

1 **WHAT IS REALLY UNCOVERED BY MIXING DIFFERENT MODEL STRUCTURES:**
2 **CONTRASTS BETWEEN LATENT CLASS AND MODEL AVERAGING**

3 **Thomas O. Hancock (Corresponding Author)**

4 Choice Modelling Centre & Institute for Transport Studies

5 University of Leeds

6 tstoh@leeds.ac.uk

7

8 **Stephane Hess**

9 Choice Modelling Centre & Institute for Transport Studies

10 University of Leeds

11 S.Hess@its.leeds.ac.uk

12

13 **ABSTRACT**

14 Latent class models have long been a tool for capturing heterogeneity across decision-makers in the
15 sensitivities to individual attributes. More recently, there has been increased interest in using them
16 to capture heterogeneity in actual behavioural processes, such as information/attribute processing
17 and decision rules. This often leads to substantial improvement in model fit and the apparent
18 finding of large clusters of individuals making choices in ways that are substantially different from
19 those used by others. Such findings have however not been without criticism given the potential
20 risk of confounding with other more model-specific heterogeneity. In this paper, we consider
21 an alternative approach for exploring the issue by contrasting the findings obtained with model
22 averaging, which combines the results from a number of separately (rather than simultaneously)
23 estimated models. We find that this leads to significant reductions in the amount of heterogeneity
24 of the type analysts have sought to uncover with latent class structures of late, with our results
25 suggesting that heterogeneity in the sensitivities to individual attributes rather than the behavioural
26 process per se could be the key factor behind the improvements gained through the adoption of
27 latent class models for heterogeneity in behavioural processes.

28 *Keywords: latent class; information processing; attribute non-attendance; decision rule heterogeneity*

29 **1. INTRODUCTION**

30 Over the last decade, there has been increasing interest by choice modellers to allow for departures
31 from traditional decision rules ([Chorus, 2014](#)) and/or the way in which individuals process the
32 information describing the alternatives ([Hensher, 2014](#)). Much of the work has looked at contrasts
33 between models using one specific alternative decision rule or process to the results against an
34 alternative model, i.e. fitting the same process to an entire sample of decision-makers. However, a
35 growing number of studies ([Hess and Rose, 2007](#); [Scarpa et al., 2009](#); [Hensher and Greene, 2010](#);
36 [Campbell et al., 2010](#); [Hole, 2011](#); [Hensher et al., 2012](#); [Hess et al., 2012](#)) have also looked at

1 allowing for heterogeneity in the actual underlying model structure across individuals in a single
2 sample. This has mainly made use of latent class structures, with two key applications, namely
3 decision rule heterogeneity and in information processing work.

4 While the work using latent class structures for heterogeneity in either decision rules or
5 information processing strategies has been shown to lead to substantial improvement ([Hess and
6 Rose, 2007](#); [Hess et al., 2012](#)) in fit and apparent meaningful insights, it has also not been without
7 criticism. In particular, concerns have been raised about the extensive risk of confounding between
8 heterogeneity in the sensitivities to individual attributes and heterogeneity in the process or model
9 structure.

10 In a traditional latent class model, the different β parameters in different classes are used
11 solely to uncover taste heterogeneity. In a latent class model that combines different structures in
12 different classes, these individual models will themselves be making use of different β parameters,
13 while in the case of ANA, they will use different combinations of the β parameters. For reasons of
14 complexity, the vast majority of applications have used just a single class per behavioural process,
15 whether that be one class for each decision rule (e.g. MNL, RRM, etc) or one class for each of the
16 combinations of considered attributes in an ANA context. Maximum likelihood estimation will
17 simply converge to those parameters that give the best mathematical fit to the data. For example,
18 imagine a situation where the decisions of all individuals in the data are best explained by a MNL
19 structure, but where there are variations across individuals in the sensitivity to for example the
20 cost attribute. If the analyst estimates a LC model with two classes, where one class uses a MNL
21 model and the other class uses a RRM model, then the only mechanism available to maximum
22 likelihood process for explaining the heterogeneity in cost sensitivities is to allocate a non-zero
23 class allocation probability to both the MNL and RRM classes, with different cost coefficients in
24 the two models. In other words, even in the absence of decision rule heterogeneity, the model
25 will *uncover* such heterogeneity if the benefit of being able to use different cost sensitivities in the
26 MNL and RRM classes outweighs the loss in fit of using a RRM model to explain the choices of
27 people who made decisions more in line with MNL. There is thus the real possibility that apparent
28 evidence of decision rule heterogeneity will be driven by heterogeneity in sensitivities rather than
29 actual decision process.

30 These concerns have found empirical support in the work of [Hess et al. \(2013b\)](#) who show
31 that the share for non-attendance classes reduces substantially when allowing for additional random
32 heterogeneity, while the work of [Hess et al. \(2016\)](#) shows that allowing for random heterogeneity
33 in the parameters of Random Utility Maximisation (RUM) and RRM models within a RUM-RRM
34 mixture model substantially reduces the extent of decision rule heterogeneity. The use in practice
35 of such latent class models allowing for different structures in different classes continues to be
36 very popular ([Boeri and Longo, 2017](#); [Dey et al., 2018](#)) despite these concerns. A key reason is
37 likely that the inclusion of additional taste heterogeneity (moving from finite latent class models
38 to continuous mixture models), as in the work of [Hess et al. \(2013b\)](#) and [Hess et al. \(2016\)](#) is
39 computationally very difficult. The same applies to the inclusion of additional classes with models
40 of the same type (e.g. using two MNL class and two RRM classes), where this also leads to a
41 proliferation in the number of parameters. Whilst there are of course other methods that can be
42 adopted to explore these and other kinds of heterogeneity, the key aim of the present paper is to

1 specifically consider a different approach to further examine the results of these latent class models
 2 which are so popular, but without increasing computational demands. We do this by highlighting
 3 how model averaging can be used as a diagnostic tool for the potential confounding between taste
 4 heterogeneity and other heterogeneity highlighted in these models.

5 Model averaging uses a sequential latent class approach, estimating first the individual candidate
 6 models at the sample level, before then combining these models in a latent class model which
 7 keeps the model-specific parameters fixed and only estimates the model weights. We illustrate this
 8 process on simulated data as well as typical stated preference data and show that model averaging
 9 can provide additional insights that could allow an analyst to reach a more informed decision as to
 10 the key drivers of heterogeneity in a model.

11 The aim of using model averaging in the present paper is to investigate potential cases of
 12 confounding in models using simultaneous estimation of different model structures. Of course,
 13 a caveat applies in that it is also possible that the presence of decision rule heterogeneity and/or
 14 heterogeneity in processing strategies can only be uncovered when estimating models in which
 15 the parameter estimates for the different subclasses are informed more by some individuals in
 16 the data than by others, as would be the case in simultaneous estimation. We address this point
 17 specifically by showing the possibility of including some models within the set M that themselves
 18 allow for heterogeneity. For example, it is straightforward to allow a given model m to be itself a
 19 LC or MMNL structure. As the model will be estimated separately rather than as part of the overall
 20 model averaging structure, we avoid the computation issues of including complex structures within
 21 an overall latent class structure.

22 The remainder of this paper is organised as follows. We first summarise the methodology in
 23 Section 2. This is followed by a simulated data experiment demonstrating how model averaging
 24 can help recover the original data generation process (Section 3). The core empirical work follows
 25 in Section 4, where we look both at attribute non-attendance and decision rule heterogeneity.
 26 Finally, some conclusions are presented in Section 5.

27 2. METHODOLOGY

28 Latent class (LC) structures have long been used as a tool for introducing heterogeneity across
 29 individual decision-makers in choice models (see [Greene and Hensher, 2003](#); [Hess, 2014](#), for
 30 background). In a latent class model, the population is probabilistically divided into S different
 31 classes, where the log-likelihood for the choices observed for a set of N decision-makers is given
 32 by:

$$33 \quad LL(\beta, \pi) = \sum_{n=1}^N \log \left(\sum_{s=1}^S \pi_{ns} \prod_{t=1}^{T_n} P_{j_{nt}^*}(\beta_s) \right), \quad (1)$$

34 where j_{nt}^* is the actual alternative observed to be chosen by person n in situation t . We have that
 35 $\pi = \langle \pi_1, \dots, \pi_N \rangle$, with π_n a vector whose element π_{ns} gives the share (probability) of individual
 36 n belonging to class s such that $\sum_{s=1}^S \pi_{ns} = 1, \forall n$ and $0 \leq \pi_{ns} \leq 1, \forall n, s$. These class allocation

1 probabilities can vary across individual decision-makers as a function of their characteristics, using
 2 a class allocation model, such that $\pi_n = f(\gamma, z_n)$ where z_n are characteristics of person n and γ is a
 3 vector of estimated parameters.

