma) \tag{24}$$

In more recent work of Agresti (2013) the more advanced models are described. The multinomial logit and probability based discrete-choice models can be further generalized by treating certain effects as

random rather than fixed.

A mixed logit model is the one in which choice probabilities are obtained by integrating the logistic expression for choice probabilities with respect to a distribution for certain model parameters. This allows heterogeneity among subjects in the size of effects. It is useful as a mechanism for inducing positive association among repeated responses with panel data. Estimates of the parameters of the mixing distribution provide information about the average effects and the extent of the heterogeneity. Individual effects can also be predicted using this technique.

The Mixed Logit is a further development and generalisation of a traditional MNL and Conditional Logit models, because both of these models may be constructed using Mixed Logit specification with a correct parametrisation. The main difference from the more simple models is that in this case it is assumed that effects vary across population and might even be correlated. The utility specification in this case is constructed identically to simple models, but the deterministic part assumes that effects vary across population:

$$U_{ij} = V_{ij} + \eta_{ij} \text{ where } V_{ij} = \alpha_j + \beta_j X_i + \gamma_i Z_j \tag{25}$$

Mathematically the random effects specification is achieved through the parameter vector $\gamma_i$, which is unobserved for each $i$. The $\gamma$ in this case is assumed to vary in the population following the continuous density $f(\gamma_i \mid \theta)$, where $\theta$ are the parameters of this distribution. The simplest choice of the distribution for the random effects is the normal distribution, which was used by Michaud, Llerena, and Joly (2012), or more precisely a multivariate normal distribution, because authors took into account the correlation between coefficients:

$$\gamma_i \sim MVN(\gamma, \Sigma) \tag{26}$$

In this case the vector of alternative specific effects can be represented as:

$$\gamma_i = \gamma + L\sigma_i \tag{27}$$

Where $\sigma_i \sim N(0, I)$, and $L$ is the lower-triangular Cholesky factor of $\Sigma$ knowing which, the actual variance-covariance matrix for random effects can be derived, as presented in Croissant (2020):

$$LL^T = V(\gamma_i) = \Sigma \tag{28}$$

Here we do not present the eventual possibility to incorporate the individual specific characteristics covariates into the given framework, because we will not use it, but such possibility is definitely worth mentioning.

Where $\beta$ are some fixed mean effects across population and $\psi$ stand for the random part with 0 mean

and some imposed variance-covariance structure, as it is technically possible to assume that only some of the effects are random.

A more advanced description of MMNL models is available in the work of McFadden and Train (2000), where some intuitions are given on the estimation techniques necessary to evaluate such complex model. The authors suggest, that numerical integration or approximation by simulation is needed to evaluate MMNL probabilities. Maximum Simulated Likelihood (MSLE) or Method of Simulated Moments (MSM) could be used to estimate the MMNL model in practice, both of which are described in the reference work (McFadden and Train 2000)

#### 1.4.2.2 Neural Networks

The second group of models focuses on more advanced and atypical modelling techniques rarely implemented by the economists in their studies, as usually this family is perceived as not offering enough insight when it comes to the effects estimation. The ML techniques are usually viewed by economists as some black boxes, which do not provide any information about the underlying process. It is quite easy to comply with their position, as even though the most advanced techniques perform better in terms of predictive power, they rarely offer any insight into the modelling process.

For this particular part we use the model's specifications described in the handbook of Hastie, Tibshirani, and Friedman (2009) with some additions and modifications, which aim at integration of this particular specification in conformity with the specifications of the econometric discrete model notation. *Neural Networks (NN)* represent an advanced class of models, being a further complexification of the *generalised additive models (GAM)*, which are a generalisation of the *generalised linear models (GLM)*, which was defined in previous subsection. This GLM is generalised through assumption that each explicative variable in $X$ can undergo some transformation, linear or not, resulting in a following GAM model:

$$g(\mu_i) = \sum_r s_r(x_{ir}) \tag{29}$$

Where $s_r(.)$ is an unspecified smooth function of predictor $x_{ir}$. In order to obtain a NN model, this structure is further developed as follows to obtain firstly a *projection pursuit regression (PPR)*:

$$f(X) = \sum_{r=m}^{M} g_m(\omega_m^T X) \tag{30}$$

The $X$ in this notation is a vector of inputs with $p$ components, and $\omega_m$ with $m \in \{1, 2, \ldots, M\}$ are unit $p$-vectors of unknown parameters. Before proceeding, we will introduce some novelties to the notation used till this point by introducing vectors $X1, X2, \ldots, XS$, where $X1$ is the output of the first layer of neural network, each element of which is some transformation (usually linear in parameters with

30

some "activation" function) of the input vector $X$. Then the simplest NN for $\Omega$ alternatives (classes) classification, with two layers, may be represented as:

$$f_j(X) = g_j(X2) \text{ with } X2_j = \psi_{0j} + \psi_k^T X1 \tag{31}$$

Where $f_j$ models the probability of a class $j$, or in more comprehensive language the probability that a given individual will choose an alternative $\omega_j$ from his choice set $\Omega_i$:

$$X1_m = \sigma(\phi_{0m} + \phi_m^T X) \tag{32}$$

While $\sigma(.)$ is an activation function and $g_k(.)$ a probability transformation function, traditionally a *softmax* function. The latter is being used as well in *multinomial logit (MNL)* models:

$$g_j(T) = \frac{e^{T_j}}{\sum_{l=1}^{\Omega} e^{Tl}} \text{ where } j, l \in \Omega \tag{33}$$

This means, that single level NN with a softmax activation layer should be identical to simple MNL model with all the coefficients varying by alternatives. $Z_m$ can be viewed as a basis expansion of the original inputs $X$ and the neural network is then a standard *linear multinomial logit (MNL)* model, using the transformations as inputs.

One of the supposed major problems for NN models in discrete choice context is the inability to take into account all the influencing factors across all the alternatives. Moreover, in this case study there is major drawback in the ambiguity among choices A and B, as they are interchangeable.

As we desire to obtain the effects assuming the alternatives A and B are identical, this means that we should impose some additional restrictions over the model. Traditional Multinomial Logistic regression (MNL) can be potentially transcribed into a NN using convolution techniques. The convolution layer operates iteratively on a given subset from the input vector, calculating one single output per $k$ inputs. In this case $k$ is denoted *kernel size*. Another parameter, which defines a convolutional layer is the *stride* ($s$), which determines how the "window" determined by kernel size should be moved over the input layer. Consequently, the output layer consists of $m$ values determined as:

$$m = \frac{n - k}{s} + 1 \tag{34}$$

Where $n$ is the length of the input vector to this layer. We may attempt to define a convolution layer with linear activation function as follows, assuming $X = X_1, \ldots, X_n$ is the input vector and $X1_1, \ldots, X1_m$ is the output vector, while $\phi = \phi_1, \ldots, \phi_k$ is the vector of weights:

$$X1_1 = \phi_1 X_1 + \phi_2 X_2 + \cdots + \phi_k X_k$$
$$\vdots \tag{35}$$
$$X1_m = \phi_1 X_{n-k} + \cdots + \phi_k X_n$$

The designed this way CNN consists of two transformation layers. The first one is 1D convolutional layer with linear activation function, which takes as input the dataset in "wide" format with 27 variables overall (9 variables for each alternative), which produces a single value as an output value for each individual for each choice set, resulting in 3 output values in total. The second layer is a restricted softmax transformation layer, which directly applies softmax transformation over the inputs, without any supplementary permutations.

The vector of inputs issued from the dataset transformed into the "wide" format can be represented as:

$$X_i \quad = \quad Buy_{i,A}, Sex_{i,A}, Age_{i,A}, \ldots, Habit_{i,C}, Price_{i,C}, Label_{i,C}, Carbon_{i,C}, LC_{i,C} \tag{36}$$

Where all values with $C$ index are set to zero in order to set the baseline alternative. The first convolutional layer can be written as:

$$V_j = \alpha_{Buy} Buy_{ij} + \beta_{Sex} Sex_{ij} + \beta_{Age} Age_{ij} + \beta_{Income} Income_{ij} + \beta_{Habit} Habit_{ij} +$$
$$+ \gamma_{Price} Price_{ij} + \gamma_{Label} Label_{ij} + \gamma_{Carbon} Carbon_{ij} + \gamma_{Label \times Carbon} Label \times Carbon_{ij} \tag{37}$$

Where $j \in \{A, B, C\}$, with $C$ denoting the "No buy" option.

We configure the convolution layer with linear activation function to move across the input vector with strides 9, producing this way a vector of length 3 as an output. This outputs of this layer may be interpreted as utilities for each alternative respectively, identically to MNL regression. The resulting design for a single convolution fold can be schematically represented as in figure 4.

The second transformation layer is a dense layer with a "softmax" activation function as described above, which has 3 coefficients for each output, because it aggregates the inputs to an identical number of outputs rescaling them in the process and producing choice probabilities. Taking a set of $V_A, V_B, V_C$ for inputs and producing a vector of probabilities $P(A), P(B), P(C)$ as outputs. The second level may be synthetized as presented in figure 5.

Finally, given the combination of these two layer we may construct the whole CNN model. We may use the following graphical representation, shown on figure 6 to visualise the resulting CNN architecture:
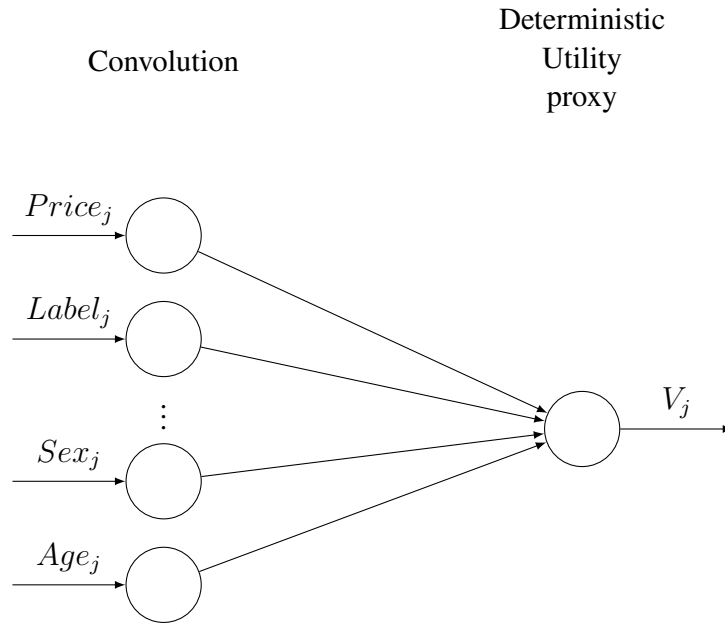
Figure 4: Convolution layer

Convolution

Deterministic
Utility
proxy

$Price_j$

$Label_j$

$Sex_j$

$Age_j$

$V_j$

Figure 5: Softmax Layer

Deterministic
Utility
proxy

Probability

$V_A$
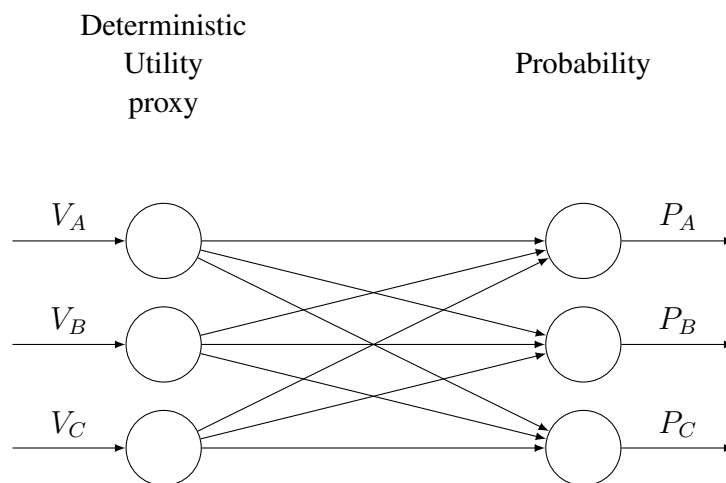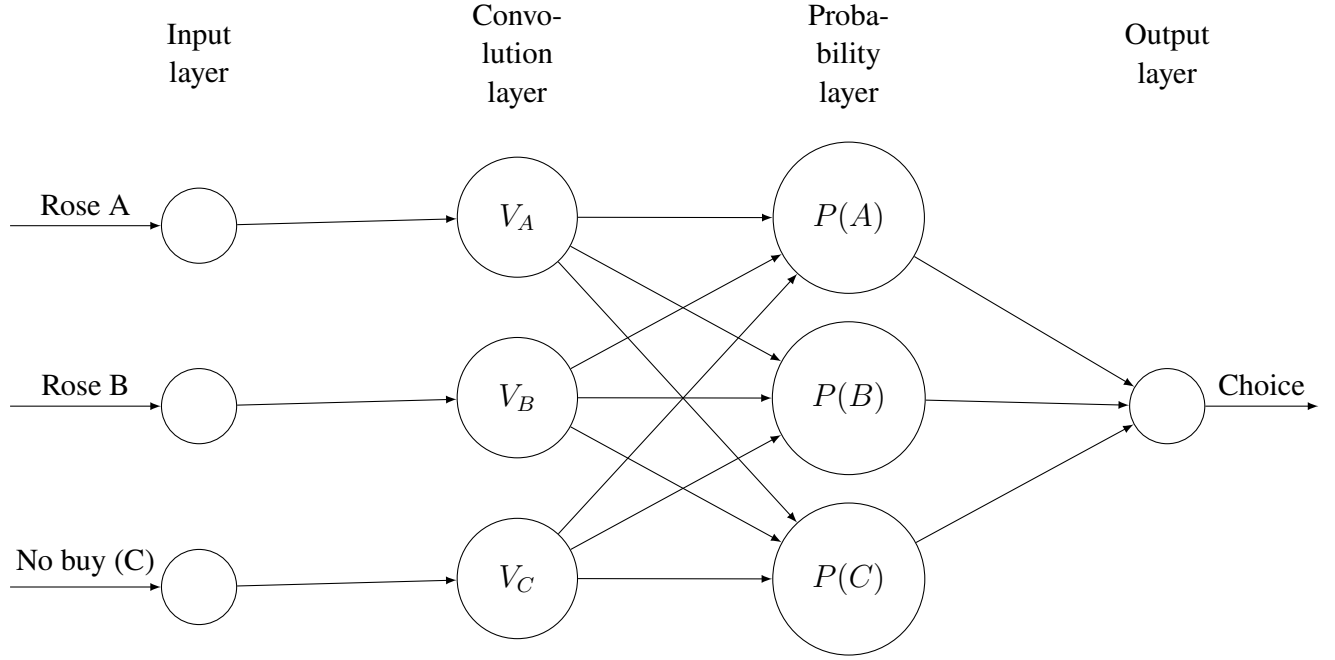
$V_B$

$V_C$

$P_A$

$P_B$

$P_C$

Figure 6: Convolution Neural Network design



The figure 6 is no more than a simplified architecture presentation for the chosen CNN design, imitating the MNL model in this particular case. Each alternative input on this graph assumes entry of the three attributes of a particular alternative, supported by five individual characteristics each, the later being specific to a particular alternative exactly as in the MNL model specification.