4 In almost all applications of latent class models, $P_{J_{nt}^*}(\beta_s)$ is of the Multinomial Logit (MNL)
 5 type. Even when this was not the case, for example using Nested Logit (NL) models inside a LC
 6 structure, the focus for the first two decades of widespread use of LC models was very much on
 7 a case where the functional form of $P_{J_{nt}^*}(\beta_s)$ is the same across classes (i.e. $s = 1, \dots, S$), with
 8 differences only in the parameters used in the classes, i.e. β_s in class s , where $\beta = \langle \beta_1, \dots, \beta_S \rangle$.
 9 This use of latent class models thus focusses on capturing what would typically be called "taste
 10 heterogeneity" while maintaining homogeneity in the underlying behavioural process across individual
 11 decision-makers.

12 Latent class models have more recently been used for heterogeneity in decision rules and
 13 information processing. While the former has received more attention, the latter work actually
 14 takes historical precedence.

15 A key interest in the field of information processing strategies (IPS) or attribute processing
 16 strategies (APS) has been the notion that some decision-makers may actually make their choices
 17 based on only a subset of the attributes that describe the alternatives at hand. This phenomenon
 18 is typically referred to as attribute non-attendance (ANA) or attribute ignoring, and an in-depth
 19 review of work in this area is given in [Hensher \(2010\)](#). The interest in this topic in the present
 20 discussions comes in the context of ways to accommodate ANA in models. A key role in this area
 21 was played by the early discussions in [Hess and Rose \(2007\)](#), who proposed the use of a latent
 22 class approach to accommodate ANA, a method since adopted by numerous other studies (e.g.
 23 [Scarpa et al., 2009](#); [Hensher and Greene, 2010](#); [Campbell et al., 2010](#); [Hole, 2011](#); [Hensher et al.,](#)
 24 [2012](#)). With this approach, different latent classes relate to different combinations of attendance
 25 and non-attendance across attributes. For each attribute treated in this manner, there exists a non-
 26 zero coefficient (to be estimated), which is used in the *attendance classes*, while the attribute is
 27 not employed in the *non-attendance classes*, i.e. the coefficient is set to zero. In a complete
 28 specification, covering all possible combinations, this would thus lead to 2^K classes, with K being
 29 the number of attributes, where a given coefficient will take the same value in all classes where
 30 that attribute is included.

31 In addition to the vector β , we now have a $S \times K$ matrix Λ , in which each row, s , contains a
 32 different combination of 0 and 1 elements, where $S = 2^K$. Next, let $A \circ B$ be the element-by-element
 33 product of two equally sized vectors A and B , yielding a vector C of the same size, where the k^{th}
 34 element of C is obtained by multiplying the k^{th} element of A with the k^{th} element of B . Using this
 35 notation, the specific values used for the taste coefficients in class s are then given by the vector
 36 $\beta_s = \beta \circ \Lambda_s$. The k^{th} element of the vector β_s is thus the k^{th} element of β if $\Lambda_{s,k} = 1$, and zero
 37 otherwise. The log-likelihood is then given by:

$$38 \quad LL(\beta, \pi) = \sum_{n=1}^N \log \left(\sum_{s=1}^S \pi_{ns} \prod_{t=1}^{T_n} P_{J_{nt}^*}(\beta_s = \beta \circ \Lambda_s) \right). \quad (2)$$

1 A different application of such heterogeneous structures in different classes has arisen in the
 2 context of decision rule heterogeneity. There has long been interest in the notion that different
 3 individuals make their decisions in different ways, going back to work in psychology in the 1970s
 4 (Montgomery and Svenson, 1976). Although structures belonging to the family of random utility
 5 models have come to dominate, it is important to recognise that alternative paradigms for decision-
 6 making have been proposed, for example the elimination by aspects model of Tversky (1972), but
 7 also more recent work based on the concepts of happiness (Abou-Zeid and Ben-Akiva, 2010) and
 8 regret (Chorus et al., 2008; Chorus, 2010). The evidence in the literature is that which paradigm
 9 works best is very much dataset specific.

10 Hess et al. (2012) put forward the hypothesis that variations in decision rules may be across
 11 decision-makers with a single dataset, not just across datasets, and propose the use of a confirmatory
 12 latent class approach in this context. Specifically, let $P_{jnt}^{(m)}(\beta_m)$ give the probability using a model
 13 of type m , with a vector of parameters β_m . The Hess et al. (2012) framework is based on the idea
 14 that different behavioural processes are used in the data. The original exposition by Hess et al.
 15 (2012) assumes that a different model type m is used in each class S , but this is not a requirement,
 16 and the same model structure could be used in more than one class. We then have:

$$17 \quad LL(\beta, \pi) = \sum_{n=1}^N \log \left(\sum_{s=1}^S \pi_{ns} \prod_{t=1}^{T_n} P_{jnt}^{m_s}(\beta_s) \right), \quad (3)$$

18 where m_s identifies the behavioural process used in class s , with β_s giving the vector of parameters
 19 used.

20 Hess et al. (2012) use the model to allow for mixtures between random utility maximisation,
 21 random regret minimisation (RRM) and elimination by aspects. They also discuss allowing for
 22 additional continuous random heterogeneity in parameters within individual classes, such that:

$$23 \quad LL(\Omega, \pi) = \sum_{n=1}^N \log \left(\sum_{s=1}^S \pi_{ns} \int_{\beta_s} \prod_{t=1}^{T_n} P_{jnt}^{m_s}(\beta_s) f(\beta_s | \Omega_s) d\beta_s \right), \quad (4)$$

24 where $\beta_{ns} \sim f(\beta_{ns} | \Omega_s)$ and $\Omega = \langle \Omega_1, \dots, \Omega_S \rangle$. In later work, Hess and Stathopoulos (2012) use an
 25 approach as in Walker and Ben-Akiva (2002) and Hess et al. (2013a), making the class allocation
 26 a function of a latent factor, which in this case also explains decision-makers' real world choices¹.

27 Model averaging, in this context, can be implemented as a sequential latent class model.
 28 Whereas a simultaneous model estimates the parameters of the class-specific models at the same
 29 time as the class allocation probabilities, a model averaging approach uses a sequential process.
 30 We first separately estimate the individual model from each class on the entire sample, before

¹At this stage, it should be noted that a latent class model mixing various decision rules is just one example of a wider set of structures that combine different models. A further possibility for example would be a model using different GEV nesting structures in different latent classes, somewhat similar in aims to the work of Ishaq et al. (2013). Finally, a separate body of work looks at using different choice sets in different classes, in the context of choice set generation work (see e.g. Swait and Ben-Akiva 1985; Ben-Akiva and Boccara 1995 and Gopinath 1995, section 2.7).

1 estimating the class allocation probabilities separately with the individual model parameters fixed.
 2 To apply model averaging, we thus first estimate a number of different individual models, where
 3 say $L_n^{(m)}(\Omega_m)$ gives the likelihood of the sequence of choices observed for person n , conditional
 4 on using model m , where this model uses a vector of parameters Ω_m . An analyst will estimate M
 5 different such models. Each model is estimated separately on the same data. Crucially, this implies
 6 that there needs to be some difference in the functional form between the different models, e.g.
 7 using different utility specifications, different mixing distributions, different attribute processing
 8 rules or indeed different decision rules. Indeed, any two models using the exact same structure
 9 will clearly converge to the same solution. Within-structure heterogeneity in sensitivities can
 10 easily be accommodated by some of the models being themselves LC or Mixed Multinomial
 11 Logit (MMNL) structures. In the context of the present paper, the set of M models would use
 12 different specifications for IPS or different specifications in terms of the underlying decision rules.
 13 The model averaging process then computes the overall likelihood for person n as the weighted
 14 average across M models, with the full sample log-likelihood given by:

$$LL(\Omega, \pi) = \sum_{n=1}^N \log \left(\sum_{m=1}^M \pi_{nm} L_n^{(m)}(\Omega_m) \right), \quad (5)$$

15 where $\sum_{m=1}^M \pi_{nm} = 1, \forall n$ and $0 \leq \pi_{nm} \leq 1, \forall n, m$. This overall log-likelihood is conditional on
 16 the vector of weights $\pi_n = \langle \pi_{n1}, \dots, \pi_{nM} \rangle$ for each person and the combined parameter estimates
 17 from the different models $\Omega = \langle \Omega_1, \dots, \Omega_M \rangle$. Crucially, a sequential estimation process is used.
 18 The parameters Ω_m are estimated separately by maximising the log-likelihood only for model m ,
 19 while the model weights are then estimated by maximising Equation 5 while keeping Ω fixed.

20 3. SIMULATED DATA ANALYSIS

21 Before testing model averaging on our stated preference dataset, we use simulated datasets to look
 22 at the contrasting insights provided by different approaches when the true process is known. We
 23 first describe how we created our 10 different datasets before applying MNL, RRM, latent class
 24 models and model averaging.

25 3.1. Generation of simulated data

26 We use an efficient design to generate 5,000 mode choice scenarios, each with four possible
 27 alternatives: car, air, rail and high-speed rail. These alternatives are described by (for respondent
 28 n , alternative i in choice scenario t) travel cost (TC_{nti}), travel time (TT_{nti}) and access time (AT_{nti}).
 29 For the MNL model, we define the utility as:

$$30 \quad V_{nti} = \delta_i + \delta_{F_i} \cdot z_{F,n} + \beta_{TT} \cdot \alpha_{TT_i} \cdot TT_{nti} + \beta_{TC} \cdot \alpha_{IE,n} \cdot TC_{nti} + \beta_{AT} \cdot AT_{nti} + \varepsilon_{nti}, \quad (6)$$

31 where δ_i and δ_{F_i} are alternative specific constants (with the constant for car normalised to zero),
 32 with δ_{F_i} only applying when the dummy variable for female respondents, $z_{F,n} = 1$ (which is the case
 33 for half of the participants). We have three marginal utility coefficients, β_{TT} , β_{TC} , β_{AT} , for travel
 34 time, travel cost and access time, respectively. We use car as the base for travel time sensitivity, and

1 apply mode-specific multipliers for travel time sensitivity through α_{TT_i} . Finally, we incorporate an
 2 income effect, which is defined as $\alpha_{IE,n} = (\frac{I_n}{2500})^{\alpha_I}$, where I_n is the income for individual n and α_I
 3 is an income elasticity. For the RRM model, we define the regret for an alternative equivalently,
 4 replacing $\beta_x \cdot X_{nti}$ with $\sum_{j \neq i}^4 \log(1 + \exp(\beta_x \cdot (X_{ntj} - X_{nti})))$. The constants which are added to the
 5 regret are also multiplied by -1 such that they have similar impacts on choice probabilities (as
 6 they do in MNL) by being applied in the same direction.

7 We use two sets of true parameter values for MNL, and two sets of true parameter values for
 8 RRM, with the parameter values given in Table 1. The first four datasets we create use the same
 9 data generation process (DGP) for all individuals and thus assume that all decision-makers use the
 10 same model with the same coefficients to make their choices (i.e. one dataset per column in Table
 11 1). A further six datasets are created (one for each pair of sets of coefficients) using a random
 12 allocation such that half of the decision-makers use one DGP, and the other half use a different
 13 DGP. The types of heterogeneity in each dataset are given below in Table 2.

TABLE 1 : Coefficient values for the four data generation processes, where the alternatives are car (C), rail (R), air (A) and high-speed rail (H), together with the choice shares for each mode. Note that $\delta_H = 0$, $\delta_{FH} = 0$ and $\alpha_{TTC} = 1$. Values are chosen for MNL1 and MNL2 such that they are very different. RRM1 and RRM2 have β -coefficients that are half the value of those of MNL1 and MNL 2 such that they have approximately the same scale.

Parameter	MNL1	RRM1	MNL2	RRM2
δ_C	-0.5000	-0.5000	1.0000	1.0000
δ_R	-1.5000	-1.5000	-0.5000	-0.5000
δ_A	-1.0000	-1.0000	1.0000	1.0000
δ_{FC}	-0.5000	-0.5000	-0.2000	-0.2000
δ_{FR}	0.5000	0.5000	1.5000	1.5000
δ_{FA}	1.0000	1.0000	-0.5000	-0.5000
α_I	-0.5000	-0.5000	-0.3000	-0.3000
β_{TT}	-0.0040	-0.0020	-0.0050	-0.0025
β_{TC}	-0.0280	-0.0140	-0.0100	-0.0050
β_{AT}	-0.0080	-0.0040	-0.0120	-0.0060
α_{TTR}	1.2500	1.2500	0.8000	0.8000
α_{TTA}	2.0000	2.0000	1.7000	1.7000
α_{TTH}	1.5000	1.5000	1.7000	1.7000
Share(Car)	53.06%	54.80%	43.12%	48.70%
Share(Rail)	9.54%	7.70%	6.36%	5.32%
Share(Air)	13.00%	11.76%	39.60%	33.66%
Share(HSR)	24.40%	25.74%	10.92%	12.32%

14 3.2. Results from simulated datasets

15 For each of the simulated datasets, we estimate five different models. The first two of these are
 16 basic MNL and basic RRM models. The next three are latent class models with two classes, using

TABLE 2 : Types of heterogeneity in each dataset

Dataset	DGP		Heterogeneity type	
	50%	50%	Taste	Decision rule
1	MNL1		no	no
2	MNL2		no	no
3	RRM1		no	no
4	RRM2		no	no
5	MNL1	MNL2	yes	no
6	MNL1	RRM1	no	yes
7	MNL1	RRM2	yes	yes
8	MNL2	RRM1	yes	yes
9	MNL2	RRM2	no	yes
10	RRM1	RRM2	yes	no

1 all three combinations of models, i.e. MNL-MNL, MNL-RRM and RRM-RRM. This means that
 2 for each dataset, we have two models that test for taste heterogeneity alone (MNL-MNL and RRM-
 3 RRM) and one model that allows for taste heterogeneity and decision rule heterogeneity (MNL-
 4 RRM). As we have all variations of datasets that include taste and/or decision rule heterogeneity
 5 covered by our 10 datasets, we aim to test the findings obtained using simultaneous LC models
 6 and sequential LC models, i.e. model averaging, and contrast these with the true DGP.

7 The results of all models are given in Table 3. We each time highlight the best BIC for each
 8 DGP, along with the noting whether there is a clear “winner” in model averaging or a more even
 9 split between models.

10 We first look at the four datasets which come from a single model without heterogeneity
 11 (i.e. models 1-4). We note that each time, the model using the correct decision rule outperforms its
 12 counterpart, while none of the three LC structures can reject the single class model. The differences
 13 in fit between MNL and RRM are small, in line with many previous findings in the literature, at
 14 least with the most common RRM implementation.

15 We notice that there are four datasets for which there is clear evidence of heterogeneity, with
 16 datasets 5, 7, 8 and 10 having improvements in model fit of at least 150 log-likelihood units by
 17 moving from a basic MNL or basic RRM model to a latent class model. These are the four datasets
 18 where we either mix two models of the same type (5 and 10) or where the models of different types
 19 also use substantially different relative taste coefficients (7 and 8, noting that the parameters for
 20 MNL1 and RRM1 are consistent with each other, as are MNL2 and RRM2). Crucially, we also
 21 see that for each of the four cases, the LC structure that is in line with the DGP performs best.
 22 However, the differences in fit are very small, again reiterating the point about similarities between
 23 MNL and RRM.

24 We finally look at the two models that combine different decision rules but use consistent
 25 parameters in the two, i.e. focussing on decision rule rather than taste heterogeneity. These are

TABLE 3 : Results from MNL, RRM and latent class models for each of the 10 simulated datasets, together with the results from model averaging.

Dataset	1	2	3	4	5	6	7	8	9	10
Data creation:	MNL1	MNL2	RRM1	RRM2	MNL1 MNL2	MNL1 RRM1	MNL1 RRM2	MNL2 RRM1	MNL2 RRM2	RRM1 RRM2
Log-likelihoods (LL(0) = -6,931.47)										
MNL	-4,816.03	-5,043.07	-4,641.41	-4,986.07	-5,357.92	-4,736.56	-5,227.40	-5,293.52	-5,037.90	-5,203.48
RRM	-4,841.62	-5,046.22	-4,631.31	-4,981.98	-5,368.51	-4,743.00	-5,239.64	-5,299.49	-5,038.91	-5,203.23
MNL-MNL	-4,804.49	-5,028.90	-4,637.99	-4,976.56	-5,101.61	-4,728.18	-5,047.94	-5,040.73	-5,021.55	-5,060.08
MNL-RRM	-4,807.08	-5,039.74	-4,622.68	-4,976.86	-5,114.58	-4,727.61	-5,046.22	-5,036.22	-5,027.76	-5,055.41
RRM-RRM	-4,828.56	-5,034.32	-4,625.27	-4,969.88	-5,115.40	-4,734.96	-5,060.76	-5,041.17	-5,027.61	-5,052.90
LC2-LC1	11.53	14.17	8.63	12.10	256.31	8.95	181.17	257.30	16.35	150.32
Bayesian Information Criterion (BIC)										
MNL	9,742.77	10,196.86	9,393.54	10,082.87	10,826.56	9,583.84	10,565.51	10,697.76	10,186.52	10,517.68
RRM	9,793.96	10,203.17	9,373.34	10,074.68	10,847.75	9,596.72	10,589.99	10,709.70	10,188.55	10,517.18
MNL-MNL	9,838.95	10,287.77	9,505.95	10,183.09	10,433.17	9,686.33	10,325.84	10,311.42	10,273.06	10,350.12
MNL-RRM	9,844.12	10,309.45	9,475.32	10,183.68	10,459.12	9,685.19	10,322.41	10,302.41	10,285.49	10,340.79
RRM-RRM	9,887.07	10,298.61	9,480.50	10,169.71	10,460.76	9,699.89	10,351.48	10,312.30	10,285.18	10,335.77
Model Averaging Share										
MNL	99.93%	87.03%	13.97%	9.57%	0.00%	71.96%	0.04%	0.00%	60.92%	0.01%
RRM	0.07%	12.97%	86.03%	90.43%	0.00%	28.04%	0.00%	0.00%	39.08%	0.00%
MNL-MNL	-	-	-	-	99.98%	-	9.55%	15.18%	-	10.29%
MNL-RRM	-	-	-	-	0.01%	-	90.34%	84.52%	-	2.04%
RRM-RRM	-	-	-	-	0.01%	-	0.07%	0.30%	-	87.67%
Model Averaging Log-likelihood										
LL	-4,816.03	-5043.00	-4631.02	-4981.93	-5,101.61	-4,735.38	-5,046.20	-5,036.08	-5,037.17	-5,052.79
Improvement	0.00	0.07	0.29	0.04	0.00	1.17	0.02	0.14	0.73	0.11