In this case the only difference between these two models is represented by the algorithm used for estimation, which can yield absolutely different results or even require some transformation of the input dataset (ie. rescaling, which is used to prevent biases in weights estimation). Consequently, the main interest of such implementation is to observe, whether or not a ML algorithm will be able to bypass the MNL model performances in the presence of heterogenous individual preferences. Different convergence rates and different iterative algorithms may result in absolutely distinct optimums for the parameters vector. The particular algorithms implemented will be discussed later, alongside the obtained results.

For NN modelling we use the advanced interface offered by Google's *Tensorflow* (Allaire and Tang 2020) with *Keras* (Allaire and Chollet 2020) back-end for *R*-language. The flexibility offered by this particular tool is astonishing compared to other neural networks implementations in proposed in *R*. This flexibility allows us to simulate exactly the architecture of a MNL model and compare this way how the different estimation techniques and algorithms perform in the identical contexts.

## 1.5    Model performance evaluation and available measures

In this subsection we are going to describe the different performance measures, attempting at the same time to shun some light on the complexity of this particular task and the multitude of different questions that are usually aborded when a problem of performance measures' choice arises.

The main problem in the case of classification context and particularly in the multiple choice classification context relates to the fact that rarely all of the models can use the same metrics for their comparison (Baldi et al. 2000). The available metrics largely depend on the output variable type, the models architecture and assumptions, the specifications, the algorithms used and, finally and most importantly, the context. As we have seen earlier, the work of Michaud, Llerena, and Joly (2012) was focused on the identification of the willingness to pay of consumers for particular environmental attributes of roses, rather than general goodness of fit of particular model, which perfectly illustrates the complexity of the posed question.

There exists a multitude of different target metrics to evaluate and compare the performances of different models. For example, one may be interested in exploration of a particular effects or the overall goodness of fit, some predictive qualities or a possibility to derive correct estimates for a particular socio-economic information. This topic was already largely explored by some of the statisticians (Japkowicz and Shah 2011) with some initial steps into producing an integrated support containing all the necessary information for applied studies. However, even given the amount of the work in reference, there is still a strong need for contextualisation and constitution of application specific methodological supports. The different possible application scenarios require sometimes absolutely different metrics. For example, econometricians rarely take into account the computational efficiency of the models, while ML researchers are rarely considering the possibility to derive the specific field specific metrics.

Nevertheless, this work aims at demonstrating the full potential of the proposed experimental framework and we are bound to demonstrate at least a fraction of its full potential, which inevitably addresses the different performance metrics used to compare the models' performance in terms of precision and predictive accuracy.

The measures available may roughly be divided into three parts following the logic of Japkowicz and Shah (2011) (for an adaptation of the vision of Japkowicz and Shah (2011) on the different measures' types see Appendix B).

- The measures that take information solely from the *confusion matrix*, which can be calculated using the estimated model over a know dataset (also denoted a test dataset). These measures are typically applied in the case of deterministic classification algorithms, but can be calculated for the probabilistic output algorithms as well.
- The measures that not only use the confusion matrix, but integrate the information about the class distribution priors and classifier uncertainty. Logically, these metrics are useful for the *scoring* classifiers" performance evaluation and could not be used with some more simple models.

- Bayesian measures to account for probabilistic classifiers and measures for regression algorithms. Bayesian measures require a probabilistic structure of the models output.

The measures may be as well separated into two different groups by their behaviour (Japkowicz and Shah 2011):

- A *monotonic* performance measures $pm(.)$, for which a strict increase (or decrease) in the value of $pm(.)$ indicates a better (or worse) classifier throughout the range of the function $pm(.)$ respectively.
- *Not strictly monotonic* can be thought of as the class-conditional probability estimate discussed in Kukar and Kononenko (2002) in the context of a multi-class problem.

As our framework can potentially treat multiple different aspects, we will not only assess the general models' performances, but explore the capacity to identify and estimate the target values of interest. Taking into account the context of the target article we will be mostly interested in exploring the willingness to pay (WTP) or the premium, that the consumer is ready to add to the observed price for a particular attribute.

### 1.5.1 Confusion matrix

Most of the performance measures for a classification task are derived from the observed entries in the confusion matrix, denoted $C$ (Japkowicz and Shah 2011, @baldi2000ar). This matrix lies in the center of most non-probabilistic performance measures for classification. A confusion matrix $C$ for a classifier defined by a function $f(.)$ over some dataset may be defined as:

$$C = c_{ij}, \ i, j \in \{1, 2, \ldots, k\} \tag{38}$$

Where $i$ is the row index and $j$ is the column index, both referring to some available alternatives for a given alternatives' set $\Omega$.

Generally, $C$ is defined with respect to some fixed learning algorithm. The confusion matrix can be extended to incorporate information for the performance of more than one algorithm, resulting in creation of a *confusion tensor*, which can be imagined as a stack of matrices. There exist specific metrics to be implemented on such *tensor*.

Given a training dataset and a test dataset, an algorithm learns on the training set, outputting a fixed classifier $f$. These datasets may be identical, as it is frequently done in economics studies. The test-set performance of $f$ is then recorded in the confusion matrix. This means that a confusion matrix, as well as its entries and the measures derived from these are defined with respect to a fixed classifier $f$ over a given dataset. Consecutively, the matrix is sometimes denoted with respect to $f$ as $C(f)$. It is a square

$k \times k$ matrix for a dataset with $k$ classes. Each element $c_{ij}(f)$ of the confusion matrix denotes the number of examples that actually have a class $i$ label and that the classifier $f$ assigns to class $j$.

In binary case these measures are simplified to four, that do not always appear in matrix form for the sake of simplification. These measures, as well as derived performance indicators are described in Baldi et al. (2000). The binary classification case is the most common setting in which the performance of the learning algorithm is measured. Also, this setting serves well for illustration purposes with regard to the strengths and limitations of the performance measures.

### 1.5.2 General performance measures

The general measures (Baldi et al. 2000) describe the performance of a given classifier $f(.)$ (or shortly $f$) over a given set of observation, taking into account all the possible classes, or choices in the discrete choice context. In other words, these measures incorporate all the information available for all the classes matches or mismatches, which offers some good general overview of a given model performances, but sometimes ignores some of the significant elements. For example, given an unbalanced dataset, where one class dominates the other, the general performance measures can have high positive values, signalling the good overall performance, while all the observations will be assigned to dominant class by the classifier.

The most known measures, which are usually implemented to assess the general performance of the algorithms or even construct "loss" functions for some learning tasks include: the empirical risk, the empirical error rate and the accuracy.

Accuracy and error rate effectively summarize the overall performance, taking into account all data classes. This is the reason why these measures are often implemented to assess general algorithms' performances and are used in the learning tasks. Moreover, they offer an insight into the generalization performance of the classifier by means of studying their convergence behaviours, which may be important for some algorithms.

Nevertheless, such general metrics have potential limitations (Japkowicz and Shah 2011). Firstly, these measures suffer from the lack of information on the varying degree of importance of different classes on the performance. What is more, as we have already pointed out, the metrics are incapacitated by the lack to produce any meaningful information in the case of skewed class distribution. This results in the situation, when as the distribution begins to skew in the direction of a particular class, the more-prevalent class dominates the measurement information in these metrics, making them biased.

**Empirical risk**

The *empirical risk* $R_N(f)$ of classifier $f$ on test set $N$, defined as:

$$R_N(f) = \frac{1}{\mid N \mid} \sum_{i=1}^{|N|} I(y_i \neq f(x_i)) \tag{39}$$

Where:

- $I(a)$ is the indicator function if predicate $a$ is true and zero otherwise;
- $f(x_i)$ is the label assigned to example $x_i$ by classifier $f$;
- $y_i$ is the true label of example $x_i$, which indicates to ome of the alternatives $\{\omega_1, \ldots, \omega_k\} \in \Omega$;
- $\mid N \mid$ is the size of the test set.

This measure describes the average loss over the data points.

**Empirical error rate**

The *empirical error rate* can be computed as follows:

$$R_N(f) = \frac{\sum_{i,j:i \neq j} c_{ij}(f)}{\sum_{i,j=1}^{\Omega} c_{ij}(f)} = \frac{\sum_{i,j=1}^{\Omega} c_{ij}(f) - \sum_{i=1}^{\Omega} c_{ii}(f)}{\sum_{i,j=1}^{\Omega} c_{ij}(f)} \tag{40}$$

This rate measures the part of the instances from the given set that are incorrectly classified by the learning algorithm $f$.

**Accuracy**

The *accuracy* describes the part of correctly classified instances in a given set and is by its nature a complement to the empirical error-rate measure. It can be computed as:

$$Acc_N(f) = \frac{1}{\mid N \mid} \sum_{i=1}^{|N|} I(f(x_i) = y_i) \tag{41}$$

Where $y_i$ is the observed class for observation $i$. Given a skew ratio $r_s$, it is possible to extend this measure and define the *skew-sensitive formulation of the accuracy*. Such modification allows partially to solve the poor measures' utility problem on a skewed class distribution dataset.

### 1.5.3 Single-class performance measures

Apart from the general performance measures, there exist some more specific performance measures, which instead of estimating the performances of the overall classifier, target some specific aspects. Usually in the modelling the consumer behaviour we may be interested in his his choice "Buy" against "No buy" beforehand, and only afterwards we are interested by his consumer habits and preferences. Among these measure we may cite:

- True- and False-Positive/Negative Rates
- Specificity

- Sensitivity
- Precision

- Recall
- Geometric means
- Likelihood Ratio (LR) [2]

- F-measured

- Skew and Cost

One of the important problems for discrete choice modelling and general classification tasks resides in the form of the greater importance of the algorithms' performance on a single class of interest. This performance on a given class can be crucial with regard to the instances of this class itself or with regard to the instances of other classes in the training data. As it was mentioned earlier, in our particular study case, we may be interested at how good the algorithm distinguishes the "Buy" and "No buy" choices.

A number of such measures can also allow us to measure the overall performance of the classifier with an emphasis on the instances of each individual class. Such precise metrics may be excessive, given a particular case study, although they offer a good substitute for more typical measures, such as the accuracy or error rate.

In this part we are going to introduce some new terminology, because contrary to the precious parts, where we had to deal with classes, here we are bound to simplify the problem to a binary case. This means that one of the classes is considered as "positive", while the rest of the alternatives is regrouped into a single "negative" class. Such transformation allows us to define new variables, which will be used later in the class-specific measures presentation. Among these values we have:

- True Positive or $TP$, which denotes the number of correctly classified observations which appertained to the "positive" class;
- True Negative or $TN$, where the number of correctly classified "negative" instances is regrouped;
- False Positive or $FP$ stands for the misclassified instances that in the dataset were encoded as "positive" class;
- False Negative or $FN$, which logically indicates the number of initially "positive" observations, which were identified as "negative" ones by the model.

All these values may be easily obtained from the confusion matrix $C$.

**True- and False- positive/negative rates, specificity and sensitivity**

The most natural metric aimed at measuring the performance of a learning algorithm on instances of a single class is arguably its *true-positive rate*. The *true-positive rate* of a classifier is also referred to as the *sensitivity* of the classifier. The complement metric to this, in the case of the two-class scenario, would focus on the proportion of negative instances is called the *specificity* of the learning algorithm. It is obtained as:

---

[2]This measure will be omitted in order to prevent the eventual confusion with Likelihood Ratio (LR) used in the MNL and MMNL models

$$TPR_i(f) = \frac{c_{ii}(f)}{\sum_{j=1}^{I} c_{ij}(f)} = \frac{c_{ii}(f)}{c_i(f)} \tag{42}$$

The *false-positive rate* of a classifier:

$$FPR_i(f) = \frac{\sum_{j:j\neq i} c_{ji}(f)}{\sum_{j,k:k\neq i} c_{jk}(f)} \tag{43}$$

Some usefull derived formulas, which are easy to compute for a binary case, are introduced hereafter. The True- and False- positive rates:

$$TPR(f) = \frac{TP}{TP + FN} = \text{Sensitivity} = 1 - FNR(f) \tag{44}$$

$$FPR(f) = \frac{FP}{FP + TN} \tag{45}$$

As well as their counterpart, the True- and False- negative rates, which are focussed on the number of correctly classified instances from a "negative class".

$$TNR(f) = \frac{FN}{TN + FP} = \text{Specificity} \tag{46}$$

$$FNR(f) = \frac{FN}{FN + TP} \tag{47}$$

**Precision and recall**

The *precision* or *positive predictive value (PPV)* of a classifier $f$ on a given class of interest $j$, denoted as well as the "positive" class, in terms of the entries of $C$, measures how *precise* the algorithm is when identifying the examples of a given class and is defined as:

$$PPV_i(f) = Prec_i(f) \frac{c_{ii}(f)}{\sum_{j=1}^{I} c_{ji}(f)} = \frac{c_{ii}(f)}{c_{.i}(f)} \tag{48}$$

For binary case we can write the following simplified definition, which should be more clear to the reader:

$$Prec(f) = PPV(f) = \frac{TP}{TP + FP} \tag{49}$$

The PPV can be complimented with the sensitivity of the classifier over this class. This measure is generally referred to as *recall*:

$$Rec(f) = \frac{TP}{TP + FN} \tag{50}$$

**Geometric means**

The *geometric means* take into account the relative balance of several performance measures for a given classifier. The most popular option is to observe simultaneously the classifier's performance on both the positive and the negative classes:

$$Gmean_1(f) = \sqrt{TPR(f) \times TNR(f)} \tag{51}$$

This implementation is of particular interest for our case study, as we will be able to compare the performances of different models across "Buy" and "No buy" options. Another popular version of the measure, which focusses on a single class of interest, can take the precision of the classifier in combination with the classifiers performance on the "positive" class into account:

$$Gmean_2(f) = \sqrt{TPR(f) \times Prec(f)} \tag{52}$$

**F-measure**

The *F-measure* as well attempts to address the issue of convenience brought on by a single metric versus a pair of metrics. It combines the information of precision and recall in a single value. More precisely, the F-measure is a weighted harmonic mean of precision and recall, with a weight $\alpha$:

$$F_\alpha = \frac{(1 + \alpha)(Prec(f) \times Rec(f))}{\alpha Prec(f) + Rec(f)} \tag{53}$$

For instance, the most comprehensive *balanced F-measure* weights the recall and precision of the classifier evenly:

$$F_1 = \frac{2(Prec(f) \times Rec(f))}{Prec(f) + Rec(f)} \tag{54}$$

In most practical cases, appropriate weights are generally not known, which results in some complications in choice of the hyper-parameter $\alpha$ of such combinations of measures.

**Class ratio**

*Class ratio* for a given class $i$, which in the consumer choice setting is usually denoted $\omega_i$ refers to the number of instances of class $i$ as opposed to those of other classes in the dataset:

$$ratio_i = r_i = \frac{\sum_j c_{ij}}{\sum_{j,j \neq i} c_{ji} + \sum_{j,j \neq i} c_{jj}} \tag{55}$$

Or for a binary case:

$$ratio_{positive} = \frac{(TP + FN)}{(FP + TN)} \tag{56}$$

Another issue worth considering when looking at misclassification is that of classifier uncertainty. This lack of classifier uncertainty information is also reflected in all the performance measures that rely solely on the confusion matrix.

### 1.5.4   Information-theoretic measures

These measures are probabilistic by their nature, as they explore the performances of the classifier with respect to the (typically empirical) prior distributions of the data. in contrast to the cost-sensitive metrics that have been introduced earlier, the *information-theoretic measures*, because of accounting for the data priors, are applicable only to probabilistic classifiers. What is more, these metrics are independent of the cost considerations and can be applied directly to the probabilistic output of a given model. These measures are extensively implemented in Bayesian learning and take their roots in physics. Among these metrics one may encounter:

- Kullback–Leibler Divergence, which estimates the difference between the entropies of the two distributions;
- Kononenko and Bratko's Information Score, which explores the likelihood of correct classification.

In this work we will present only the first among these two.

**Kullback–Leibler Divergence**

Let the true probability distribution over the labels be denoted as $p(y)$. Let the posterior distribution generated by the learning algorithm after seeing the data be denoted by $P(y \mid f)$. Because $f$ is obtained after looking at the training samples $x \in S$, this empirically approximates $P(y \mid x)$, the conditional posterior distribution of the labels. Then the *Kullback–Leibler divergence* (KLD or KL) can be utilized to quantify the difference between the estimated posterior distribution and the true underlying distribution of the labels:

$$KLD[p(y) \mid\mid P(y \mid f)] = \int p(y) ln p(y) dy - \int p(y) ln P(y \mid f) dy \tag{57}$$

$$KLD[p(y) \mid\mid P(y \mid f)] = - \int p(y) ln \frac{P(y \mid f)}{p(y)} dy \tag{58}$$

42

The KLD divergence basically just finds the difference between the entropies of the two distributions $P(y \mid f)$ and $p(y)$. This measure is also known as *relative entropy* (see Baldi et al. (2000)) for more information.

$$KLD[p(y) \mid\mid P(y \mid f)] = -\sum_{x \in S} p(y) ln \frac{P(y \mid f)}{p(y)} dy = \sum_{x \in S} p(y) ln \frac{p(y)}{P(y \mid f)} dy \tag{59}$$

The KLD value is equal to zero if and only if the posterior distribution is the same as the prior, when the perfect fit is achieved, meaning that the classifier perfectly mimics the true underlying distribution of the labels.

Even though the KLD measures the difference between the posterior distribution obtained by the learner from the true distribution so there is a significant drawback to it. The KLD needs the knowledge of the true underlying prior distribution of the labels, which is rarely, if at all, known in any practical application. In practice the estimated priors are used, although in the experimental framework where a synthetic dataset is used, we may theoretically impose some "true" structure over the choice distribution.

### 1.5.5   Case specific metrics

The article of Michaud, Llerena, and Joly (2012) focuses on the WTP for roses and derivation of the premiums for particular alternative attributes of interest. This focus allows authors to explore the consumer attitude towards the alternative specific environmental attributes. Consequently, as we try to follow the logic introduced in the article, we are going to attempt to derive the WTP and premiums for attributes as well. However, before introducing the notion of the WTP and premium, we should firstly describe the procedure of derivation of the marginal effects, as the WTP and premiums are expressed using the marginal effects.

In the conventional MNL models the coefficients $\beta_{rj}$ can be interpreted as the marginal effect of variable $X_r$ on the log odds-ratio of alternative $j$ to the baseline alternative. The marginal effect of $X_r$ on the probability of choosing a specific alternative $j$ can be expressed as:

$$ME_{rj} = \frac{\Delta P(Y_i = \omega_j)}{\Delta X_r} \tag{60}$$

Consequently, for the MNL model, the marginal effect of $X_r$ on alternative $j$ not only takes into account the parameters specific to $j$ alternative, but the ones of all other alternatives as well. The equation can be written in this case as:

$$\frac{\Delta P(Y_i = \omega_j)}{\Delta X_i} = P(Y_i = \omega_j)[\beta_{j1} - \sum_{l=0}^{k} P(Y_i = \omega_l)\beta_{j1}] \tag{61}$$

The parameters such as WTP and premiums are more easy to interpret. They can be estimated directly or can be obtained from the marginal utility by dividing it by the effect estimate of a price, taken as a non random parameter. The resulting ratio can afterwards be interpreted as a monetary value. The WTP as it was described in the context presentation, taking into account the case specific relative utility functions can be represented as:

$$WTP = \frac{\frac{\Delta V}{\Delta BUY}}{\frac{\Delta V}{\Delta Price}} = \frac{-\alpha_{Buy}}{\beta_{Price}} \tag{62}$$

The premiums for a given attribute $X_r$ ($Label$, $Carbon$ or their cross-product $LC$), can therefore be expressed as:

$$WTP = \frac{\frac{\Delta V}{\Delta X_r}}{\frac{\Delta V}{\Delta Price}} \tag{63}$$

### 1.5.6   Selection of measures to implement

In this work we are going to explore only a selection of the described above most popular performance metrics, that are the most interesting given the context of the study. Moreover, in our application we are limited in the number of measures we can explore.

In the first place we are interested by the WTP for roses and the premiums associated with particular alternative specific attributes. These theoretical values could be easily derived for all the three explored models and they will allow us to compare, how close are the derived values from the theoretical input values, which were defined on the dataset generation step.

Secondly, it is important to assess the overall goodness of fit over the whole dataset for the selected models. For this particular task the most suited measure is the *accuracy*. This way we will be able to observe the ratio of the overall correctly classified instances. We may implement the KLD estimator for overall goodness of fit, based on the probability distributions, because all the models predict the probabilities for the available alternatives.

We may be interested as well in comparing the performances of the given models in terms of distinguishing the "Buy" choice, irrelevant of the alternative, and the "No buy" choice. This is a particularly interesting question, because in the different choice settings and over the datasets generated under different theoretical assumptions. For this purpose the most interesting choice will be to select the F-measure or a Geometric mean of the TPR and TNR.

Finally, we are going to observe the performance of these different models in terms of computational efficiency in resources consumption. For this task we will observe the computation times for given models[3]. The obtained results will be discussed at the end of this work.

---

[3]This measure is one of the most complex, because it accounts at the same time for different models, different estimation

44

# 2    Model comparison in practice: an application

This section is designed to present the results of the designed theory-testing framework as well as to offer some more detailed view on the adopted procedure. It will respect the following structure. First of all we will start by a presentation of two generated datasets and their comparison with the original dataset obtained through a controlled experiment by Michaud, Llerena, and Joly (2012). Then we will discuss the technical implementation of the models to test and the resulting estimations over two of the simulated datasets. Finally the target performance metrics will be constructed for all the models' performances over both of the datasets and we will compare the obtained estimates with the input values, assessing this way the biases suffered during estimation.

## 2.1    Simulating individual choices

Based on the article of Michaud, Llerena, and Joly (2012) we generate a synthetic dataset assuming the utility function is as described in the paper with some minor changes and adjustments. We have already delimited the scope of study and delimited our area of interest to the exploration of different models performance given the theoretical structure of consumer preferences for the alternative specific attributes. For simplicity we relax some of the assumptions made in the paper in order to generate two different datasets. For the first dataset we assume that estimations made in the paper and the derived utility functions are correct and reflect the real world situation. For the second one, we relax some of the advanced assumptions and regenerate a simplified version, which will allow us to contrast the performances of different models in different choice context assuming different nature of choice functions.

In both situations the utility functions are determined as in paper: we use the exact means for the coefficients estimates, assuming they are correct. The relative utility's deterministic part for each individual is defined by the following function, which was presented in a more detailed way in previous section:

$$
\begin{aligned}
V_{ij} = \alpha_{i,Buy} + \beta_{Buy,Sex}Sex_i + \beta_{Buy,Age}Age_i + \beta_{Buy,Salary}Salary_i + \beta_{Buy,Habit}Habit_i + \\
+ \gamma_{Price}Price_{ij} + \gamma_{i,Label}Label_{ij} + \gamma_{i,Carbon}Carbon_{ij} + \gamma_{i,LC}LC_{ij} \quad (64)
\end{aligned}
$$

Where $LC = Label \times Carbon$. The random component of the relative utility $U_{ij}$ is defined as identically and independently distributed random variable $\epsilon_{ij}$ issued from the Gumble distribution parametrised with $(0, 1)$. The mean effects for the components of the deterministic part are given as presented in the table 4a

The only difference between the two generated datasets is in the specification of the randomness of

algorithms, different numerical implementation in the statistical software and different PC configuration. It is valid in this particular case, because all models were estimated using the same hardware and software set-up.

Table 4: The assumed relative utility function parameters

(a) Mean effects

| | *Effects* |
|---|---|
| | *Means* |
| **Individual characteristics ($\beta$)** | |
| Sex | 1.420 |
| Age | 0.009 |
| Salary | 0.057 |
| Habit | 1.027 |
| **Alternatives' attributes ($\gamma$)** | |
| Price | $-1.631$ |
| Buy | 2.285 |
| Label | 2.824 |
| Carbon | 6.665 |
| LC | $-2.785$ |

(b) Variance-covariance structure

| | *Effects* | |
|---|---|---|
| | Fixed | Random |
| **Variance** | | |
| Buy | 0 | 3.202 |
| Label | 0 | 2.654 |
| Carbon | 0 | 3.535 |
| LC | 0 | 2.711 |
| **Covariance** | | |
| Buy:Label | 0 | -0.54 |
| Buy:Carbon | 0 | -4.39 |
| Buy:LC | 0 | 6.17 |
| Label:Carbon | 0 | 8.77 |
| Label:LC | 0 | -2.33 |
| Carbon:LC | 0 | -4.82 |

these coefficients as they may vary or not across population. It means, that the first dataset is generated assuming the variance-covariance matrix for correlated random coefficients is composed with 0's only and the resulting multivariate normal distribution produces exact means for the coefficients. The second dataset is generated using the exact estimates of the variance-covariance matrix as provided in the article. The assumed parameters for effects distributions are represented in the table 4b.

Additionally we impose some supplementary constraints to our data due to the limitations of the simulation tool. Particularly, the individual characteristics are supposed to be not correlated, which can be explained by the fact that the context of a controlled experiment offers a possibility to control this particular feature. Obviously, this is not optimal decision, as naturally the age, sex, income and environmental habits of individuals should be correlated. Unfortunately, the original article does not provide information about the characteristics' variance-covariance matrix.

The targeted features and requirements to the resulting dataset are numerous and they make a contrast compared to the initial empirical dataset.

The simulated dataset allows us to explore significant number of choice sets for numerous artificial individuals, which ensures statistical validity for obtained results and permits us to use advanced estimation algorithms (such as neural networks, for example). It means that we generate a large sample with exhaustive number of choice sets, in which all the possible combinations of alternative attributes are represented. Here by *attributes* we understand the binary factors describing rose's labelling and carbon footprint and ignore the price, the latter being added afterwards using randomisation techniques. This choice is similar to the experimental design described in the Michaud, Llerena, and Joly (2012) work

and is easily explained when we take a closer look at the number of choice sets for different specifications. In simulated datasets it is traditional to use Full-Factorial (FF) experimental design as it uncovers completely the full potential of simulation tools: it allows to observe all the possible combinations of factors affecting some process and fully explore their implications. In our case, a simple full factorial design for a binary choice context has 28 combinations of factors (seven levels of prices, two levels for eco-label and two levels for Carbon imprint), but a complete full factorial design for a choice context with two alternatives implies 784 different combinations (as we have two alternatives each having 28 possible variants), which is unrealistic in a standard experimental study context and risks to be too demanding in terms of calculation times.

The dataset should be equilibrated with relatively identical number of choices for all three alternatives. In the field experiment the authors managed to achieve satisfying result with 67.5% of "Buy" choices and 32.5% for "Not to buy" choices, although the "A" and "B" alternatives showed different properties. The resulting observed descriptive statistics derived from the data proposed by Michaud, Llerena, and Joly (2012) are presented in table 5. The table focusses on the choice "Buy" descriptive statistics, ignoring the "No buy" option, for which all the attributes are considered to be equal to 0. The $p$-values are the results of the two subsets ("A" and "B") comparison[4].

Table 5: Alternatives' descriptive statistics by group, correlated random effects

|  | A (N=1186) | B (N=1186) | Total (N=2372) | p value |
|---|---|---|---|---|
| **Choice** |  |  |  | $< 0.001$ |
| Mean (SD) | 0.517 (0.500) | 0.159 (0.366) | 0.338 (0.473) |  |
| Range | 0.000 - 1.000 | 0.000 - 1.000 | 0.000 - 1.000 |  |
| **Price** |  |  |  | 0.418 |
| Mean (SD) | 2.990 (0.881) | 3.020 (0.893) | 3.005 (0.887) |  |
| Range | 1.500 - 4.500 | 1.500 - 4.500 | 1.500 - 4.500 |  |
| **Carbon** |  |  |  | $< 0.001$ |
| Mean (SD) | 0.167 (0.373) | 0.832 (0.374) | 0.500 (0.500) |  |
| Range | 0.000 - 1.000 | 0.000 - 1.000 | 0.000 - 1.000 |  |
| **Label** |  |  |  | 0.837 |
| Mean (SD) | 0.502 (0.500) | 0.497 (0.500) | 0.500 (0.500) |  |
| Range | 0.000 - 1.000 | 0.000 - 1.000 | 0.000 - 1.000 |  |

Of particular interest in the table 5 to us is the unbalanced structure of the resulting dataset. The $Carbon$ imprint of the different alternatives has not identical properties, which leads to different $Choice$ statistics, where the alternative with higher carbon imprint is chosen less frequently. In the original study such difference was not dangerous, because only the "Buy" option was compared against "No Buy" one. However, in case of the NN modelling such unbalanced dataset may lead to erroneous results,

---

[4]$\chi^2$ test is used for discrete variables, while *Anova* is implemented for continuous ones.

where the more popular alternative will always have a higher choice probability. The distribution inside the "Buy" group for different alternatives ("A" and "B") should be quasi-identical, producing equally distributed three groups of choices each nearing 33.3%. Even if this property is not as important for a traditional MNL model, we are interested to observe the same choice structure in our artificial dataset, because it may highly affect the performance of more advanced models, such as NN for example.

### 2.1.1 Generated dataset presentation

In this section we will discuss the resulting datasets simulated under the listed above assumptions.

For our dataset we choose to generate 160000 observations, for 1000 individuals, each facing 16 different choice sets 10 times. The 16 choice sets include all the possible combinations of two roses ("A" and "B") described by two environmental attributes, while prices are randomly assigned within the choice sets. The prices are assumed to be uniformly distributed over the choice sets, following a discrete uniform distribution. The prices vary among the different replications. This procedure resulted in sufficiently large dataset, which in the same time was not difficult to treat without implementation of Big Data specific techniques.

The original experimental design used to generate the choice sets assumed no branding for the alternatives to avoid any undesired bias in the results. Theoretically this design should have provided an equilibrated dataset with no correlation between different attributes, although the size of the final dataset might have affected the results. In our case we assume that individuals have no additional information about the roses in choice sets except the three observed attributes. As in the original work we assign insignificant labels "A" and "B" to the roses within choice sets, which is done mostly for convenience and has no impact on the individuals' decisions.