1 models 6 and 9. For these two cases, the improvements in model fit for the latent class structures
2 are very modest, to the point of not passing a statistical test for rejecting the single class model. In
3 fact, the improvements found for datasets 6 and 9 are comparable to those found by datasets 1-4.
4 This finding again suggests that MNL and RRM are very similar models when using parameters
5 that are consistent across structures. A small gain in log-likelihood by moving to latent class
6 models results in a worse BIC value, and thus we observe that the best BIC for datasets without
7 taste heterogeneity is found by either the basic MNL model (for datasets generated by a single
8 MNL or by MNL and RRM) or the basic RRM model (for datasets generated by a single RRM).
9 We can thus reject the latent class models in these cases and apply model averaging over only the
10 basic MNL and basic RRM models. For models with taste heterogeneity, however, we apply model
11 averaging across all latent class models as well as the MNL and RRM models.

12 The results from model averaging are also given in Table 3. For the first four datasets, we
13 observe a large share (minimum 86%) given to the model that created the data, as we would expect.
14 For the two remaining datasets without taste heterogeneity (datasets 6 and 9), we observe a more
15 equal split in the shares, thus implying the presence of decision rule heterogeneity. The results for
16 the datasets with taste heterogeneity also clearly match the data generation process. MNL-MNL
17 obtains almost the full share for the dataset created by two MNLS, and RRM-RRM obtains 88%
18 of the share for the dataset created by two RRMs. Additionally, we observe large shares for the
19 MNL-RRM latent class model (90% and 85%, respectively) for the two datasets created by MNL
20 and RRM models where there was also taste heterogeneity. It is crucial here to note that applying
21 model averaging across all latent class models in all cases may result in unclear results, as latent
22 class models without substantial improvements in model fit may only be picking up noise in the
23 data. Thus averaging across these models will highlight models that only explain noise, rather than
24 the underlying data generation process. Whilst we may see a greater improvement in model fit
25 through applying model averaging across larger sets of candidate models utilising a wider range of
26 decision rules, the results here demonstrate that model averaging can help differentiate some cases
27 where taste heterogeneity alone exists. Additionally, model averaging helps us when taste and
28 decision rule heterogeneity exist, with only small improvements in model fit for the MNL-RRM
29 model (1.72 and 4.51 log-likelihood units in DGP 7 and 8, respectively) over the MNL-MNL
30 model, but clear advantages for the MNL-RRM model in terms of the model averaging shares.

31 4. ANALYSIS ON SP DATA

32 This section presents our work on a typical SP dataset. We first give an overview of the data
33 before looking separately at the case of attribute non-attendance (cf. Section 4.2) and decision rule
34 heterogeneity (cf. Section 4.3).

35 4.1. Data

36 Our main analysis relies on a SC dataset from Hess and Stathopoulos (2013) where public transport
37 commuters living in the UK each make ten choices between three routes. A total of 368 participants
38 completed the survey resulting in 3,680 choices. Each choice task involves an invariant reference
39 trip and two hypothetical alternatives (where the invariant trip is chosen 35.19% of the time and

1 the new alternatives have shares of 34.27% and 30.54%, respectively). The invariant trip for each
 2 individual is based on averaging trip attributes across 10 regular trips corresponding to a week
 3 of commuting, with the attributes of the hypothetical alternatives being pivoted around those of
 4 the invariant trip. These choice tasks were generated with a D-efficient experimental design using
 5 NGene. A total of 60 choice scenarios were blocked into groups of 10. Further details for the
 6 dataset are given by [Hess and Stathopoulos \(2013\)](#). Each alternative is described by travel time
 7 (in minutes), fare (in £), rate of crowded trips, rate of delays (both out of 10 trips), the average
 8 length of delays (across delayed trips) and the presence of a delay information service (either not
 9 available, available at a small fixed cost, or free). This dataset has previously been used for decision
 10 rule heterogeneity ([Hess and Stathopoulos, 2013](#)) as well as for ANA work ([Hess et al., 2013b](#)),
 11 making it an ideal case study for the present paper.

12 4.2. Attribute non-attendance work

13 We first look at the case of ANA, where we adopt a specification in line with [Hess et al. \(2013b\)](#).

14 4.2.1. Specification

15 We start by estimating a simple MNL model, where we use a logarithmic transform on the fare
 16 attribute given earlier evidence of strong non-linearity. This model uses five marginal utility
 17 parameters for the continuous attributes, two parameters for the dummy coded delay information
 18 system, and two alternative specific constants (ASC).

19 We next move to the latent class model for attribute non-attendance. We use a model with
 20 2^K classes, with all combinations of attendance and non-attendance for the K parameters. The
 21 probability for class s is given by π_s , with $0 \leq \pi_s \leq 1$ and $\sum_{s=1}^S \pi_s = 1$. Rather than imposing
 22 constraints in estimation, an easier approach is to use $\pi_s = \frac{e^{\delta_s}}{\sum_{m=1}^S e^{\delta_m}}$, with one δ_m , i.e. the parameter
 23 used in the class allocation probabilities, being fixed to zero. Nevertheless, this specification still
 24 involves estimating $2^K - 1$ separate δ terms, of which many will be very negative, equating to
 25 very small class probabilities. In the context of the applications presented in this paper, we make
 26 the simplifying assumption that attendance versus non-attendance is independent across attributes
 27 (with probabilities that vary across attributes but are constant across individuals), by instead setting

$$28 \quad \pi_s = \prod_{k=1}^K (\Lambda_{s,k} (1 - P_{NA,k}) + (1 - \Lambda_{s,k}) P_{NA,k}), \quad (7)$$

29 where $\Lambda_{s,k}$ gives the entry in Λ relating to attribute k in class s , where this is 1 only if attribute k is
 30 attended to in class s . With this specification, we only need to estimate K separate δ elements (as
 31 $P_{NA,k}$ is the probability of non-attendance to attribute k , thus $P_{NA,k} = \frac{e^{\delta_k}}{e^{\delta_k} + 1}$), as opposed to $2^K - 1$,
 32 leading to significant reductions in the number of parameters.

33 We finally look at the estimation of our model averaging structure. For this, we first estimate
 34 128 individual models, corresponding to all possible combinations of attribute attendance and non-
 35 attendance, i.e. going from a model with all 9 model parameters (all 7 attributes are attended to)

1 to one with the two alternative specific constants only (none of the attributes attended to). We
 2 then estimate the model averaging structure, meaning that we keep the parameters for each of the
 3 128 models at the estimates from the individual model estimation process and only estimate the
 4 weights for model averaging. We again use multiplicative class allocation probabilities, as in the
 5 LC model.

6 4.2.2. Results

7 The results for the simple MNL model are shown in Table 4 where all estimates are of the expected
 8 sign.

TABLE 4 : MNL results for public transport route choice

LL(0)	-4,042.89	
LL(final)	-3,366.95	
ρ^2	0.1672	
adj. ρ^2	0.1650	
	Estimate	Rob.t.ratio(0)
ASC_1	0.3841	5.76
ASC_2	0.1608	3.26
β_{tt}	-0.0467	-9.47
$\beta_{\log\text{-fare}}$	-5.9726	-18.89
β_{crowding}	-0.2198	-8.51
$\beta_{\text{rate of delays}}$	-0.2411	-9.82
$\beta_{\text{average delay}}$	-0.0421	-5.35
$\beta_{\text{info system charged}}$	-0.0833	-1.04
$\beta_{\text{info system free}}$	0.3370	5.06

9 The results for the confirmatory LC model are shown in Table 5. We see an improvement in log-
 10 likelihood by 308.16 units for 7 additional parameters. This is highly significant and in line with
 11 previous findings when using such a confirmatory latent class model for ANA. We also see that
 12 the marginal utility parameters, which now only apply in the attendance classes, have increased
 13 substantially compared to the base model. This is in line with the notion that the MNL model
 14 would find an intermediary value between 0 for the non-attenders and a positive value for those
 15 attending to the attribute. However, the implied rates of non-attendance are unrealistically high,
 16 exceeding 50% for all attributes except fare.

17 For model averaging, we initially estimated seven class allocation weights as in the LC model but
 18 find that for the first four attributes, the constants go towards $-\infty$, suggesting a zero probability of
 19 ANA. The results of the model averaging work are shown in Table 6. We see that this model now
 20 only offers a marginally better log-likelihood than the MNL model in Table 4, much in contrast
 21 with the LC model in Table 5. No formal statistical test is used here as model averaging is not a
 22 process of simultaneously estimating all the parameters for all the models on a single dataset. In

TABLE 5 : Confirmatory latent class model for attribute non-attendance

LL(0)	-4,042.89	
LL(final)	-3,058.79	
ρ^2	0.2434	
adj. ρ^2	0.2395	
	Estimate	Rob.t.ratio(0)
ASC_1	0.8416	10.32
ASC_2	0.329	4.23
β_{tt}	-0.1841	-5.64
$\beta_{\log\text{-fare}}$	-14.6889	-14.37
β_{crowding}	-1.1524	-7.16
$\beta_{\text{rate of delays}}$	-1.1307	-5.62
$\beta_{\text{average delay}}$	-0.3966	-4.85
$\beta_{\text{info system charged}}$	2.3264	3.37
$\beta_{\text{info system free}}$	2.0433	7.23
$\delta_{NA,tt}$	0.3232	1.11
$\delta_{NA,\log\text{-fare}}$	-0.5142	-3.43
$\delta_{NA,\text{crowding}}$	0.7767	3.3
$\delta_{NA,\text{rate of delays}}$	0.7363	2.43
$\delta_{NA,\text{average delay}}$	1.1917	4.02
$\delta_{NA,\text{info system charged}}$	3.1776	3.82
$\delta_{NA,\text{info system free}}$	0.9874	3.61
	Implied rate of NA	
	Estimate	Rob.t.ratio(0)
travel time	0.5801	8.18
fare	0.3742	10.65
crowding	0.685	13.49
rate of delays	0.6762	10.21
average delay	0.767	14.48
info system charged	0.96	30.05
info system free	0.7286	13.47

1 addition to the earlier finding of zero weight for any classes that imply non-attendance of either
 2 time, fare, crowding or the rate of delays, we also see low rates for the average delay and the free
 3 information system, with a higher rate for the charged system. A number of other statistics are
 4 valuable. First, we can rank the 128 models by log-likelihood and we note that the 8 models that
 5 obtain the best individual log-likelihoods are also the only 8 models that contribute to the model
 6 average. The two best fitting models also contribute the most to the model averaging, though in
 7 reverse order (models with rankings 2 and 1). Finally, for each individual person in the data, we
 8 can see which of the 128 models best explains their choices. Doing this, we see that out of the 368
 9 individuals in the data, only 95 have their choices explained the best way by one of these 8 models,

TABLE 6 : Model averaging for ANA work

LL(final) = -3,363.28, LL(0) = -4,042.89		Implied rate of NA			
	Estimate	Rob.t.ratio(0)	Estimate	Rob.t.ratio(0)	
$\delta_{NA,average\ delay}$	-1.9099	-1.95	average delay	0.129	1.17
$\delta_{NA,ch\ inf\ sys}$	0.0844	0.05	info system charged	0.5211	1.22
$\delta_{NA,free\ inf\ sys}$	-1.1531	-2.04	info system free	0.2399	2.33

Information for 8 retained models									
	LL	1	2	3	4	5	6	7	8
ranking out of 128 candidates	-3,367.75	-3,366.95	-3,400.98	-3,390.17	-3,391.85	-3,391.62	-3,424.48	-3,416.22	
providing best fit for N respondents	12	17	14	14	9	8	12	9	
MA share	34.50%	31.71%	10.89%	10.01%	5.11%	4.70%	1.61%	1.48%	

	attribute included								
travel time	YES	YES	YES	YES	YES	YES	YES	YES	YES
fare	YES	YES	YES	YES	YES	YES	YES	YES	YES
crowding	YES	YES	YES	YES	YES	YES	YES	YES	YES
rate of delays	YES	YES	YES	YES	YES	YES	YES	YES	YES
average delay	YES	YES	YES	YES	YES	NO	NO	NO	NO
info system charged	NO	YES	NO	YES	NO	NO	YES	YES	NO
info system free	YES	YES	NO	NO	YES	YES	NO	YES	NO

	est. (rob. t-rat)									
ASC ₁	0.41 (6.46)	0.38 (5.76)	0.40 (6.24)	0.32 (4.77)	0.39 (6.15)	0.38 (5.61)	0.38 (5.91)	0.31 (4.61)		
ASC ₂	0.16 (3.29)	0.16 (3.26)	0.16 (3.29)	0.16 (3.17)	0.18 (3.59)	0.18 (3.58)	0.17 (3.46)	0.17 (3.41)		
β_t	-0.05 (-9.48)	-0.05 (-9.47)	-0.05 (-9.73)	-0.05 (-9.63)	-0.05 (-9.35)	-0.05 (-9.34)	-0.05 (-9.59)	-0.05 (-9.49)		
$\beta_{log-fare}$	-5.95 (-18.86)	-5.97 (-18.89)	-5.77 (-18.19)	-5.90 (-18.62)	-5.87 (-18.81)	-5.88 (-18.81)	-5.68 (-18.12)	-5.80 (-18.51)		
$\beta_{crowding}$	-0.22 (-8.50)	-0.22 (-8.51)	-0.22 (-8.59)	-0.22 (-8.61)	-0.22 (-8.46)	-0.22 (-8.46)	-0.22 (-8.55)	-0.22 (-8.56)		
$\beta_{rate\ of\ delays}$	-0.24 (-9.76)	-0.24 (-9.82)	-0.24 (-9.82)	-0.24 (-9.95)	-0.27 (-10.94)	-0.27 (-10.98)	-0.26 (-11)	-0.27 (-11.17)		
$\beta_{average\ delay}$	-0.04 (-5.32)	-0.04 (-5.35)	-0.04 (-5.29)	-0.04 (-5.51)	0	0	0	0		
$\beta_{ch\ inf\ sys}$	0	-0.08 (-1.04)	0	-0.27 (-3.67)	0	-0.04 (-0.57)	0	-0.24 (-3.24)		
$\beta_{free\ inf\ sys}$	0.36 (5.96)	0.34 (5.06)	0	0	0.36 (5.91)	0.35 (5.22)	0	0		

1 where a remarkable 104 out of the 128 models have at least one individual where they are the best
2 performing model.

3 Overall, the findings from this analysis are much in contrast with those from the confirmatory
4 latent class model in that very little evidence of ANA is found. In addition, there is very little
5 variation in the remaining parameters across classes. Of course, the counter-argument could be
6 that the model averaging approach cannot retrieve ANA as it is based on individual models that
7 each apply a homogeneous approach to all individuals. However, some reassurance can be obtained
8 from the fact that the model averaging results are in line with the findings by [Hess et al. \(2013b\)](#)
9 which find evidence of ANA only for the average delay attribute and for the delay information
10 attribute after allowing for random heterogeneity in their models. It is thus doubtful whether
11 additional insights would be obtained with more flexibility for the individual models, such as by
12 including random heterogeneity. A possible step in that direction would be to estimate one latent
13 class model for each of the 128 candidates, i.e. allowing for heterogeneity within a model that
14 assumes a given ANA strategy.

15 4.3. Decision rule heterogeneity work

16 We next turn to decision rule heterogeneity, which has been the key interest in applying latent class
17 structures for process heterogeneity in recent years.

18 4.3.1. Specification

19 To maximise the possibility of finding such heterogeneity, we specifically choose to consider five
20 very different decision rules, namely:

21 **Multinomial logit (MNL):** We assume that the utility a respondent n obtains from alternative i
22 (out of J alternatives) in choice task t is:

$$V_{nti} = U_{nti} + \varepsilon_{nti}, \quad (8)$$

23 where V_{nti} and ε_{nti} are the deterministic and random components of utility respectively. The
24 assumption of a type I extreme value distribution for ε_{nti} then gives us the usual MNL choice
25 probabilities:

$$P_{MNL,nti} = \frac{e^{V_{nti}}}{\sum_{j=1}^J e^{V_{ntj}}}. \quad (9)$$

26 **Random regret minimisation (RRM):** We base our random regret minimisation (RRM) model
27 on the updated specification of [Chorus \(2010\)](#). Thus, the deterministic regret for respondent
28 n for alternative i in choice task t is given by:

$$R_{nti} = \delta_{RRM,i} + \sum_{j \neq i}^J \sum_{k=1}^K \ln(1 + e^{\beta_k(x_{ntjk} - x_{ntik})}) \quad (10)$$