It is interesting to explore the statistical properties of the resulting datasets: the original one (Original), gathered by Michaud, Llerena, and Joly (2012) and made available in anonymised format by Iragaël Joly; and the two generated artificial datasets, assuming homogeneous (Generated FE) and heterogeneous (Generated RE) preferences respectively of the individuals for the environmental attributes. First of all, we may observe the individuals descriptive statistics for three datasets in the table 6.

Even though the $p$-values show no evident differences between the simulated datasets and the original one, except for the $Age$ variable, we observe the differences in the means. This is explained by the implemented dataset generation procedure. The variables in the original dataset are integers, assuming continuous nature of the real world variables. When synthesizing the dataset, we assume the quasi continuous variables, such as $Age$ and $Salary$ (denoted as $Income$ in original work) to be issued from normal distribution with parameters as figuring in the descriptive statistics for the original dataset, and only afterwards we convert the resulting values to integers. The binary variables $Sex$ and $Habit$ are generated with random draws from Bernoully distribution and consequently produce more realistic results. This procedure leads to potential biases in the resulting datasets, which is true not only for the individual variables, but for the alternatives' attributes as well.

Table 6: Individuals' characteristics descriptive statistics by dataset

|  | Fixed Effects (N=1000) | Random Effects (N=1000) | Target (N=102) | p value |
|---|---|---|---|---|
| **Sex** |  |  |  | 0.851 |
| Mean (SD) | 0.506 (0.500) | 0.515 (0.500) | 0.490 (0.502) |  |
| Range | 0.000 - 1.000 | 0.000 - 1.000 | 0.000 - 1.000 |  |
| **Habit** |  |  |  | 0.182 |
| N-Miss | 0 | 0 | 1 |  |
| Mean (SD) | 0.683 (0.466) | 0.657 (0.475) | 0.604 (0.492) |  |
| Range | 0.000 - 1.000 | 0.000 - 1.000 | 0.000 - 1.000 |  |
| **Salary** |  |  |  | < 0.001 |
| Mean (SD) | 2.750 (1.476) | 2.671 (1.438) | 2.147 (1.222) |  |
| Range | 1.000 - 6.000 | 1.000 - 6.000 | 1.000 - 6.000 |  |
| **Age** |  |  |  | 0.255 |
| Mean (SD) | 41.862 (13.685) | 42.161 (13.820) | 39.755 (18.895) |  |
| Range | 18.000 - 84.000 | 18.000 - 84.000 | 18.000 - 85.000 |  |

Table 7: Alternatives' descriptive statistics by dataset

|  | Fixed Effects (N=320000) | Random Effects (N=320000) | Target (N=2372) | p value |
|---|---|---|---|---|
| **Price** |  |  |  | 0.002 |
| Mean (SD) | 2.936 (0.958) | 2.936 (0.958) | 3.005 (0.887) |  |
| Range | 1.500 - 4.500 | 1.500 - 4.500 | 1.500 - 4.500 |  |
| **Carbon** |  |  |  | 0.999 |
| Mean (SD) | 0.500 (0.500) | 0.500 (0.500) | 0.500 (0.500) |  |
| Range | 0.000 - 1.000 | 0.000 - 1.000 | 0.000 - 1.000 |  |
| **Label** |  |  |  | 0.999 |
| Mean (SD) | 0.500 (0.500) | 0.500 (0.500) | 0.500 (0.500) |  |
| Range | 0.000 - 1.000 | 0.000 - 1.000 | 0.000 - 1.000 |  |

Secondly, we may as well observe the alternative specific descriptive statistics. They are presented in table 7. In this table we present the cumulative statistics for the "Buy" option, including both rose "A" and rose "B" properties, while 160000 entries (1186 entries for the original dataset) describing the "No buy" alternative are omitted, because their attributes are reduced to zeros in order to achieve identifiability of the models (a complete presentation of descriptive statistics par dataset and stratified by alternative may be found in Appendix C). The distributions of $Carbon$ footprint and Eco-$Label$ attributes follows perfectly the ones inside the original dataset, although the prices differ. This particular divergence, may be explained by the procedure implemented to assign prices to the alternatives inside choice sets, because the random generator algorithms different across statistical programs and potentially the procedures implemented in $R$ and $SAS$ are not identical.

What is more interesting, is the difference in the $Choice$ statistics. We may be interested in comparing the statistics for different classes in our sample to ensure that they are not biased in favour of label "A" or label "B", as in this case it risks to bias the estimates. For the artificial dataset the ratio of choices per "Buy" alternative is higher than 40% and reaches 47.3% for the fixed effect utility (table 8), while for the random effects specification the numbers are lower, reaching only 42% in mean for two classes (table 9). This particular observation is rather interesting as it demonstrates how the heterogeneous effects for alternatives' features the consumer decisions.

We will start with a close examination of the fixed effects dataset, where we can see, that prices are not equally distributed among the different choices.

Table 8: Alternatives' descriptive statistics by group, fixed coefficients

|  | A (N=160000) | B (N=160000) | Total (N=320000) | p value |
|---|---|---|---|---|
| **Choice** |  |  |  | $< 0.001$ |
| Mean (SD) | 0.427 (0.495) | 0.518 (0.500) | 0.473 (0.499) |  |
| Range | 0.000 - 1.000 | 0.000 - 1.000 | 0.000 - 1.000 |  |
| **Price** |  |  |  | $< 0.001$ |
| Mean (SD) | 3.069 (0.979) | 2.803 (0.917) | 2.936 (0.958) |  |
| Range | 1.500 - 4.500 | 1.500 - 4.500 | 1.500 - 4.500 |  |

The unbalanced prices potentially bias our dataset and we can see how the option with inferior mean prices is chosen less frequently. Even thought this differences do not affect the MNL and MMNL models, which calculate average effects for all the alternatives, there may be an impact over the performances of the NN models performances.

For the dataset with correlated random effects of the alternative specific variables, we observe an identical situation in table 9. The class with lower average prices is chosen more rarely by the consumers, while the overall choices are less frequent due to the presence of stochastic individual preferences for particular alternatives' attributes.

Table 9: Alternatives' descriptive statistics by group, correlated random effects

|  | A (N=160000) | B (N=160000) | Total (N=320000) | p value |
|---|---|---|---|---|
| **Choice** |  |  |  | < 0.001 |
| Mean (SD) | 0.382 (0.486) | 0.462 (0.499) | 0.422 (0.494) |  |
| Range | 0.000 - 1.000 | 0.000 - 1.000 | 0.000 - 1.000 |  |
| **Price** |  |  |  | < 0.001 |
| Mean (SD) | 3.069 (0.979) | 2.803 (0.917) | 2.936 (0.958) |  |
| Range | 1.500 - 4.500 | 1.500 - 4.500 | 1.500 - 4.500 |  |

We may conclude the preliminary datasets study and comparison with the main impression that two artificial datasets may be assumed to be quasi-identical. The slight differences in prices, captured by statistical tests may be considered insignificant in comparison with the biases present in the original dataset. What is more, even if the biases were more significant, the models' specification, which assumes no variable specific coefficients for choice A and B would have lead to the correct estimates, exactly as it was done by Michaud, Llerena, and Joly (2012). The heterogeneous preferences result in less probable decisions to buy a rose in the population, which should definitely impact the performances of our models. Now it rests to verify how well the number of selected models will be able to derive the target values for the relative utility function.

## 2.2 Modelling consumer choices under different assumptions

This part of the work aims at presenting the results of the estimation for our selection of the econometric and ML models. We should particularly underline the fact, that this section does not focus on the performances of the models as they will be discussed more in detail latter. There is still a double objective for this section, as before presentation of the obtained results, we should discuss the methods and techniques, which were implemented in order to estimate the models, presented earlier.

The estimation procedure and choice of the estimation algorithms as well as their numeric implementation in the statistical software are important in the context of model performance comparison. The different estimation procedures may lead to different results and different conclusions.

We consecutively estimate the chosen models over the two datasets: with and without the presence of heterogeneous preferences of the individuals for the environmental attributes. Then we compare the estimates with the target values we have used previously as inputs in defining the relative utility functions.

### 2.2.1 Estimation procedures

In this section we will discuss the different techniques implemented in order to estimate the different models, which were described in the first theoretical part of this work. The different algorithms may result in discrepancy in seminally identical mathematical models. This particular difference will be demonstrated in comparison of the MNL results and the estimates obtained through estimation of a CNN model imitating MNL model. What is more, different models can provide different insights into the real world state. For example, MMNL model should account for heterogeneity in consumer preferences in the presence of random alternative specific effects.

The econometric models focused on inference and understanding of the underlying effects are usually estimated over the full dataset as there is no question about the precision of the obtained results, but rather the statistical power achieved in idenfication of the effects. We will follow the same approach in order not to face the different question related to external validity and verification of the estimated model, as well as the questions related to verification and testing of the models' performances over some external dataset.

In this part of the work we will firstly present the different estimation techniques, starting with *maximum likelihood* (Cosslett 1981) estimator for the MNL, as well as it's algorithmic implementation within *R*, and the *Adam* algorithm (Kingma and Ba 2014) traditionally used to estimate the NN models. Afterwards, we will discuss the results of the estimations we obtain over the generated datasets, presented in the previous part.

#### 2.2.1.1 Maximum-Likelihood for MNL and MMNL

The MNL and MMNL models, both are estimated by the maximum likelihood method. In this technique the estimator is used to derive the parameters, which were the most likely to produce the observed results (observed dataset).

Assuming we face probabilities defined by some function $f(.)$ parametrized $\theta$, the joint probability density may be defined as:

$$\mathcal{L}(\theta) = \prod_{j}^{\Omega} P_i(j \mid \theta) \text{ with } \theta : max_\theta \mathcal{L}(\theta) \tag{65}$$

This function is also known as likelihood function. The log-likelihood is obtained through a $log$ transformation of the likelihood function:

$$L(\theta) = \sum_{j}^{\Omega} log(P_i(j \mid \theta)) \text{ with } \theta : min_\theta L(\theta) \tag{66}$$

As we can see the obtained function is then minimised by adjusting $\theta$ in order to obtain the optimal

parameters. The optimisation problem is non-linear and requires an implementation of some iterative technique to be solved. Under "general conditions" they are consistent, asymptotically efficient and asymptotically normally distributed (McFadden 2001).

Speaking about the algorithmic implementation within the statistical software, the optimization is performed by iteratively updating the vector of parameters by the amount given by *step $\times$ direction*. The *step* in this case is a positive scalar and the *direction* is given by:

$$D = H^{-1} \times g \tag{67}$$

Where $g$ represents the gradient, while $H^{-1}$ is an estimate of the inverse of the Hessian matrix.

In this procedure the main question is the choice and estimation procedure of $H^{-1}$, which has several possible definitions. For example, *Broyden–Fletcher–Goldfarb–Shanno (BFGS)* (Broyden 1970) algorithm may be implemented, which is an iterative method for solving unconstrained non-linear optimization problems. This algorithm updates $H^{-1}$ at each iteration using the variations of the vector of parameters and the gradient. The initial value of the matrix in this particular case is the inverse of the outer-product of the gradient. The initial step equals to 1 and, if the new value of the function is inferior to the previous value, it is divided by two, until a higher value is obtained. This iterative procedure stops when the gradient is sufficiently close to 0, which is achieved through comparison of the $g \times H^{-1} \times g$ product with the *tolerance* argument. An alternative stopping condition is achieved by introduction of the maximum number of iterations for the algorithm, which ensures the impossibility to fall into an eternal loop. We may summarise this algorithm as follows, as described in *mlogit* package documentation by Croissant (2020):

1. The likelihood for the baseline model is calculated (assuming all the parameters are 0);
2. The function is then evaluated, assuming a step equals to one;
3. If the value of likelihood function is lower than the baseline value, the step is divided by two until the likelihood increases;
4. The gradient $g$ is then computed;

The authors of *mlogit* package insist that this method is more efficient than other functions available in *R* at this time. The codes used to estimate MNL model are available in Appendix (Appendix D.1).

For the MMNL model there exists an interesting modification for the algorithm, because we need to estimate the random effects variances and covariances in case of correlated random effects. The parameters are not directly introduced inside the likelihood function, but rather the elements of the Choleski decomposition of the covariance matrix are used. The Choleski decomposition matrix $L$ is defined in this case as follows:

$$L = \begin{bmatrix} chol_{11} & 0 & 0 & 0 \\ chol_{12} & chol_{22} & 0 & 0 \\ chol_{13} & chol_{23} & chol_{33} & 0 \\ chol_{14} & chol_{24} & chol_{34} & chol_{44} \end{bmatrix} \tag{68}$$

Where indices correspond to the random effects variables, for example, in our case study we have four parameters: $Buy$ dummy variable, eco-$Label$, $Carbon$ footprint and the $LC$, which stands for the $Label$ and $Carbon$ cross-product. Once the estimates of the matrix elements are obtained a variance-covariance matrix can be obtained:

$$LL^T = \begin{bmatrix} \sigma_1^2 & \sigma_{21} & \sigma_{31} & \sigma_{41} \\ \sigma_{12} & \sigma_2^2 & \sigma_{32} & \sigma_{42} \\ \sigma_{13} & \sigma_{23} & \sigma_3^2 & \sigma_{43} \\ \sigma_{14} & \sigma_{24} & \sigma_{34} & \sigma_4^2 \end{bmatrix} = \Sigma \tag{69}$$

Where $\sigma_i^2$ stands for the variance of effect $i$ and $\sigma_{ij}$ represents the covariance between two random parameters $i$ and $j$. The codes used to estimate MMNL model may be found in Appendix (Appendix D.2).

### 2.2.1.2 Backpropagation algorithm for NN

For the estimation of the NN model we benefit from the flexibility offered by *Keras* (Allaire and Chollet 2020), which is a high-level NN API developed with a focus on the speed of computation, offering at the same time an astonishing level of control over the models. The port of *Keras* inside *R* offered by Allaire and Chollet (2020) allows us to correctly specify our model, devised to imitate the structure of the traditional MNL. This particular ML library offers a choice of different model estimation algorithms, ranging from the *state-of-the-art* to the most recent and advanced techniques. In this particular application it was decided to implement the *Adam* algorithm (Kingma and Ba 2014), which can be considered as rather outdated estimation method by the standards of ML field, because it was introduced only in 2014.

*Adam* is an algorithm for first-order gradient-based optimization of stochastic objective functions, based on adaptive estimates of lower-order moments. This method was proved to be computationally efficient, as well as to have low memory requirements. It is invariant to diagonal rescaling of the gradients, and is well suited for problems that are large in terms of data or parameters. This algorithm is also considered appropriate for non-stationary objectives, as well as the problems with very sparse gradients. What is more, one of the particular advantage for us is that the hyper-parameters do not typically require advanced tuning.

Historically, the *Adam* algorithm is an extension to the *stochastic gradient descent* or *SGD* (Kiefer, Wolfowitz, and others 1952) method. The latter is an iterative method for optimizing a differentiable or

sub-differentiable objective functions. It can be considered as a stochastic approximation of the *gradient descent* (GD) optimization, because it replaces the actual gradient, which is typically calculated from the entire data set, by an estimate, which is calculated from a randomly selected subset of the data. In high-dimensional optimization problems this technique reduces the computational complexity and hence the computation time, resulting in faster iterations. SGD has a single learning rate, denoted $\alpha$ by convention, for all weight updates during training. The learning rate is considered to fixed through an entire estimation procedure as well. The two latter features, are sometimes regarded as disadvantage of the particular estimation technique and may not be suitable in all the contexts.