1 where $k = 1, \dots, K$ is an index across attributes, β_k is an attribute-specific coefficient for
 2 attribute k and $\delta_{RRM,i}$ is an alternative specific constant. With the error component of regret
 3 also being given by a type I extreme value distribution, the corresponding RRM probabilities
 4 for a respondent n choosing alternative i in choice task t is:

$$P_{RRM,nti} = \frac{e^{-R_{nti}}}{\sum_{j=1}^J e^{-R_{ntj}}} \quad (11)$$

5 **Decision field theory (DFT):** DFT is a dynamic, stochastic model where the preferences for alternatives
 6 update over the course of the decision-making process (Busemeyer and Townsend, 1992).
 7 Whilst it has not traditionally been used to model travel behaviour, work by Roe et al.
 8 (2001) extended the model such that it can be applied to multi-attribute, multi-alternative
 9 choice tasks and Berkowitsch et al. (2014) demonstrated that DFT can outperform logit and
 10 probit models for consumer choice tasks. More recently, Hancock et al. (2018) and Hancock
 11 et al. (2020c) have demonstrated that DFT can outperform standard models for typical travel
 12 behaviour data. It is mathematically a very different model to MNL or RRM as it assumes
 13 that a decision-maker stochastically considers the different attributes of the alternatives over
 14 the course of a decision-making process for a single choice task. In each of these steps, a
 15 single attribute is considered - note that this is different from ANA in that all attributes are
 16 attended, but one at a time, and with different attributes being attended to different numbers
 17 of times as a function of their importance. This results in the preference values updating
 18 iteratively:

$$P_\tau = S \cdot P_{\tau-1} + V_\tau, \quad (12)$$

19 where P_τ is a column vector containing the preference values of each alternative i at time τ .
 20 S is a feedback matrix with memory and sensitivity parameters² and V_τ is a valence vector,
 21 which is the preference accumulated in the step τ and will depend on which attribute is
 22 attended to in that step. The valence vector can be described by:

$$V_\tau = C \cdot M \cdot W_\tau + \varepsilon_\tau \quad (13)$$

23 where C is a contrast matrix used to rescale the values such that they total zero, M is the
 24 matrix of attribute values and $W_\tau = [0..1..0]'$ with entry $k = 1$ if and only if attribute k is
 25 the attribute being attended to by the decision-maker at preference updating step τ . A DFT
 26 model thus estimates a weight, w_k , for the likelihood of attending to attribute k . As the
 27 error term, ε_τ is drawn from a normal distribution with mean 0 (and a variance which is an
 28 estimated parameter), the preference values P_τ converge to a multivariate normal distribution.
 29 To calculate the probabilities of alternatives under decision field theory we thus simply
 30 require the expectation and covariance of P_τ (ξ_τ and Ω_τ , respectively). Hence the probability
 31 of choosing alternative i from a set of J alternatives at time τ is (dropping the index for choice
 32 set and individual):

$$P_{DFT,i} \left[\max_{j \in J} P_\tau[j] = P_\tau[i] \right] = \int_{X>0} \exp \left[-(X - \Gamma)' \Lambda^{-1} (X - \Gamma) / 2 \right] / (2\pi |\Lambda|^{0.5}) dX \quad (14)$$

²We implement the feedback matrix equivalently to Hancock et al. (2020c).

with $X = [P_\tau[j] - P_\tau[1], \dots, P_\tau[j] - P_\tau[J]]'$, $\Gamma = L\xi_\tau$, $\Lambda = L\Omega_\tau L'$ and L a matrix comprised of a column vector of 1s and a negative identity matrix of size $J - 1$ where J is the number of alternatives. The column vector of 1s is placed in the i^{th} column where i is the chosen alternative. The DFT model utilised in the empirical tests in this paper is based on the version in [Hancock et al. \(2020c\)](#), which incorporates attribute-specific scaling factors.³

Quantum amplitude model (QAM) Our quantum model is based on the quantum amplitude model defined by [Hancock et al. \(2020a\)](#). Under a quantum model, the possible choice alternatives can be represented by a set of orthogonal vectors which make up the basis for a multidimensional Hilbert space ([Bruza et al., 2015](#)). A decision-maker's opinion or 'belief state' can then be represented by another vector within this space. The action of making a choice is represented by a projection from the belief state vector onto the vector representing the chosen alternative (see figures in [Hancock et al. 2020a](#)). Setting the belief state vector to be of unit length results in the set of squared 'lengths' (amplitudes) of these projection onto each of the vectors representing the possible alternatives summing to one. This means that, rather than the probability of alternatives being based on utilities with some error terms, they are based on the belief state vector. Each alternative has an amplitude rather than a utility, and the error is inherent in the belief state vector itself, which directly captures the fact that a decision-maker may choose any of the alternatives. Evidence in favour of an alternative increases the relative 'lengths' of the vector for the alternative, resulting in the probability for this alternative increasing. Consequently, the belief state vector captures the fact that a decision-maker may choose different alternatives without the need for an additional error term. There are many reasons for adopting quantum frameworks ([Bruza et al. 2015](#)), but most importantly in the context of travel behaviour analysis, they are flexible models whereby different value functions can be integrated into the structure together with 'quantum rotations' or 'changes in perspectives' for a change in choice context. In particular, [Hancock et al. \(2020a\)](#) demonstrate that the quantum amplitude model can outperform traditional choice models for standard travel behaviour data. We thus include it as an additional model in this paper to further broaden the possibility of finding decision rule heterogeneity, given that it is, like DFT, a mathematically very different choice model to MNL or RRM.

Under QAM, there are a number of possible functions that can be used to define the amplitudes for each alternative, with [Hancock et al. \(2020a\)](#) giving a detailed overview and empirical tests of these possibilities. In this case, we consider amplitudes based on 'softplus' (random regret-like) functions. For alternative i (for respondent n in choice task t), we define the amplitude:

$$\psi_{nti} = \left(\delta_{QAM,i} + I_0 + \sum_{j \neq i} \sum_{k=1}^K wt_{ij} \cdot \ln(1 + e^{\beta_k(x_{ntik} - x_{ntjk})}) \right) / \sqrt{\mathcal{N}}, \quad (15)$$

where $\delta_{QAM,i}$ are alternative-specific constants, I_0 is a constant that has the same value across all alternatives, wt_{ij} is a weight for the relative importance of the comparison between alternatives i and j and β_k is a coefficient for attribute k as before for RRM. \mathcal{N} is a normalisation

³For a full description of decision field theory, how it can be applied and how the different parameters in the model work, readers should consult [Hancock et al. \(2018\)](#) and [Hancock et al. \(2020c\)](#).

1 factor, which ensures that the probabilities with which each alternative is chosen sum to one.
 2 It is obtained from the sum of the squared moduli:

$$\mathcal{N} = \sum_i^J \left| \left(\delta_{QAM,i} + \sum_{j \neq i}^J \sum_{k=1}^K w_{t_{ij}} \cdot \ln(1 + e^{\beta_k(x_{ntik} - x_{ntjk})}) \right) \right|^2, \quad (16)$$

3 where $i = 1, \dots, J$ is an index across the possible alternatives. Once these amplitudes have
 4 been calculated, the application of the normalisation means that the probability for each
 5 alternative⁴ can be defined simply as:

$$P_{QAM,jnt} = \Psi_{jnt}^2. \quad (17)$$

6 **Relative advantage maximisation (RAM)** In RAM (Leong and Hensher, 2014), the utility for
 7 respondent n in choice task t is:

$$U_{nti} = \delta_{RAM,i} + \sum_{k=1}^K \beta_k \cdot x_{ntik} + \sum_{j \neq i} RA(i, j), \quad (18)$$

8 which is equivalent to a multinomial logit model with the addition of the comparison of
 9 relative advantages $RA(i, j)$ of alternative i with each of the other alternatives. This relative
 10 advantage is then defined as:

$$RA(i, j) = \frac{A(i, j)}{A(i, j) + D(i, j)}, \quad (19)$$

11 where the advantages are calculated $A(i, j) = \ln(1 + e^{\beta_k(x_{ntik} - x_{ntjk})})$ and the disadvantages
 12 $D(i, j) = \ln(1 + e^{\beta_k(x_{ntjk} - x_{ntik})})$.

13 4.3.2. Results

14 For our SP dataset, we first apply the five different models individually, obtaining the results given
 15 in Table 7. We see that DFT obtains the best log-likelihood ahead of QAM, while the performance
 16 of the three logit-style models is noticeably poorer and comparatively more similar. As a first
 17 step, we look at model averaging across these five individual models with different decision rules,
 18 where the resulting shares and fit are shown in Table 7. We see that the model average leads to
 19 a further small improvement in model fit over the best fitting individual model, i.e. DFT, where
 20 this model also obtains by far the largest share in the model average. As with earlier examples,
 21 the shares are not necessarily proportional to the model fit of the individual model, and we see
 22 that RRM obtains a substantially larger share than QAM, despite having poorer overall individual
 23 log-likelihood. This again shows that some models can work well for some people even if they
 24 obtain a lower overall fit to the sample.

25 In practice, the estimation of a latent class model with five separate classes all using individual
 26 decision rules is computationally challenging and most applications rely on just combining a

⁴Note that as we are using real values only in the specification of the amplitudes, squaring these amplitudes for the probability suffices. If we had complex valued amplitudes, we would require the use of complex conjugates.