Once we have briefly presented its original predecessor, we may pass directly to *Adam* algorithm description as well as the procedures, which influenced its creation. The chosen method combines the advantages of two other extensions of SGD, which are:

- *Adaptive Gradient Algorithm* (AdaGrad), where the per-parameter learning rate is maintained fixed, which is suitable for sparse data learning problems (Duchi, Hazan, and Singer 2011);
- *Root Mean Square Propagation* (RMSProp), for which the per-parameter learning rates are adapted based on the average of recent values of the gradients for the weight. Tthis is an unpublished method supported by the community, more information may be found in Bengio and CA (2015).

This properties make the algorithm especially well performing on a non-stationary problems, including noisy data, as well as any other problem types. Instead of adapting the parameter learning rates based on the mean values (first moments) as in RMSProp, *Adam* uses of the average of the second moments of the gradients as well.

The *Adam* is configured using a following set of hyper-parameters (for more details on numerical implementation and working with *Keras* see Appendix D.3):

- $alpha$, which stands for the learning rate or step size, designing the proportion at which the weights are updated. Traditionally, as it was proposed by authors Kingma and Ba (2014), the value of $\alpha$ equals to $1e-8$, but in our application we approach it to the values used in the DFGS algorithm, assuming that $\alpha = 1e-1 = 0.1$. Large values results in faster initial learning rate, before it is updated, while inferior values slow learning significantly and require more runs;
- $beta_1$, describes the exponential decay rate for the first moment estimates. We assume this value to be fixed to the defaults of *Keras*, which is $\beta_1 = 0.9$;
- $beta_2$ is the exponential decay rate for the second moment estimates, which is by default $\beta_2 = 0.999$;
- $\epsilon$ is the last hyper-parameter, which is a very small number to prevent any division by zero in the algorithm implementation. In *Keras* this value is $\epsilon = 1e-8$.

### 2.2.2 Estimation results presentation

The comparison of the estimates obtained by the different models over different datasets can be done in two steps. First of all, we are interested in the observed mean effects over the datasets, because the possibility to correctly identify the means for the coefficients is of utmost importance for the analysis, regardless of the assumption on the heterogeneity of these effects. Then we are going to explore the additional dimension, provided by the MMNL estimates, which comprises the estimates for the variance-covariance matrix of the correlated random effects. The estimates obtained directly are the entries of the Choleski decomposition matrix and need to be transformed in order to observe the variances and covariances.

The results for the means estimates are regrouped in the table 10 on page 57. Now we can pass to the discussion of the obtained results and demonstrate the differences of the performances observed for different algorithms. We are going to start with the discussion of the estimates obtained with more traditional to econometric field MNL and MMNL models. Effectively, the MNL model allows us to obtain the exact estimates, due to the fast convergence rate and the relative simplicity of the problem. What is more, and what is of particular interest for us, it is how the MMNL model performs on the MNL specific dataset with fixed effects. The estimates obtained with the MMNL model for the fixed effects dataset demonstrate quasi-identical estimates as traditional MNL model, which nearly all the Choleski decomposition matrix element estimates statistically insignificant to zeros. Observing the estimates obtained from the two models we may rightfully conclude, that there is no evident danger in implementing a MMNL model in place of a MNL model on the fixed effects dataset, because the obtained estimates will point out the absence of the heterogeneous preferences in such case. The only disadvantage of the models misspecification in this case resides in the significantly increased estimation time, which requires significantly more iteration in order to estimate correctly the variance-covariance matrix elements and, consequently, the estimation complexity.

On the contrary, in the case of presence of the correlated random effects in the preferences of the population the estimates are significantly biased for the MNL model. Moreover, the estimates obtained with the MMNL model are not identical to the input parameters, which were used during the simulation step. In this situation the MNL model tends to significantly underestimate the effects of all the characteristics and attributes for the choice situation. This can potentially lead to a notorious bias in case we were using incorrect model specification during a field experiment data exploration.

The results for the Choleski matrix entries estimates are regrouped into a single table 11 on page 58 Based on these estimates, we can comment as well the potential inefficiency of the implemented algorithm, even if it is one of the best available to us. Even though the estimates of the means obtained with MMNL in the presence of the random effects are close to the theoretical ones, the estimates of the variance-covariance matrix elements are rather close, but not perfectly calculated. Which is important, as we had a rather large dataset compared to the datasets typically collected during field studies: 1000 individuals with 10 replications of 16 choice sets situations for each totalling to 160000 choice situa-

Table 10: Estimation results: mean effects

| | Fixed effects | | | Random effects | | | Target |
|---|---|---|---|---|---|---|---|
| | MNL | MMNL | CNN | MNL | MMNL | CNN | |
| **Characteristics** | | | | | | | |
| Sex | 1.401*** | 1.400*** | 1.369 | 0.712*** | 1.297*** | 0.719 | 1.420 |
| | (0.031) | (0.031) | | (0.016) | (0.024) | | |
| Age | 0.009*** | 0.009*** | 0.010 | 0.007*** | 0.010*** | 0.005 | 0.009 |
| | (0.001) | (0.001) | | (0.001) | (0.001) | | |
| Salary | 0.048*** | 0.048*** | 0.060 | 0.066*** | 0.120*** | 0.062 | 0.057 |
| | (0.010) | (0.010) | | (0.005) | (0.008) | | |
| Habit | 1.070*** | 1.071*** | 1.056 | 0.361*** | 0.641*** | 0.343 | 1.027 |
| | (0.030) | (0.030) | | (0.016) | (0.024) | | |
| **Attributes** | | | | | | | |
| Price | −1.626*** | −1.628*** | −1.618 | −0.886*** | −1.586*** | −0.886 | −1.631 |
| | (0.010) | (0.010) | | (0.006) | (0.010) | | |
| Buy | 2.311*** | 2.313*** | 2.228 | 0.662*** | 2.180*** | 0.665 | 2.285 |
| | (0.065) | (0.066) | | (0.036) | (0.054) | | |
| Label | 2.815*** | 2.817*** | 2.810 | 1.279*** | 1.922*** | 1.277 | 2.824 |
| | (0.022) | (0.022) | | (0.015) | (0.023) | | |
| Carbon | 6.654*** | 6.662*** | 6.634 | 3.259*** | 5.430*** | 3.250 | 6.665 |
| | (0.032) | (0.033) | | (0.016) | (0.030) | | |
| LC | −2.781*** | −2.782*** | −2.765 | −1.546*** | −2.663*** | −1.558 | −2.785 |
| | (0.028) | (0.028) | | (0.019) | (0.030) | | |

*Note:* $^{*}p<0.1$; $^{**}p<0.05$; $^{***}p<0.01$

tions. While in the original field study 102 individuals with only 2 replications of 6 choice sets were present, mounting to 1224 observations.

This situation demonstrates the existing trade-off between the need to correctly specify the model from the start and the potential computation inconveniences in the case of implementation of a more complex model in case of uncertainty. In other words, the scientists always face the choice either to simply use more complex model, which requires more data, calculation time and resources, or to perform an extensive theoretical study beforehand in order to correctly specify and delimit the model from the start.

Table 11: Estimation results: standard deviations and covariances

| | Fixed effects | Random effects | Target |
|---|---|---|---|
| | MMNL | MMNL | |
| **Standard deviations** | | | |
| Buy | 0.095 | 2.960*** | 3.202 |
| | (0.061) | (0.028) | |
| Label | 0.031 | 2.687*** | 2.654 |
| | (0.077) | (0.023) | |
| Carbon | 0.164* | 3.734*** | 3.535 |
| | (0.076) | (0.026) | |
| LC | 0.145* | 2.851*** | 2.711 |
| | (0.071) | (0.031) | |
| **Covariances** | | | |
| Buy:Label | −0.948 | −0.311*** | −0.54 |
| | (5.116) | (0.026) | |
| Buy:Carbon | −0.886 | −0.565*** | −4.39 |
| | (1.954) | (0.026) | |
| Label:Carbon | 0.891 | 0.959*** | 8.77 |
| | (1.578) | (0.003) | |
| Buy:LC | 0.669 | 0.789*** | 6.17 |
| | (0.501) | (0.005) | |
| Label:LC | −0.576 | −0.490*** | −2.33 |
| | (4.423) | (0.032) | |
| Carbon:LC | −0.568 | −0.651*** | −4.82 |
| | (1.604) | (0.030) | |

*Note:* $^*p<0.1$; $^{**}p<0.05$; $^{***}p<0.01$

Now we can switch to the discussion of the estimates obtained with *Adam* estimated CNN model, identical in structure to the MNL model. For reminder, we use convolution layers to calculate the relative deterministic utilities for the population $V_j$ for three alternatives, which are then converted to probabilities using a softmax dense layer with predefined unit weights for corresponding neurons.

As for the CNN estimates, the table 10 demonstrates, that the obtained estimates are technically identical to the means, we could see in the previous part for the MNL model estimates. These results demon-

strate the flexibility of the NN models and the hypothetical possibility to implement them in place of traditional econometric models with only inconvenience being the relative complexity to obtain the variances for the weights estimates, as non known to us method allows this. This only inconvenience renders impossible to analyse the statistical significance for the obtained weight estimates, which can be seen only over the marginal effects graphs for particular variable on the probability, but this is other discussion's topic. For now, the most important part is that the CNN imitation of the MNL models, estimated with a high learning rate ($\alpha = 1e - 1$) *Adam* algorithm, allows to obtain correct estimates for the means of the theoretical utility function, assuming the variables were chosen correctly.

Because of the nature of the constructed CNN model latter performs similarly to the traditional MNL model. This situation implies that the proposed CNN algorithm is, identically to MNL model, unable to identify correct parameters and consequently derive the true means for the underlying coefficients of the relative utility function in the presence of heterogeneous preferences among individuals. Nevertheless, given the flexibility of the NN it is theoretically possible to device an algorithm imitating the MMNL model's behaviour or even propose some alternative modelling techniques which will be able to supply, not directly through estimated weights but rather after a supplementary study, the correct estimates for marginal effects of the attributes on the choice probabilities.

To summarise this section, we can underline the successful implementation of the chosen mathematical models over the artificially created datasets simulating different choice situations. The effects identified by all of the models are close to the target values, although there exists clear evidence that the MMNL models perform significantly better in mean effects identification in all the contexts. At the same time the MNL model and its synthetically recreated NN counterpart underestimate the coefficient of the given relative utility functions in presence of the correlated random parameters in the individual utilities.

## 2.3 Performance evaluation and comparison

This section comprises the results we managed to achieve in the exploration of different performance metrics and provides insights on the functioning of the discussed mathematical models in a given context. As we have seen in the previous part, where the effects' estimates were provided, all of the models are able to provide some estimates for the retaliate utility function parameters in different discrete choice set-ups. The most simple models performed well on the dataset defined by the homogeneous preferences in the population for environmental attributes, underestimating the effects in the presence of preference heterogeneity. In the same time the more complex MMNL model performed sufficiently well in both behavioural set-ups, although it demonstrated some potential problems with the algorithmic implementation.

### 2.3.1 Overall precision

First of all we focus our attention on the general performance metrics, describing how well the estimated models fit the predicted outcomes over an original dataset. As we have discussed earlier we use only some of the available measures in an attempt not to make this work too cumbersome. The retained performance metrics are: accuracy, describing the overall goodness of fit over observed choices of the subjects; and more complex KDL measure, which compares the distributions instead of more simple metrics, which use only the information available in the confusion matrix.

We can observe the values of these general performance measures, describing overall performance of a given classifier in the table 12. The table regroups the metrics' values for all the estimated models.

Table 12: General performance measures

|  | Fixed effects | | | Random effects | | |
|---|---|---|---|---|---|---|
|  | MNL | MMNL | CNN | MNL | MMNL | CNN |
| **Overall measures** | | | | | | |
| Accuracy | 0.863 | 0.863 | 0.723 | 0.725 | 0.863 | 0.721 |
| **Probabilistic measures** | | | | | | |
| KLD | 0.623 | 0.623 | 0.328 | 0.349 | 0.625 | 0.317 |

As we have underlined earlier we observe quite natural situation when the best model in terms of overall performance is the model, which was used in the data generation step. This situation perfectly demonstrates the potential bias, which is explained by our choice of the artificial data-generation algorithm. Nevertheless, it should be noted, that the MNL and MMNL models perform equally well on the fixed effects dataset, where the preferences for the environmental attributes are homogeneous. This fact supports our initial hypothesis that an implementation of a more complex model is preferred when the real effects are unknown to the researcher.

Focusing our attention on the CNN model observe that the *Adam* algorithm did not outperform the *BFGS* procedure. This observation may be explained by the data-generation set-up, where the generative algorithm favoured the MNL model, rather than *Adam*. The latter not supporting the fine tuning over the error distribution.

We can observe the results for the resources efficiency we managed to obtain, which are regrouped in the table 13. Even though we present all the time values, we are mostly interested with the "user" and "system" time values. The first one indicates the CPU time charged for the execution of user instructions of the calling process, while the second one stand for the CPU time spent for execution by the system on behalf of the calling process.

The more advanced *Adam* algorithm easily bypasses the algorithms available in the *mlogit* package, although this boost in efficiency goes at the cost of lower overall performance and goodness of fit. At

Table 13: Ressources efficiency

| | Fixed effects | | | Random effects | | |
|---|---|---|---|---|---|---|
| | MNL | MMNL | CNN | MNL | MMNL | CNN |
| User | 20.910 | 452.414 | 17.433 | 18.722 | 2066.934 | 16.806 |
| System | 0.153 | 1.712 | 0.714 | 0.004 | 16.112 | 0.415 |
| Total | 21.068 | 454.192 | 8.412 | 18.726 | 2083.221 | 7.604 |

the same time, the MMNL implementation is far less efficient and takes 128 times more time, than CNN model. This situation clearly illustrates us how the precision and flexibility come at higher costs.

### 2.3.2 Alternative specific metrics

We proceed with a look at some more specific measures. The table 14 regroups response specific metrics, that describe the precision of model in predicting only one target class of the dataset. These metrics are mostly used when we are interested in some in-depth insight into the model performance and allow to identify the models which perform the best over a single class of interest. Given the context of Michaud, Llerena, and Joly (2012) study we are interested in identifying the algorithm which predicts the best "buy" (A and B alternatives) against "not buy" (C) alternative, providing at the same time some information about the alternative chosen. In order to evaluate the performance at this dimension we use Geometric mean and the F-measure performance estimators.

Table 14: Variable specific performance measures, fixed effects data

| | Fixed effects | | | Random effects | | |
|---|---|---|---|---|---|---|
| | C | A | B | C | A | B |
| **Geometric mean** | | | | | | |
| MNL | 0.454 | 0.848 | 0.868 | 0.432 | 0.696 | 0.693 |
| MMNL | 0.454 | 0.849 | 0.867 | 0.452 | 0.848 | 0.867 |
| CNN | 0.443 | 0.697 | 0.698 | 0.447 | 0.697 | 0.700 |
| **F-measure** | | | | | | |
| MNL | 0.318 | 0.834 | 0.873 | 0.282 | 0.666 | 0.704 |
| MMNL | 0.318 | 0.834 | 0.873 | 0.316 | 0.833 | 0.873 |
| CNN | 0.291 | 0.665 | 0.706 | 0.294 | 0.665 | 0.707 |

In the table 14 we are interested with the entries in the columns corresponding to the "No buy" alternative (C). For the dataset with fixed effects across the population, the MNL and MMNL models perform identically according to both of the selected measures. The CNN model falls behind the econometrics models on the fixed effects dataset, although situation changes in the presence of heterogeneous effects.

61

In the more complex case scenario, when the individuals have varying across population preferences towards one or another attribute, the CNN model outperforms the simple MNL model in detecting "No buy" decisions for given choice sets, which is rather interesting, because the overall model's performance is still inferior to the MNL, as it was shown in table 14.

### 2.3.3 Willingness to pay and premiums

Here we should present the most important results comparing the estimates for the WTP, as well as the premiums for particular attributes derived for different models. The Premium to pay for a rose's particular attribute as it was described previously can be represented as:

$$Premium = \frac{\frac{\delta V}{\delta X_k}}{\frac{\delta V}{\delta Price}} \tag{70}$$

At the same time, the WTP for a rose may be seen as the ratio of two corresponding coefficients of dummy variable and price. The table 15 presents the estimated WTP and premiums for the models, which output fixed coefficient estimates, without taking into account the randomness of the individual effects. In other words, this table regroups the results, which do not require bootstrapping for confidence interval estimation.

Table 15: WTP and Premiums obtained with MNL and CNN

|        | Fixed effects | | Random effects | | Target |
|--------|-------|-------|-------|-------|--------|
|        | MNL   | CNN   | MNL   | CNN   |        |
| WTP    | 1.421 | 1.377 | 0.747 | 0.751 | 1.401  |
| Label  | 1.731 | 1.737 | 1.445 | 1.442 | 1.731  |
| Carbon | 4.091 | 4.101 | 3.679 | 3.669 | 4.086  |
| LC     | 4.112 | 4.129 | 3.378 | 3.352 | 4.110  |

For the estimation of the WTP and the premiums for more complex models (the MMNL in our case) we use the same procedure, as was implemented by Michaud, Llerena, and Joly (2012). Because the random parameters are assumed to be correlated in the MMNL model's specification, the estimated standard deviations and confidence intervals are obtained using the Krinsky and Robb parametric bootstrapping method (Krinsky and Robb 1986). This procedure consists of generating of multiple random draws from a multivariate normal distribution and using the obtained results to obtain the confidence interval estimates. Exactly as in the original study we generate 1000 draws from a multivariate normal distribution ($MNV(\mu, \Sigma)$), with the coefficient estimates as means $\mu$ and the estimated variance-covariance matrix of the random parameters as $\Sigma$.

The obtained results are then summarised as follows in the table 16

Table 16: WTP and Premiums obtained with MMNL

| | Statistics | | | | | |
|---|---|---|---|---|---|---|
| | Mean | St. Dev. | Min | Pctl(25) | Pctl(75) | Max |
| **Fixed effects** | | | | | | |
| WTP | 1.416 | 0.058 | 1.233 | 1.377 | 1.455 | 1.613 |
| Label | 1.732 | 0.019 | 1.672 | 1.720 | 1.745 | 1.791 |
| Carbon | 4.097 | 0.103 | 3.730 | 4.026 | 4.166 | 4.434 |
| LC | 4.116 | 0.098 | 3.741 | 4.051 | 4.182 | 4.421 |
| **Random effects** | | | | | | |
| WTP | 1.360 | 1.887 | $-4.239$ | 0.073 | 2.662 | 7.893 |
| Label | 1.243 | 1.667 | $-3.867$ | 0.104 | 2.330 | 6.638 |
| Carbon | 3.467 | 2.323 | $-4.026$ | 1.880 | 5.043 | 11.671 |
| LC | 3.036 | 3.240 | $-7.430$ | 0.908 | 5.160 | 14.259 |
| **Target** | | | | | | |
| WTP | 1.418 | 1.973 | $-4.474$ | 0.058 | 2.798 | 6.706 |
| Label | 1.735 | 1.611 | $-2.652$ | 0.653 | 2.849 | 6.709 |
| Carbon | 4.076 | 2.134 | $-1.774$ | 2.608 | 5.543 | 11.217 |
| LC | 4.106 | 3.379 | $-6.304$ | 1.913 | 6.439 | 14.612 |

*Note:*  The estimates are obtained with 1000 draws from MNV distribution

Comparing the estimates to the input values we observe that the variance of the WTP and Premiums estimates, estimated over a fixed effects dataset, do not potentially affect the conclusion one can derive from the results. The values stay positive with the 75% interval within 0.2€ of the mean estimate. Assuming the model is not re-estimated and adjusted after the insignificant estimators are obtained for Choleski matrix elements, the results remain valid.

We may conclude, that given sufficiently large dataset the implementation of more complex model is preferable, because it will allow to control for unknown parameters without adding a risk of obtaining biased results. The more simple models, should be preferred in a more restricted context. They allow to obtain the valid results only in the case of correct theoretical assumptions, biasing the estimates in other conditions. Consequently, in the presence of uncertainty about the presence of heterogeneity in the customer choice modelling questions there is a strong interest to implement a more complex model, readjusting it afterwards if needed.

# Conclusion

In this work we have introduced the reader to the problematic of the different modelling paradigms in application to the consumer choice studies. By means of an experimental theory-testing framework we demonstrate the complexity of the model performance evaluation problematic, showing the eventual bottlenecks and the questions to be answered on all the levels of data exploration procedure. The correct specification of the theoretical assumptions, the dataset generation, the model choice as well as the performance measure choice were studied. The main objective to propose a comprehensive methodology for theory-testing framework creation was accomplished, illustrating the devised frameworks' potential over an economic question issued from real world.

Two different consumer choice situation were explored, issued from the setting delimited by Michaud, Llerena, and Joly (2012). The discrete choice context allowed us to compare how the presence of heterogeneous preferences for environmental attributes affected the possibility to identify correctly the underlying utility functions, as well as to derive the WTP and premiums for the attributes. The implementation of artificial dataset simulation techniques proved its potential in creation of fully controlled data samples, providing two consistent datasets constructed under RUM assumptions. Given the data, we could observe, how taste heterogeneity affected the population's choice distribution and the resulting datasets, as well as their impact on models' performances.

A total of three models, issued from alien disciplines such as econometrics (MNL and MMNL) and ML (CNN-MNL), were implemented over the generated artificial datasets. We could demonstrate the differences and similarities between the traditional econometrics models and such ML techniques as NN. The econometric models allowed us to observe the potential biases that researchers risk to induce using the simplest models in unjustified context. The ML model made it possible to demonstrate, how different approaches to optimisation and algorithmic solutions influence the obtained results. Moreover, the framework demonstrated, that ML models could be used instead of the traditional econometrics techniques under correct specification, as technically NN are able to approximate any other more simple linear or non-linear model. All of the models demonstrated good overall performance given the homogeneous individual preferences setting, while only the most complex MMNL model achieved sufficient results in presence of taste heterogeneity.

The multidimensionality of the explored situation allows us to tear several solutions from this work in terms of model performances in presence of heterogeneous preferences. The MMNL models demonstrated a better adaptivity for the different datasets and consequently a better adaptiveness in all the cases. This family of models showed a great tolerance for the eventual misspecification in the assumptions of the presence of random effects. On the contrary, the MNL models produced biased estimates in the presence of the random effects in population, which indicates a great danger and signal the importance of the correct specifications a preliminary data studies to be performed before the models estimation. The only observed difference was in the way the resulting approximation was unable to directly estimate the variance for the linear part coefficients, which is not initially the main focus of the

NN models. However, the marginal effects could still be derived for the individual characteristics or the alternative specific attributes, assuming a correct approximation was used, which does not inflate the overall variance for the marginal effects.

Nevertheless, there exist potential biases that require particular attention and caution in future research. The implemented data-generation procedure risks to bias the results in favour of the econometrics models, which were used to simulate the data. Speaking about the models, we have observed that the adaptiveness and flexibility of the MMNL model comets at some costs in resources efficiency. The time, computation power and the data amount needed to achieve satisfying results are significantly higher than for the other models.

This work demonstrates only a fraction of the full potential of the theory-testing framework. Many extensions and generalisations should be performed before it could be used at scale. For example, it is particularly interesting to introduce an extension which will provide the possibility to explore and compare how different behavioural theories (RUM, RRM, QDM) affects the estimation results. Even more, with this methodology it becomes possible to explore the effects of non-additive utility presence or the behaviour of populations with mixed behaviours presence. Another extension concerns the implemented mathematical models and consists in incorporating the most recent developments in the ML field into the framework, enabling users to implement such models as decision trees or more advanced NN. Last, but not the least, the framework could be complemented with a methodological tool-set for hypothesis testing using the advantages of a controlled experiment data collection.

To summarise, we conclude that the experimental framework has proven its importance for the empirical and theoretical studies and has demonstrated its potential. There clearly exist a strong need for a more extensive study and development of this framework to provide the research community with a hypothesis testing tool-set, which could be used in the context of the consumer choice modelling. The exploration of potential biases and theory-testing will allow us to establish a comprehensive and consistent methodology to be implemented latter in empirical work and controlled experiments in particular.

# Table of contents

68

# List of figures

# List of tables

# Bibliography

Agrawal, Ajay, Joshua Gans, and Avi Goldfarb. 2019. *The Economics of Artificial Intelligence: An Agenda*. Book. National Bureau of Economic Research; University of Chicago Press. https://doi.org/https://doi.org/10.7208/chicago/9780226613475.001.0001.

Agresti, Alan. 2007. *An Introduction to Categorical Data Analysis, Second Edition*.

———. 2013. *Categorical Data Analysis, Third Edition*.

Allaire, JJ, and François Chollet. 2020. *Keras: R Interface to 'Keras'*. https://CRAN.R-project.org/package=keras.

Allaire, JJ, and Yuan Tang. 2020. *Tensorflow: R Interface to 'Tensorflow'*. https://CRAN.R-project.org/package=tensorflow.

Allaire, JJ, Yihui Xie, Jonathan McPherson, Javier Luraschi, Kevin Ushey, Aron Atkins, Hadley Wickham, Joe Cheng, Winston Chang, and Richard Iannone. 2018. *Rmarkdown: Dynamic Documents for R*. https://CRAN.R-project.org/package=rmarkdown.

Anderson, Simon P, Andre De Palma, and Jacques-Francois Thisse. 1992. *Discrete Choice Theory of Product Differentiation*. MIT press.

Athey, Susan. 2018. "The Impact of Machine Learning on Economics." Book. In *The Economics of Artificial Intelligence: An Agenda*, by Ajay Agrawal, Joshua Gans, and Avi Goldfarb, 507–47. National Bureau of Economic Research; University of Chicago Press. https://doi.org/https://doi.org/10.7208/chicago/9780226613475.001.0001.

Athey, Susan, and Guido W. Imbens. 2019. "Machine Learning Methods That Economists Should Know About." *Annual Review of Economics* 11 (1): 685–725. https://doi.org/10.1146/annurev-economics-080217-053433.

Ayodele, Taiwo Oladipupo. 2010. "Types of Machine Learning Algorithms." *New Advances in Machine Learning*. InTech, 19–48.

Baayen, Harald, Shravan Vasishth, Reinhold Kliegl, and Douglas Bates. 2017. "The Cave of Shadows: Addressing the Human Factor with Generalized Additive Mixed Models." *Journal of Memory and Language* 94: 206–34. https://doi.org/https://doi.org/10.1016/j.jml.2016.11.006.

Baldi, Pierre, Søren Brunak, Yves Chauvin, Claus A. F. Andersen, and Henrik Nielsen. 2000. "Assessing the accuracy of prediction algorithms for classification: an overview." *Bioinformatics* 16 (5): 412–24. https://doi.org/10.1093/bioinformatics/16.5.412.

Baltagi, Badi. 2008. *Econometric Analysis of Panel Data*. John Wiley & Sons.

Bengio, Yoshua, and MONTREAL CA. 2015. "Rmsprop and Equilibrated Adaptive Learning Rates for Nonconvex Optimization." *Corr Abs/1502.04390*.

Bernard, John C, and Daria J Bernard. 2009. "What Is It About Organic Milk? An Experimental Analysis." *American Journal of Agricultural Economics* 91 (3). Oxford University Press: 826–36.

Bhat, Chandra R. 1995. "A Heteroscedastic Extreme Value Model of Intercity Travel Mode Choice." Suggested. *Transportation Research Part B: Methodological* 29 (6): 471–83. https://EconPapers. repec.org/RePEc:eee:transb:v:29:y:1995:i:6:p:471-483.

Bouscasse, Hélène, Iragaël Joly, and Jean Peyhardi. 2019. "A new family of qualitative choice models: An application of reference models to travel mode choice." *Transportation Research Part B: Methodological* 121 (C): 74–91. https://doi.org/10.1016/j.trb.2018.12.010.

Brathwaite, Timothy, Akshay Vij, and Joan L Walker. 2017. "Machine Learning Meets Microeconomics: The Case of Decision Trees and Discrete Choice." *arXiv Preprint arXiv:1711.04826*.

Breiman, Leo, and others. 2001. "Statistical Modeling: The Two Cultures (with Comments and a Rejoinder by the Author)." *Statistical Science* 16 (3). Institute of Mathematical Statistics: 199–231.

Brock, William, and Steven Durlauf. 2003. "Multinomial Choice with Social Interactions." NBER Technical Working Papers 0288. National Bureau of Economic Research, Inc. https://EconPapers. repec.org/RePEc:nbr:nberte:0288.

Broyden, Charles G. 1970. "The Convergence of Single-Rank Quasi-Newton Methods." *Mathematics of Computation* 24 (110): 365–82.

Cascetta, Ennio. 2009. *Transportation Systems Analysis: Models and Applications*. Vol. 29. Springer Science & Business Media.

Chen, Bryant, and Judea Pearl. 2013. "Regression and Causation: A Critical Examination of Six Econometrics Textbooks." *Real-World Economics Review, Issue*, no. 65: 2–20.

Chorus, Caspar G. 2010. "A New Model of Random Regret Minimization." *European Journal of Transport and Infrastructure Research* 10 (2).

Cosslett, Stephen R. 1981. "Maximum Likelihood Estimator for Choice-Based Samples." *Econometrica: Journal of the Econometric Society*. JSTOR, 1289–1316.

Coussement, Kristof, Dries F. Benoit, and Dirk Van den Poel. 2010. "Improved Marketing Decision Making in a Customer Churn Prediction Context Using Generalized Additive Models." *Expert Systems with Applications* 37 (3): 2132–43. https://doi.org/https://doi.org/10.1016/j.eswa.2009.07.029.

Croissant, Yves. 2020. *Mlogit: Multinomial Logit Models*. https://CRAN.R-project.org/package= mlogit.

Denuit, Michel, and Donatien Hainaut. 2019. *Effective Statistical Learning Methods for Actuaries Iii: Neural Networks and Extentions*. Springer.

Denuit, Michel, and Julien Trufin. 2019. *Effective Statistical Learning Methods for Actuaries I: GLMs and Extentions*. Springer.