TABLE 7 : Results from different individual models applied to the SP dataset

Model	Type	Log-likelihood	BIC	MA Share
1	MNL	-3,360.43	6,803	0.00%
2	RRM	-3,363.91	6,810	17.67%
3	DFT	-3,317.18	6,749	76.54%
4	QAM	-3,336.44	6,771	5.70%
5	RAM	-3,354.55	6,791	0.08%
Model averaging		-3,312.40		

1 couple of different rules. We therefore look at the estimation of 15 different latent class structures
 2 with two classes each, picking all combinations of two model structures with replacement, thus
 3 also allowing for five models where the two classes are of the same type, i.e. looking for taste
 4 heterogeneity alone. Table 8 gives the log-likelihoods of these models. For all 15 models, a
 5 likelihood ratio test against the corresponding model (in the case of single decision rule) or two
 6 corresponding models (in the case of two decision rules) clearly rejects the base model. This
 7 provides evidence of taste heterogeneity (in the case of single structure models) and would typically
 8 be seen as evidence of decision rule heterogeneity in the case of the models with two different
 9 structures in the two classes.

TABLE 8 : Results from latent class models applied to SP dataset

Model	Class 1	Class 2	Log-likelihood	BIC	MA1 Share	MA2 Share
1	MNL	MNL	-3,113.13	6,399	0.0%	0.0%
2	MNL	RRM	-3,102.66	6,378	0.0%	-
3	MNL	DFT	-3,099.84	6,380	6.7%	-
4	MNL	QAM	-3,106.76	6,394	0.0%	-
5	MNL	RAM	-3,100.79	6,374	0.0%	-
6	RRM	RRM	-3,106.33	6,385	16.0%	22.8%
7	RRM	DFT	-3,086.79	6,354	11.7%	-
8	RRM	QAM	-3,096.35	6,373	0.0%	-
9	RRM	RAM	-3,104.22	6,381	0.0%	-
10	DFT	DFT	-3,077.79	6,361	52.8%	62.3%
11	DFT	QAM	-3,085.28	6,376	0.0%	-
12	DFT	RAM	-3,085.38	6,351	0.0%	-
13	QAM	QAM	-3,095.71	6,380	12.8%	14.9%
14	QAM	RAM	-3,094.59	6,370	0.0%	-
15	RAM	RAM	-3,100.27	6,373	0.0%	0.0%
MA 1 (all 15)			-3,071.46			
MA 2 (1,6,10,13,15)			-3,071.65			

10 Most existing applications compare a model combining multiple different decision rules to

1 a set of single class models using the individual rules. This comparison is of course likely to
2 be biased in the presence of taste heterogeneity. Crucially, the improvements to be made from
3 combining different structures depend on their individual performance. For example, we see that
4 for DFT, which is the best performing individual model in Table 7, combining the model with a
5 different structure does not reach as high a log-likelihood as a structure with two separate DFT
6 classes, although a better BIC may be obtained. On the other hand, for those models that perform
7 less well individually, combining them with a different structure gives a better log-likelihood than
8 a model with two classes using the same structure. This already suggests that the results from
9 the latent class structure point more towards taste heterogeneity than decision rule heterogeneity.
10 In fact, when looking at pairs of decision rules, we see only two cases in favour of decision rule
11 heterogeneity, i.e. where a model combining two decision rules is outperforming the two LC
12 models that use the same model type in both classes. The MNL-RRM model outperforms RRM-
13 RRM by 3.67 log-likelihood units and outperforms MNL-MNL by 10.47 units. Additionally,
14 QAM-RAM has a better log-likelihood than either QAM-QAM or RAM-RAM.

15 Further evidence is given in the model averaging results in Table 8. Overall, we observe
16 that averaging across all 15 latent class models results in a slightly better log-likelihood than the
17 best performing individual model (DFT-DFT). Single model type classes also obtain 81.6% of the
18 weight in model averaging. We also observe that averaging across the 5 models that only capture
19 taste heterogeneity alone obtains a very similar log-likelihood in model averaging. These findings
20 again highlight the importance of within-model taste heterogeneity, at least for this data.

21 We explore the most common example of decision rule heterogeneity (MNL-RRM) in more
22 detail by also considering the outputs for the parameter estimates, in comparison to a model average
23 performed on MNL and RRM. The results for this are shown in Table 9. For each model, we have
24 coefficients for travel time (TT), the log of fare (LFare⁵), rate of crowding (Crowd), length of
25 delays (Delay), rate of delays (Rate), a reliability level (Rel), created by calculating the expected
26 length of delays, and the provision of a charged delay information service (Inf) or a free service
27 (InfF). Finally, we include two alternative specific constants for the first two alternatives.

28 Table 9 gives model fit as well as estimates for the above parameters for both a latent class model
29 and a model averaging approach. The model averaging approach separately runs MNL and RRM
30 models before then estimating a class allocation parameter individually. Crucially, the model
31 averaging approach does not result in a significant improvement over a MNL model on its own,
32 with an improvement of just 0.07 log-likelihood units. As a contrast, the latent class approach
33 results in a vast improvement in model fit (258 units). At face value, this would again suggest
34 decision rule heterogeneity, although the fit is not much better than for the MNL-MNL or RRM-
35 RRM models. Most significantly, it appears that the fare parameter estimates (highlighted in red)
36 are very different between the two classes. In contrast with the model averaging results, and
37 given the poor class-specific model fit for the RRM class (compared to the RRM-RRM model),
38 we believe that this finding shows that a substantial share of the improvements obtained by the
39 MNL-RRM model are due to heterogeneity in the cost sensitivity rather than heterogeneity in
40 the decision rules. This means that the classes individually have a very poor fit (as they cannot
41 explain all individuals) but when combined into a latent class approach, the result is a model with

⁵Note that we use a logarithmic transform of the fare rather than the fare itself as a cost damping effect is observed.

TABLE 9 : A detailed example of model averaging compared to a simultaneous latent class approach using MNL and RRM

	Latent Class - 1 model 21 pars, estimated simultaneously		Model averaging - 3 models 2*10 pars, then 1 for MA	
	Class 1:MNL	Class 2:RRM	Class 1: MNL	Class 2: RRM
Class LL: Log-likelihood	-3,645.30	-4,431.55	-3,360.43	-3,363.91
	-3,102.66		-3,360.36	
δ_{alt1}	0.64 (6.42)	0.04 (0.27)	0.39 (5.85)	0.27 (4.17)
δ_{alt2}	0.25 (2.81)	0.20 (1.13)	0.16 (3.3)	0.17 (3.38)
β_{TT}	-0.05 (-6.74)	-0.05 (-6.79)	-0.05 (-9.5)	-0.03 (-9.58)
β_{Lfare}	-3.21 (-6.1)	-11.32 (-7.58)	-6.00 (-18.87)	-4.11 (-17.66)
β_{Crowd}	-0.31 (-7.41)	-0.15 (-2.89)	-0.22 (-8.58)	-0.15 (-8.59)
β_{Delay}	-0.06 (-1.27)	-0.05 (-1.29)	-0.03 (-3.24)	-0.02 (-3.06)
β_{Rate}	-0.34 (-4.82)	-0.09 (-1.76)	-0.19 (-5.96)	-0.12 (-5.82)
β_{Rel}	-0.05 (-3.22)	0.00 (0.06)	-0.06 (-2.64)	-0.04 (-2.71)
β_{Inf}	-0.10 (-0.82)	-0.16 (-1.09)	-0.09 (-1.13)	-0.05 (-0.95)
β_{InfF}	0.54 (5.84)	0.05 (0.47)	0.33 (4.95)	0.22 (4.85)
π_m	59.30% (10.89)	40.70%	87.70% (2.7)	12.30%

- 1 far superior model fit. Together with the poor improvement from model averaging, these results
- 2 suggest that most of the model improvement is due to taste rather than decision rule heterogeneity.

3 5. CONCLUSIONS

4 In this paper, we revisit the use of latent class models to capture heterogeneity across decision-
5 makers in behavioural processes such as attribute non-attendance and decision rule heterogeneity.
6 These approaches have been very popular in recent years and have often been shown to produce
7 significant gains in fit over simpler models. We first argue that many such findings may be due
8 to an unfair comparison with models not allowing for any heterogeneity and that the findings
9 may in fact be driven by heterogeneity in the sensitivities to individual attributes rather than the
10 presence of other phenomena. We have contrasted the findings obtained from such latent class
11 models with those obtained using model averaging which combines the evidence from a number
12 of separately estimated models. This latter approach of course leads to inferior model fit compared
13 to a simultaneous latent class model as model averaging is based on combining different sample
14 level models, i.e. using parameters that are appropriate at the sample level, but our findings provide
15 some evidence that suggests that these bigger improvements may indeed be in part due to effects
16 other than those that analysts seek to uncover. This is especially the case when showing that
17 equivalent (or near equivalent) gains in model fit can be obtained from LC models that use the
18 same structure in each class, thus only allowing for taste heterogeneity.