De Palma, André, Robin Lindsey, Emile Quinet, and Roger Vickerman. 2011. *A Handbook of Transport Economics*. Edward Elgar Publishing.

Donoho, David. 2017. "50 Years of Data Science." *Journal of Computational and Graphical Statistics* 26 (4). Taylor & Francis: 745–66.

Duchi, John, Elad Hazan, and Yoram Singer. 2011. "Adaptive Subgradient Methods for Online Learning and Stochastic Optimization." *Journal of Machine Learning Research* 12 (7).

Fiebig, Denzil, Michael Keane, Jordan Louviere, and Nada Wasi. 2010. "The Generalized Multinomial Logit Model: Accounting for Scale and Coefficient Heterogeneity." *Marketing Science* 29 (3): 393–421. https://EconPapers.repec.org/RePEc:inm:ormksc:v:29:y:2010:i:3:p:393-421.

Fürnkranz, J., and E. Hüllermeier. 2010. *Preference Learning*. Springer Verlag, Berlin.

Garrow, Tudor D.; Lee, Laurie A.; Bodea. 2010. "Generation of Synthetic Datasets for Discrete Choice Analysis." *Transportation* 37 (2): 183–202. https://doi.org/10.1007/s11116-009-9228-6.

Greene, William H. 2008. "The Econometric Approach to Efficiency Analysis." *The Measurement of Productive Efficiency and Productivity Growth* 1 (1): 92–250.

Harrison, Glenn W, Ronald M Harstad, and E Elisabet Rutström. 2004. "Experimental Methods and Elicitation of Values." *Experimental Economics* 7 (2). Springer: 123–40.

Hastie, Trevor, Robert Tibshirani, and Jerome Friedman. 2009. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer Science & Business Media.

Heinzen, Ethan, Jason Sinnwell, Elizabeth Atkinson, Tina Gunderson, and Gregory Dougherty. 2020. *Arsenal: An Arsenal of 'R' Functions for Large-Scale Statistical Summaries*. https://CRAN.R-project.org/package=arsenal.

Hlavac, Marek. 2018. *Stargazer: Well-Formatted Regression and Summary Statistics Tables*. https://CRAN.R-project.org/package=stargazer.

Japkowicz, Nathalie, and Mohak Shah. 2011. *Evaluating Learning Algorithms: A Classification Perspective*. Cambridge University Press. https://doi.org/10.1017/CBO9780511921803.

Jebara, Tony. 2004. *Machine Learning: Discriminative and Generative*. Springer Science & Business Media.

Kiefer, Jack, Jacob Wolfowitz, and others. 1952. "Stochastic Estimation of the Maximum of a Regression Function." *The Annals of Mathematical Statistics* 23 (3). Institute of Mathematical Statistics: 462–66.

Kingma, Diederik P, and Jimmy Ba. 2014. "Adam: A Method for Stochastic Optimization." *arXiv Preprint arXiv:1412.6980*.

Kirk, Roger E. 2012. "Experimental Design." *Handbook of Psychology, Second Edition* 2. Wiley Online Library.

Kotsiantis, Sotiris, I. Zaharakis, and P. Pintelas. 2006. "Machine Learning: A Review of Classification and Combining Techniques." *Artificial Intelligence Review* 26 (November): 159–90. https://doi.org/10.1007/s10462-007-9052-3.

Krinsky, Itzhak, and A Leslie Robb. 1986. "On Approximating the Statistical Properties of Elasticities." *The Review of Economics and Statistics*. JSTOR, 715–19.

Kuhfeld, Warren F. 2003. *Marketing Research Methods in Sas.* Citeseer.

Kukar, Matjaž, and Igor Kononenko. 2002. "Reliable Classifications with Machine Learning." In *European Conference on Machine Learning*, 219–31. Springer.

Leong, Waiyan, and David A. Hensher. 2015. "Contrasts of Relative Advantage Maximisation with Random Utility Maximisation and Regret Minimisation." *Journal of Transport Economics and Policy (JTEP)* 49 (1): 167–86. https://www.ingentaconnect.com/content/lse/jtep/2015/00000049/00000001/art00010.

Louviere, Jordan J, David A Hensher, and Joffre D Swait. 2000. *Stated Choice Methods: Analysis and Applications*. Cambridge university press.

McCausland, William J., and A.A.J. Marley. 2013. "Prior Distributions for Random Choice Structures." *Journal of Mathematical Psychology* 57 (3): 78–93. https://doi.org/https://doi.org/10.1016/j.jmp.2013.05.001.

McFadden, Daniel. 1974. "The Measurement of Urban Travel Demand." *Journal of Public Economics* 3 (4): 303–28. https://doi.org/https://doi.org/10.1016/0047-2727(74)90003-6.

———. 2001. "Economic Choices." *The American Economic Review* 91 (3). American Economic Association: 351–78. http://www.jstor.org/stable/2677869.

McFadden, Daniel, and Kenneth Train. 2000. "Mixed Mnl Models for Discrete Response." *Journal of Applied Econometrics* 15 (5). Wiley Online Library: 447–70.

Michaud, Celine, Daniel Llerena, and Iragael Joly. 2012. "Willingness to pay for environmental attributes of non-food agricultural products: a real choice experiment." *European Review of Agricultural Economics* 40 (2): 313–29. https://doi.org/10.1093/erae/jbs025.

———. 2013. "Willingness to Pay for Environmental Attributes of Non-Food Agricultural Products: A Real Choice Experiment." *European Review of Agricultural Economics* 40 (2). Oxford University Press: 313–29.

Microsoft, and Steve Weston. 2020. *Foreach: Provides Foreach Looping Construct*. https://CRAN.R-project.org/package=foreach.

Molina, Mario, and Filiz Garip. 2019. "Machine Learning for Sociology." *Annual Review of Sociology* 45 (1): 27–45. https://doi.org/10.1146/annurev-soc-073117-041106.

Mullainathan, Sendhil, and Jann Spiess. 2017. "Machine Learning: An Applied Econometric Approach." *Journal of Economic Perspectives* 31 (2): 87–106. https://doi.org/10.1257/jep.31.2.87.

Munizaga, Marcela A., and Ricardo Alvarez-Daziano. 2005. "Testing Mixed Logit and Probit Models by Simulation." *Transportation Research Record* 1921 (1): 53–62. https://doi.org/10.1177/0361198105192100107.

Paredes, Miguel, Erik Hemberg, Una-May O'Reilly, and Chris Zegras. 2017. "Machine Learning or Discrete Choice Models for Car Ownership Demand Estimation and Prediction?" In *2017 5th Ieee International Conference on Models and Technologies for Intelligent Transportation Systems (Mt-Its)*, 780–85. IEEE.

Pigozzi, Gabriella, Alexis Tsoukiàs, and Paolo Viappiani. 2016. "Preferences in Artificial Intelligence." *Annals of Mathematics and Artificial Intelligence* 77 (3-4). Springer Verlag: 361–401. https://doi.org/10.1007/s10472-015-9475-5.

R Core Team. 2018a. *Foreign: Read Data Stored by 'Minitab', 'S', 'Sas', 'Spss', 'Stata', 'Systat', 'Weka', 'dBase', ...* https://CRAN.R-project.org/package=foreign.

———. 2018b. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. https://www.R-project.org/.

Revelt, David, and Kenneth Train. 1998. "Mixed Logit with Repeated Choices: Households' Choices of Appliance Efficiency Level." *Review of Economics and Statistics* 80 (4). MIT Press: 647–57.

Rose, John M, and Michiel CJ Bliemer. 2006. "Constructing Efficient Stated Choice Experimental Designs." *Transport Reviews* 29 (5). Taylor & Francis: 587–617.

Rose, John M, Michiel CJ Bliemer, David A Hensher, and Andrew T Collins. 2008. "Designing Efficient Stated Choice Experiments in the Presence of Reference Alternatives." *Transportation Research Part B: Methodological* 42 (4). Elsevier: 395–406.

Talebijamalabad, Amirreza, Nikita Gusarov, and Iragael Joly. 2020. *Sdcm: Simulation of Discrete Choice Models*.

Train, Kenneth E. 2009. *Discrete Choice Methods with Simulation*. Cambridge university press.

Tsoukiàs, Alexis, and Paolo Viappiani. 2013. "Tutorial on preference handling." In *ACM Conference on Recommender System (RecSys)*, 497–98. Hong Kong, China. https://doi.org/10.1145/2507157.2508065.

Tsoumakas, Grigorios, and Ioannis Katakis. 2007. "Multi-Label Classification: An Overview." *International Journal of Data Warehousing and Mining (IJDWM)* 3 (3): 1–13. https://EconPapers.repec.org/RePEc:igg:jdwm00:v:3:y:2007:i:3:p:1-13.

Varian, Hal R. 2014. "Big Data: New Tricks for Econometrics." *Journal of Economic Perspectives* 28 (2): 3–28. https://doi.org/10.1257/jep.28.2.3.

Vitetta, Antonino. 2016. "A Quantum Utility Model for Route Choice in Transport Systems." *Travel Behaviour and Society* 3. Elsevier: 29–37.

Wickham, Hadley. 2019. *Tidyverse: Easily Install and Load the 'Tidyverse'*. https://CRAN.R-project.org/package=tidyverse.

Xie, Yihui. 2020a. *Knitr: A General-Purpose Package for Dynamic Report Generation in R*. https://CRAN.R-project.org/package=knitr.

———. 2020b. *Tinytex: Helper Functions to Install and Maintain Tex Live, and Compile Latex Documents*. https://CRAN.R-project.org/package=tinytex.

Yukalov, Vyacheslav I, and Didier Sornette. 2017. "Quantum Probabilities as Behavioral Probabilities." *Entropy* 19 (3). Multidisciplinary Digital Publishing Institute: 112.

Zielesny, Achim. 2011. *From Curve Fitting to Machine Learning*. Vol. 18. Springer.

# Appendices

## A Taxonomies of statistical models

Figure 7: Taxonomy as proposed by Hastie and Tibshirani (2009), reduced form



Figure 8: Taxonomy as proposed by Ayodele (2010)

Figure 9: Taxonomy as proposed by Agresti (2013), based on data types

# B Performance measures positioning

Figure 10: Performance measures as described by Japkowicz (2011)

All measures

Confusion matrix based

Additional information

Alternative information

Deterministic classifiers

Scoring classifiers

Continuous Probabilistic classifiers

Multicriteria measures

Multiclass focused

Single-class focused

Graphical measures

Summary statistics

Distance error measures

Information Theoretic measures

Accuracy
Error rate

TP-TN ratios
F-measure
Geometric means
...

ROC
PR
DET
Lift
and Cost curves

AUC

RMSE

KLD
KB IR
BIR

# C Descriptive statistics

## C.1 Comparing datasets over A alternative

Table 17: Alternatives' descriptive statistics by dataset, stratified by alternative

| Alternative | | Fixed Effects (N=320000) | Random Effects (N=320000) | Target (N=2372) | p value |
|---|---|---|---|---|---|
| A | **Alternative** | | | | |
| | A | 160000 (100.0%) | 160000 (100.0%) | 1186 (100.0%) | |
| | B | 0 (0.0%) | 0 (0.0%) | 0 (0.0%) | |
| | **Choice** | | | | < 0.001 |
| | Mean (SD) | 0.427 (0.495) | 0.382 (0.486) | 0.517 (0.500) | |
| | Range | 0.000 - 1.000 | 0.000 - 1.000 | 0.000 - 1.000 | |
| | **Price** | | | | 0.022 |
| | Mean (SD) | 3.069 (0.979) | 3.069 (0.979) | 2.990 (0.881) | |
| | Range | 1.500 - 4.500 | 1.500 - 4.500 | 1.500 - 4.500 | |
| | **Carbon** | | | | < 0.001 |
| | Mean (SD) | 0.500 (0.500) | 0.500 (0.500) | 0.167 (0.373) | |
| | Range | 0.000 - 1.000 | 0.000 - 1.000 | 0.000 - 1.000 | |
| | **Label** | | | | 0.993 |
| | Mean (SD) | 0.500 (0.500) | 0.500 (0.500) | 0.502 (0.500) | |
| | Range | 0.000 - 1.000 | 0.000 - 1.000 | 0.000 - 1.000 | |
| | **Price by group** | | | | < 0.001 |
| | 1.5 | 16000 (10.0%) | 16000 (10.0%) | 82 (6.9%) | |
| | 2 | 24000 (15.0%) | 24000 (15.0%) | 223 (18.8%) | |
| | 2.5 | 27000 (16.9%) | 27000 (16.9%) | 214 (18.0%) | |
| | 3 | 23000 (14.4%) | 23000 (14.4%) | 175 (14.8%) | |
| | 3.5 | 22000 (13.8%) | 22000 (13.8%) | 187 (15.8%) | |
| | 4 | 21000 (13.1%) | 21000 (13.1%) | 219 (18.5%) | |
| | 4.5 | 27000 (16.9%) | 27000 (16.9%) | 86 (7.3%) | |

## C.2 Comparing datasets over B alternative

Table 18: Alternatives' descriptive statistics by dataset, stratified by alternative

| Alternative | | Fixed Effects (N=320000) | Random Effects (N=320000) | Target (N=2372) | p value |
|---|---|---|---|---|---|
| B | **Alternative** | | | | |
| | A | 0 (0.0%) | 0 (0.0%) | 0 (0.0%) | |
| | B | 160000 (100.0%) | 160000 (100.0%) | 1186 (100.0%) | |
| | **Choice** | | | | < 0.001 |
| | Mean (SD) | 0.518 (0.500) | 0.462 (0.499) | 0.159 (0.366) | |
| | Range | 0.000 - 1.000 | 0.000 - 1.000 | 0.000 - 1.000 | |
| | **Price** | | | | < 0.001 |
| | Mean (SD) | 2.803 (0.917) | 2.803 (0.917) | 3.020 (0.893) | |
| | Range | 1.500 - 4.500 | 1.500 - 4.500 | 1.500 - 4.500 | |
| | **Carbon** | | | | < 0.001 |
| | Mean (SD) | 0.500 (0.500) | 0.500 (0.500) | 0.832 (0.374) | |
| | Range | 0.000 - 1.000 | 0.000 - 1.000 | 0.000 - 1.000 | |
| | **Label** | | | | 0.985 |
| | Mean (SD) | 0.500 (0.500) | 0.500 (0.500) | 0.497 (0.500) | |
| | Range | 0.000 - 1.000 | 0.000 - 1.000 | 0.000 - 1.000 | |
| | **Price by group** | | | | < 0.001 |
| | 1.5 | 25000 (15.6%) | 25000 (15.6%) | 108 (9.1%) | |
| | 2 | 28000 (17.5%) | 28000 (17.5%) | 192 (16.2%) | |
| | 2.5 | 26000 (16.2%) | 26000 (16.2%) | 158 (13.3%) | |
| | 3 | 30000 (18.8%) | 30000 (18.8%) | 204 (17.2%) | |
| | 3.5 | 18000 (11.2%) | 18000 (11.2%) | 232 (19.6%) | |
| | 4 | 23000 (14.4%) | 23000 (14.4%) | 195 (16.4%) | |
| | 4.5 | 10000 (6.2%) | 10000 (6.2%) | 97 (8.2%) | |