19 In practice, an analyst should of course attempt to simultaneously allow for all different types
20 of heterogeneity whilst remaining aware of potential confounding. This would however require
21 the use of latent class structures with many different classes and quickly become computationally

1 and empirically infeasible. While we do not suggest that researchers abandon the use of latent
2 class structures to investigate heterogeneity in behavioural processes, we urge for some caution
3 in interpretation and suggest that model averaging can provide a useful tool for checking the
4 likely validity of their insights. In this context, model averaging can provide useful tests. Small
5 improvements in model fit from model averaging point towards other sources of heterogeneity,
6 as our previous work (Hancock et al., 2020b) demonstrates that large improvements from model
7 averaging can be found when averaging over very similar models.

8 As a closing comment, the findings in the application looking at decision rule heterogeneity
9 are particularly insightful. They suggest that there is more scope for heterogeneity in parameters
10 across individuals conditional on a specific model structure rather than heterogeneity across individuals
11 in the model structure itself. In many ways, this is not surprising given that datasets, especially
12 from stated choice surveys, are relatively homogeneous in the structure of the choice sets and
13 explanatory variables. The models that work best are more likely to be dataset-specific rather
14 than person-specific. More work is of course required, including testing using simulated datasets.
15 This is especially important with a view to looking into the ability of model averaging to uncover
16 heterogeneity of the type analysts increasingly attempt to uncover with latent class structures.

17 ACKNOWLEDGEMENTS

18 The authors would like to acknowledge the financial support by the European Research Council
19 through the consolidator grant 615596-DECISIONS.

REFERENCES

- Abou-Zeid, M. and Ben-Akiva, M. (2010). A model of travel happiness and mode switching. In Hess, S. and Daly, A., editors, *Choice Modelling: The State-of-the-Art and the State-of-Practice*, pages 289–305. Emerald Publishing, UK.
- Ben-Akiva, M. and Boccara, B. (1995). Discrete choice models with latent choice sets. *International Journal of Research in Marketing*, 12(1):9–24.
- Berkowitsch, N. A., Scheibehenne, B., and Rieskamp, J. (2014). Rigorously testing multialternative decision field theory against random utility models. *Journal of Experimental Psychology: General*, 143(3):1331.
- Boeri, M. and Longo, A. (2017). The importance of regret minimization in the choice for renewable energy programmes: Evidence from a discrete choice experiment. *Energy Economics*, 63:253–260.
- Bruza, P. D., Wang, Z., and Busemeyer, J. R. (2015). Quantum cognition: a new theoretical approach to psychology. *Trends in cognitive sciences*, 19(7):383–393.
- Busemeyer, J. R. and Townsend, J. T. (1992). Fundamental derivations from decision field theory. *Mathematical Social Sciences*, 23(3):255–282.
- Campbell, D., Lorimer, V., Aravena, C., and Hutchinson, W. G. (2010). Attribute processing in environmental choice analysis: implications for willingness to pay. 84th Annual Conference, March 29-31, 2010, Edinburgh, Scotland 91718, Agricultural Economics Society.
- Chorus, C. G. (2010). A new model of random regret minimization. *EJTIR*, 10 (2), 2010.

- Chorus, C. G. (2014). Capturing alternative decision rules in travel choice models: a critical discussion. In *Handbook of choice modelling*. Edward Elgar Publishing.
- Chorus, C. G., Arentze, T. A., and Timmermans, H. J. (2008). A random regret-minimization model of travel choice. *Transportation Research Part B: Methodological*, 42(1):1–18.
- Dey, B. K., Anowar, S., Eluru, N., and Hatzopoulou, M. (2018). Accommodating exogenous variable and decision rule heterogeneity in discrete choice models: Application to bicyclist route choice. *PloS one*, 13(11):e0208309.
- Gopinath, D. (1995). *Modeling Heterogeneity in Discrete Choice Processes: Application to Travel Demand*. PhD thesis, MIT, Cambridge, MA.
- Greene, W. H. and Hensher, D. A. (2003). A latent class model for discrete choice analysis: contrasts with mixed logit. *Transportation Research Part B: Methodological*, 37(8):681–698.
- Hancock, T. O., Broekaert, J., Hess, S., and Choudhury, C. F. (2020a). Quantum probability: a new method for modelling travel behaviour. *Submitted*.
- Hancock, T. O., Hess, S., and Choudhury, C. F. (2018). Decision field theory: Improvements to current methodology and comparisons with standard choice modelling techniques. *Transportation Research Part B: Methodological*, 107:18–40.
- Hancock, T. O., Hess, S., Daly, A., and Fox, J. (2020b). Using a sequential latent class approach for model averaging: benefits in forecasting and behavioural insights. *Submitted*.
- Hancock, T. O., Hess, S., Marley, A. A. J., and Choudhury, C. F. (2020c). An accumulation of preference: two alternative dynamic models for understanding transport choices. *Submitted*.
- Hensher, D. (2014). Attribute processing as a behavioural strategy in choice making. In *Handbook of choice modelling*. Edward Elgar Publishing.
- Hensher, D. A. (2010). Attribute processing, heuristics and preference construction in choice analysis. In Hess, S. and Daly, A. J., editors, *State-of Art and State-of Practice in Choice Modelling: Proceedings from the Inaugural International Choice Modelling Conference*, chapter 3, pages 35–70. Emerald, Bingley, UK.
- Hensher, D. A. and Greene, W. H. (2010). Non-attendance and dual processing of common-metric attributes in choice analysis: a latent class specification. *Empirical Economics*, 39(4):413–426.
- Hensher, D. A., Rose, J. M., and Greene, W. H. (2012). Inferring attribute non-attendance from stated choice data: implications for willingness to pay estimates and a warning for stated choice experiment design. *Transportation*, 39(2):235–245.
- Hess, S. (2014). 14 latent class structures: taste heterogeneity and beyond. In *Handbook of choice modelling*, pages 311–329. Edward Elgar Publishing Cheltenham.
- Hess, S., Beck, M., and Crastes dit Sourd, R. (2016). Can a better model specification avoid the need to move away from random utility maximisation? Transportation Research Board (TRB) 96th Annual Meeting.
- Hess, S. and Rose, J. M. (2007). *A latent class approach to recognising respondents' information processing strategies in SP studies*. paper presented at the Oslo Workshop on Valuation Methods in Transport Planning, Oslo.
- Hess, S., Shires, J., and Jopson, A. (2013a). Accommodating underlying pro-environmental attitudes in a rail travel context: Application of a latent variable latent class specification. *Transportation Research Part D*, 25:42–48.
- Hess, S. and Stathopoulos, A. (2012). Linking the decision process to underlying attitudes and perceptions: a latent variable latent class construct. *paper presented at the 13th International Conference on Travel Behaviour Research, Toronto*.

- Hess, S. and Stathopoulos, A. (2013). A mixed random utility - random regret model linking the choice of decision rule to latent character traits. *Journal of Choice Modelling*, 9:27–38.
- Hess, S., Stathopoulos, A., Campbell, D., O’Neill, V., and Caussade, S. (2013b). It’s not that I don’t care, I just don’t care very much: confounding between attribute non-attendance and taste heterogeneity. *Transportation*, 40(3):583–607.
- Hess, S., Stathopoulos, A., and Daly, A. J. (2012). Allowing for heterogeneous decision rules in discrete choice models: an approach and four case studies. *Transportation*, 39(3):565–591.
- Hole, A. R. (2011). A discrete choice model with endogenous attribute attendance. *Economics Letters*, 110(3):203–205.
- Ishaq, R., Bekhor, S., and Shiftan, Y. (2013). A flexible model structure approach for discrete choice models. *Transportation*, 40(3):60–624.
- Leong, W. and Hensher, D. A. (2014). Relative advantage maximisation as a model of context dependence for binary choice data. *Journal of choice modelling*, 11:30–42.
- Montgomery, H. and Svenson, O. (1976). On decision rules and information processing strategies for choices among multiattribute alternatives. *Scandinavian Journal of Psychology*, 17(1):283–291.
- Roe, R. M., Busemeyer, J. R., and Townsend, J. T. (2001). Multialternative decision field theory: A dynamic connectionist model of decision making. *Psychological Review*, 108(2):370.
- Scarpa, R., Gilbride, T., Campbell, D., and Hensher, D. A. (2009). Modelling attribute non-attendance in choice experiments for rural landscape valuation. *European Review of Agricultural Economics*, 36(2):151–174.
- Swait, J. and Ben-Akiva, M. (1985). Incorporating random constraints in discrete models of choice set generation. *Transportation Research Part B*, 21(2):91–102.
- Tversky, A. (1972). Elimination by aspects: a theory of choice. *Psychological Review*, 79:281–299.
- Walker, J. and Ben-Akiva, M. (2002). Generalized random utility model. *Mathematical Social Sciences*, 43(3):303–343.