# D *R* code for implemented models

## D.1 MNL model

```r
# Transform dataset to mlogit format
mnl_data = data %>%
    mlogit.data(
        choice = "Choice",
        alt.var = "Alternative",
        shape = "long", # Long format
        alt.levels = c("C", "A", "B") # Define order of alternatives
    )


# Function
utility = Choice ~ Sex + Age + Salary + Habit + # Individual characteristics
    Price + Buy + Label + Carbon + LC + 0 | 0 # Alternatives attributes


# Estimate MNL model
mnl_novar = mlogit(
        utility,
        data = mnl_data,
        reflevel = "C", # The No-buy option is the baseline
        print.level = 3, # Print estimation details
        iterlim = 1000
    )
```

## D.2 MMNL model

```r
# Transform dataset to mlogit format
mmnl_data = data %>%
    mlogit.data(
        choice = "Choice",
        alt.var = "Alternative",
        id = "ID", # Set individuals' index
        chid = "CHID", # Set choice sets index
        shape = "long",
        alt.levels = c("C", "A", "B")
```

```
    )

# Function
utility = Choice ~ Sex + Age + Salary + Habit + # Individual characteristics
    Price + Buy + Label + Carbon + LC + 0 | 0 # Alternatives attributes

# Estimate MMNL model
mmnl = mlogit(
        utility,
        data = mmnl_data,
        reflevel = "C", # The No-buy option is the baseline
        correlation = TRUE, # Include covariance (and not variance only)
        rpar =  c( # Normality assumption and four parameters
            "Buy" = "n",
            "Label" = "n",
            "Carbon" = "n",
            "LC" = "n"
        ),
        panel = TRUE, # Estimate dataset as panel
        print.level = 3, # Print estimation details
        iterlim = 1000
    )
```

### D.3 CNN model with *Adam* algorithm

```
# Used libraries
library(tidyverse)
library(tensorflow)
library(keras)

# Define optimization algorithm to be used
adam_own = optimizer_adam(
    lr = 1e-1, # We adjust the learning rate, keeping the rest as defaults
    beta_1 = 0.9,
    beta_2 = 0.999,
    epsilon = NULL,
    decay = 0,
```

```
    amsgrad = FALSE,
    clipnorm = 6, # We limit as well the max value for weights
    clipvalue = NULL
)


# Set hyperparameters
## The number of epochs is a hyperparameter that defines the number times
## that the learning algorithm will work through the entire training
## dataset.
epoch = 50
## The batch size is a hyperparameter that defines the number of samples
## to work through before updating the internal model parameters.
batch = 16000


# Limit softmax weights
## (keras uses dense layer transformation inside softmax layer by default)
softmax_weights = list(
    matrix(
        c(  1,  0,  0,
            0,  1,  0,
            0,  0,  1),
        nrow = 3
    )
)


# Setup CNN model
model_cnn = keras_model_sequential() %>%
    # We reshape the dataset, as 1D convolution requires 3D tensor as input
    layer_reshape(
        target_shape = c(27, 1),
        input_shape = 27,
        trainable = FALSE
    ) %>%
    # 1D convolution layer
    layer_conv_1d(
        filters = 1L, # Dimentions of the output space
        kernel_size = 9L, # Number of parameters
        strides = 9L, # Strides of convolution equal to parameters side
```

```r
    # The starting value is 0 to ensure reproducibility
    kernel_initializer = "zeros",
    # The constant is not added, because we already have "Buy" dummy
    use_bias = FALSE,
    # We want a linear activation function
    activation = "linear",
    input_shape = c(27, 1)
) %>%
# An inverse transformation into a 2D tensor for softmax implementation
layer_flatten(
    data_format = "channels_first"
) %>%
# Softmax layer
layer_dense(
    units = 3, # Number of units equal to categories (3 utilities)
    use_bias = FALSE, # The bias constant is not estimated
    weights = softmax_weights,
    trainable = FALSE, # This layer is fixed
    activation = "softmax" # Softmax layer (to obtain probabilities)
) %>%
# Learning algorith definition
compile(
    loss = "categorical_crossentropy", # Choice of loss function
    optimizer = adam_own, # Parametrised Adam
    metrics = c("accuracy") # Target metrics
) %>%
# Training the model
fit(
    X_train, Y_train, # To train the model we use 80% of our dataset
    epochs = epoch,
    batch_size = batch,
    validation_data = list(X_test, Y_test) # 20% for validation
)
```

# Annexes

## I Simulation tool for performance comparison of discrete choice models

**Author:** Amirreza Talebijamalabad, M1 SIE (Grenoble INP)

**Under supervision of:** Iragaël Joly, HDR (GAEL, UGA, Grenoble INP)

**Available at:** https://github.com/Amirreza-96/sdcm

An experimental design is a plan which identifies the independent, dependent, and nuisance variables and indicates the way in which the randomization and statistical aspects of an experiment are to be carried out. Speaking of experimental design, we need to bear randomization, replication and blocking in our minds as three key elements of the experimental design (Kirk 2012). Randomization as a rather new concept in design of experiments, plays a pivotal role in distribution of idiosyncratic characteristics and variables' levels so that they do not selectively bias the outcome of the experiment. For example, in our designs, we applied randomization to avoid dominant alternatives as much as we can. Replication is the observation of two or more experimental units under the same conditions. Replication enables us to validate the proposed model and ensures the precise effects. Usually, in simulation, we run a very long replication or we make relatively many replications but small in dimensions, which we choose to replicate once but large enough. Blocking, on the other hand, is an experimental procedure for isolating variation attributable to a nuisance variable. Also, making blocks, we can randomly assign respondents to the choice sets or control the number of respondents in order to intimate the real scenarios; However, blocking is not of great importance when we are talking in the realm of simulation since we can control variations and variables.

Stated choice experiments present sampled respondents with a number of different choice situations, each consisting of a universal but finite set of alternatives defined on a number of attribute dimensions. Respondents are then asked to specify their preferred alternatives given a specific hypothetical choice context. In simulation, since the respondents are artificial, it will not be wiseful to sample the population, instead, replicationg the processes would be fruitful as the population will be generated repeatedly which is more close to reality. Moreover, to simulate the choice making process, based on decision rules such as utility maximization, utilities are calculated to reveal the choices of the individuals. SC data requires that the analyst designs the experiment in advance by assigning attribute levels to the attributes that define each of the alternatives which respondents are asked to consider(Rose et al. 2008).

To generate experimental designs for SC studies we need to find out how to allocate the attribute levels to the design matrix. Traditionally, researchers have relied on the principle of orthogonality to populate the choice situations shown to respondents. The orthogonality of an experimental design relates to the correlation structure between the attributes of the design. however, this class of designs may not be statistically efficient, as they do not take the SC model specification into account. These models are optimal for the linear models and assure the researcher that multicollinearity does not exist in design.

Considering this, it is assumed that such designs can be used for the non-linear models by linear arrangements(Kuhfeld 2003). It is important to note however, that the orthogonality of a design suggests nothing about whether two or more attributes are cognitively correlated in the minds of the respondents (e.g. price and quality attributes). As such, orthogonality is purely a statistical property of the design and not a behavioural property imposed upon the experiment(Rose and Bliemer 2006). Moreover, by entreing non-design attributes such as socio-demographic variables, any covariate within the dataset will unlikely be orthogonal, not only amongst themselves, but also with the design attributes. For example, if age, gender and income are added as variables in an analysis, correlations are not only likely to exist for these variables, but given that the variables described are constant over all choice situations within individual respondents, correlations between these variables and other attributes of the design are also likely to exist. Simulation tool should allow us to enter or not such soci-demographic variables to the simulation process so that at least we have some control on correlations. Furthermore, more advanced data generation methods should be applied to generate correlated data with specific precision. In this research, we have made a very conventional and widespread design so called full factorial design. It contains all of the possible levels of factors, and allows us to estimate all of the main effects and two-way interactions. Main effects are independent of the levels of other attributes, however; interactions involve two or more factors in which, effect of one factor depends on the level of another(Kuhfeld 2003). Furthermore, there are fractional orthogonal designs known as efficient designs providing ratherly small but efficient designs, also, there are algorithms to determine the correlations between columns as orthogonality is violated in these designs. It would be excellent if various kinds of designs were available in simulator.

Michaud, Llerena, and Joly (2013) conducted an empirical work to figure out consumers' willingness to pay, and a price premium for two environmental attributes of a non-food agricultural product(Roses). In this research there are two unlabelled alternatives Rose A and B and one no choice alternative. The two attributes, Label and Carbon, have two levels which make four combinations, hence six pairs of alternatives can be drawn from these combinations. Price ranges from 1.5 to 4.5 and is randomly assigned to the combinations of two other attributes. Finally, each respondent is faced with twelve choice sets(24 combinations of all attributes or 12 questions), hence, considerring no choice mode, there are three alternatives in each question. Trying to simulate the paper's results, we made a design with the same attributes and attribute levels. We sample put alternatives two by two in choice sets (16 choice sittuation), hence, respondents are faced twice with six pairs of alternatives, but the price is randomly assigned to each of the choice sets. Moreover, as no-choice mode does not effect the design, we do not add this mode to the design but finally, when it comes to utility comparison and decision process, this alternative is taken into account. Furthermore, we have not put interaction variable in the design since as no-choice mode, it does not affect the combinations of the design. To add, it makes the tool more flexible if we allow the user to decide about these two options.

Michaud, Llerena, and Joly (2013) considered four socioeconomic characteristics as well as sex, age, income and organic purchase habit. Since no information were available in regard to these features'

correlation, we assumed that they are uncorrelated, and made each feature independently. This is a limitation for data generation process. Simulator must enable the user to specify whether data is correlated or not. Moreover, it should allow the user to enter the inputs and specifications as well as distributions and their parameters. In order to generate sex data, we draw samples out of a uniform distribution with parameters $a = 0, b = 1$, then we assume that there is a $0.49$ chance that a respondent is female. Hence, if the random number is in range $(0, 0.49)$, hypothetical individual is female, otherwise, is male. The same procedure applies to the habit feautre. If the random number is among $(0, 0.35)$, organic habit is assumed to be zero. In order to generate age feature, the best distribution that we can draw samples which exactly could resemble the real data is truncated normal distribution. However, we just have the tnormal distribution's parameters and we need to have the underlying normal distribution's parameters(mean, and std.) to be able to draw samples. To tackle this, we solve a system of non-linear equations utilising numerical methods(Newton Raphson method) to find the underlying normal dist. parameters. Another way to generate such data is to draw samples from a log normal distribution. We still have a problem with this way since we need to have data ranging from 18 to 85, nevertheless positive values are generated. And finally, we simply take draws from normal distribution with the same parameters. As future improvements, simulation tool should be able to generate data based on theoretical distributions or empirical ones. Hence, some curve fitting procedures to find the distribution best fitting the real datashould be installed in the tool.

So far, we have made the design and socioeconomic features for artificial individuals. Now, we need to specify utilities per each individual, and finally, due to the RUM model, we select the alternative with highest utility per each individual per each choice set. As a future improvement, we suggust that simulation tool should be able to simulate decision making process based on different approaches for example, regret minimization. In order to calculate utilities, we took parameters from the paper (*a priori*). All of the terms mentioned in the paper including ASC are used. We take no-choice as reference alternative. Firstly, a matrix of $4 \times 1000$ is constructed for socioeconomic characteristics parameters. Each column indicates a person , and all of the columns are similar since the parameters are constant for all of the people. Then, this matrix is pointly multiplied with the matrix of socioeconomic charachteristics matrix, and finally the sum of each column is the socioeconomic utility of each person and a vector of utility is achieved. Secondly, a matrix of $1000 \times 5$ is made to contain the parameters corresponding to the alternatives(price, label, carbon, label-carbon, constant). each row of this matrix is drawn from multivariate normal distribution, $\mu + L \times R$ where $\mu$ is a vector of means of parameters, L is derived from Cholesky decomposition($L \times L' = \sigma^2$) and $R$ is a vector of $K$ draws from a $N(0, 1)$. Finally, this matrix is multiplied by the inverse of design matrix which results in a matrix of $1000 \times 36$ in which each element shows the utility of an alternative for an individual. Consequently, we add up socioeconomic utility to each of the columns of this matrix. This makes the observed utility. In regard to unobserved utility, a matrix of $1000 \times 36$ is containing the draws of $Gumbel(0, 1)$, then we add this matrix to the previous one and this brings about the utility matrix. For the choice selection process, columns of utility matrix are compared pair by pair and also the max of these each of these paires is

compared with an element of $Gumbel(0, 1)$ to specify whether the individual buys or not. Finally, we suggust that tool decode and clean the data. One important issue is the difference between real data and simulated data which arises from ommited variables. For example, when it comes to reality, time is a very important factor affecting the choices made by respondents, but when it comes to simulation, time is meaningless for artificial individuals. These issues also need to be taken into account specially when we are comparing estimation results of these two types of data.

## II Reproducible research

This work was accomplished with implementation of the most advanced reproducible research techniques. First of all, a version control system (*git*) was used to track the changes and modification in the working tree from the start of the internship. The collaboration with other participants was organised through *GitHub*, where a common repository was maintained to store the data and document, as well as to keep every element of the code or text available to everyone. The report generation was automated with the use of a simplified markup language with embedded executable *R* code. For heavy tasks, such as data generation, model estimation or big data exploration separate source code files were used.

This short documents aims to introduce the reader to used research methodology, that was used in this work and during the internship. The used tool-set will be introduced.

*Git* is one of the version control tools alongside SVN and Mercurial-SCM, which allows to easily control changes and modifications within text documents. Unfortunately the proposed functionality does not function with more complex proprietary formats such as Word or image based documents, such as PDF. Consequently, this tool is not practical only for working with simple text documents: it remains absolutely impractical for working with typical office tasks. Several text editors for developers among which RStudio, VSCode, Atom and many other provide possibilities to integrate *git* functionality directly into the editor and drastically optimise the workflow. This makes interacting with *git* much more comfortable than through the command line or a standalone *git* client.

*GitHub* is an open source cooperative platform for developers offered by Microsoft making it easier to work with the *git* version control service. The platform has an entire ecosystem of extensions, expanding git functionality, as well as a set of project management and communication tools. In total this platform offers:

- A cloud space to host the working files and publish the results;
- A web interface to interact with *git* from browser or through a standalone app;
- A platform facilitating collaboration with other users, which gradually approaches in the functionality to a social network;
- An integrated project management system.

To write the scientific report it was decided to use the *LaTeX* complete markup language. There are several distributions of LaTeX, one of the verified versions to integrate well with *R* being *tinytex* (which is available as *tinytex* package in CRAN repositories). However, even if *LaTeX* produces well structured documents that are easy to manage, there exists the problem of its complexity in extending its functionalities. Consequently, it was decided to use an intermediary simplified markup language, which is easy to use and does not require advanced knowledge of *LaTeX*: Markdown. It allows to write documents with simple syntax, which could be later transformed into PDF, HTML and Word documents using the *pandoc* converter.

Finally, to embed the *R* code inside the document to automatically generate the figures and tables, we used *RMarkdown*, which is an extension for *Markdown* integrating *R* language inside. Such set-up ensured, that the documents will be easy to share and modify, preserving at the same time all their functionality. This work offers all the necessary elements to be fully reproducible.

**The resulting compedium is available at:** https://github.com/nikitagusarov/performance_exploration